Pacific Journal of Mathematics

Volume 263 No. 2

June 2013

PACIFIC JOURNAL OF MATHEMATICS

msp.org/pjm

Founded in 1951 by E. F. Beckenbach (1906-1982) and F. Wolf (1904-1989)

EDITORS

V. S. Varadarajan (Managing Editor) Department of Mathematics University of California Los Angeles, CA 90095-1555 pacific@math.ucla.edu

Paul Balmer Department of Mathematics University of California Los Angeles, CA 90095-1555 balmer@math.ucla.edu

Daryl Cooper Department of Mathematics University of California Santa Barbara, CA 93106-3080 cooper@math.ucsb.edu

Jiang-Hua Lu Department of Mathematics The University of Hong Kong Pokfulam Rd., Hong Kong jhlu@maths.hku.hk Don Blasius Department of Mathematics University of California Los Angeles, CA 90095-1555 blasius@math.ucla.edu

Robert Finn Department of Mathematics Stanford University Stanford, CA 94305-2125 finn@math.stanford.edu

Sorin Popa Department of Mathematics University of California Los Angeles, CA 90095-1555 popa@math.ucla.edu

Paul Yang Department of Mathematics Princeton University Princeton NJ 08544-1000 yang@math.princeton.edu

PRODUCTION

Silvio Levy, Scientific Editor, production@msp.org

SUPPORTING INSTITUTIONS

ACADEMIA SINICA, TAIPEI CALIFORNIA INST. OF TECHNOLOGY INST. DE MATEMÁTICA PURA E APLICADA KEIO UNIVERSITY MATH. SCIENCES RESEARCH INSTITUTE NEW MEXICO STATE UNIV. OREGON STATE UNIV. STANFORD UNIVERSITY UNIV. OF BRITISH COLUMBIA UNIV. OF CALIFORNIA, BERKELEY UNIV. OF CALIFORNIA, DAVIS UNIV. OF CALIFORNIA, LOS ANGELES UNIV. OF CALIFORNIA, RIVERSIDE UNIV. OF CALIFORNIA, SAN DIEGO UNIV. OF CALIF., SANTA BARBARA Vyjayanthi Chari Department of Mathematics University of California Riverside, CA 92521-0135 chari@math.ucr.edu

Kefeng Liu Department of Mathematics University of California Los Angeles, CA 90095-1555 liu@math.ucla.edu

Jie Qing Department of Mathematics University of California Santa Cruz, CA 95064 qing@cats.ucsc.edu

UNIV. OF CALIF., SANTA CRUZ UNIV. OF MONTANA UNIV. OF OREGON UNIV. OF SOUTHERN CALIFORNIA UNIV. OF UTAH UNIV. OF WASHINGTON WASHINGTON STATE UNIVERSITY

These supporting institutions contribute to the cost of publication of this Journal, but they are not owners or publishers and have no responsibility for its contents or policies.

See inside back cover or msp.org/pjm for submission instructions.

The subscription price for 2013 is US \$400/year for the electronic version, and \$485/year for print and electronic. Subscriptions, requests for back issues and changes of subscribers address should be sent to Pacific Journal of Mathematics, P.O. Box 4163, Berkeley, CA 94704-0163, U.S.A. The Pacific Journal of Mathematics is indexed by Mathematical Reviews, Zentralblatt MATH, PASCAL CNRS Index, Referativnyi Zhurnal, Current Mathematical Publications and the Science Citation Index.

The Pacific Journal of Mathematics (ISSN 0030-8730) at the University of California, c/o Department of Mathematics, 798 Evans Hall #3840, Berkeley, CA 94720-3840, is published monthly except July and August. Periodical rate postage paid at Berkeley, CA 94704, and additional mailing offices. POSTMASTER: send address changes to Pacific Journal of Mathematics, P.O. Box 4163, Berkeley, CA 94704-0163.

PJM peer review and production are managed by EditFLOW® from Mathematical Sciences Publishers.

PUBLISHED BY



http://msp.org/ © 2013 Mathematical Sciences Publishers

REALIZATIONS OF BC_r -GRADED INTERSECTION MATRIX ALGEBRAS WITH GRADING SUBALGEBRAS OF TYPE B_r , $r \ge 3$

SANDEEP BHARGAVA AND YUN GAO

We study intersection matrix algebras $im(A^{[d]})$ that arise from affinizing a Cartan matrix A of type B_r with d arbitrary long roots in the root system Δ_{B_r} , where $r \geq 3$. We show that $im(A^{[d]})$ is isomorphic to the universal covering algebra of $so_{2r+1}(\alpha, \eta, C, \chi)$, where α is an associative algebra with involution η , and C is an α -module with hermitian form χ . We provide a description of all four of the components α , η , C, and χ .

1. Introduction

Peter Slodowy [1984; 1986] discovered that matrices like

$$M = \begin{bmatrix} 2 & -1 & 0 & 1 \\ -1 & 2 & -1 & 1 \\ 0 & -2 & 2 & -2 \\ 1 & 1 & -1 & 2 \end{bmatrix}$$

encode the intersection form on the second homology group of Milnor fibers for germs of holomorphic maps with an isolated singularity at the origin. These matrices were like the generalized Cartan matrices of Kac–Moody theory in that they had integer entries, 2's along the diagonal, and M_{ij} was negative if and only if M_{ji} was negative. What was new, however, was the presence of positive entries off the diagonal. Slodowy called such matrices generalized intersection matrices:

Definition 1 [Slodowy 1986]. An $n \times n$ integer-valued matrix M is called a *generalized intersection matrix* (gim) if the following conditions are satisfied. whenever $1 \le i, j \le n$ with $i \ne j$:

$$M_{ii} = 2;$$

 $M_{ij} < 0 \iff M_{ji} < 0;$
 $M_{ij} > 0 \iff M_{ji} > 0.$

Funding from the National Sciences and Engineering Research Council of Canada is gratefully acknowledged.

MSC2010: primary 17B65, 17B70, 17B05; secondary 17B67, 16W10.

Keywords: Lie algebras, intersection matrix algebras.

Slodowy used these matrices to define a class of Lie algebras that encompassed all the Kac–Moody Lie algebras:

Definition 2 [Slodowy 1986; Berman and Moody 1992]. Given an $n \times n$ generalized intersection matrix $M = (M_{ij})$, define a Lie algebra over \mathbb{C} , called a *generalized intersection matrix* (gim) *algebra* and denoted by gim(M), with

generators:
$$e_1, ..., e_n, f_1, ..., f_n, h_1, ..., h_n$$
,
relations:
(R1) For $1 \le i, j \le n$,
 $[h_i, e_j] = M_{ij}e_j$,
 $[h_i, f_j] = -M_{ij}f_j$,
 $[e_i, f_i] = h_i$.
(R2) For $M_{ij} \le 0$,
 $[e_i, f_j] = 0 = [f_i, e_j]$,
 $(ad e_i)^{-M_{ij}+1}e_j = 0 = (ad f_i)^{-M_{ij}+1}f_j$.
(R3) For $M_{ij} > 0, i \ne j$,
 $[e_i, e_j] = 0 = [f_i, f_j]$,
 $(ad e_i)^{M_{ij}+1}f_j = 0 = (ad f_i)^{M_{ij}+1}e_j$.

If the M that we begin with is a generalized Cartan matrix, then the 3n generators and the first two groups of axioms, (R1) and (R2), provide a presentation of the Kac–Moody Lie algebras [Gabber and Kac 1981; Kac 1990; Carter 2005].

Slodowy [1986] and, later, Berman [1989] showed that the gim algebras are also isomorphic to fixed point subalgebras of involutions on larger Kac–Moody algebras. So, in their words, the gim algebras lie both "beyond and inside" Kac–Moody algebras.

Further progress came in the 1990s as a byproduct of work on the classification of root-graded Lie algebras [Berman and Moody 1992; Benkart and Zelmanov 1996; Neher 1996], which revealed that some families of intersection matrix (im) algebras, which are quotient algebras of gim algebras, were universal covering algebras of well-understood Lie algebras. For instance the im algebras that arise from multiply affinizing a Cartan matrix of type A_r , with $r \ge 3$, are the universal covering algebras of sl(a), where a is the associative algebra of noncommuting Laurent polynomials in several variables (the number of indeterminates depends on how many times the original Cartan matrix is affinized). A handful of other researchers also began engaging these new algebras. For example, Eswara Rao, Moody, and Yokonuma [Rao et al. 1992] used vertex operator representations to show that im algebras were nontrivial. Gao [1996] examined compact forms of im algebras arising from conjugations over the complex field. Peng [2002] found relations between im algebras and the representations of tilted algebras via

Ringel–Hall algebras. Berman, Jurisich, and Tan [Berman et al. 2001] showed that the presentation of gim algebras could be put into a broader framework that incorporated Borcherds algebras.

The chief objective of this paper is to continue advancing our understanding of gim and im algebras. We construct a generalized intersection matrix $A^{[d]}$ by adjoining d long roots to a base of a root system of type B_r , where $r \ge 3$. This is exactly the analogue of the affinization process in which a single root is adjoined to a Cartan matrix of a finite-dimensional Lie algebra to arrive at a generalized Cartan matrix and, eventually, an affine Kac–Moody algebra. The matrix $A^{[d]}$ is used to define a gim algebra gim $(A^{[d]})$. Since gim $(A^{[d]})$ may possess roots with mixed signs, we quotient out by an ideal r that is tailor-made to capture all such roots. The quotient algebra is called the intersection matrix algebra and is denoted by im $(A^{[d]})$.

We show that $\mathfrak{im}(A^{[d]})$ is a BC_r -graded Lie algebra, which, in turn, allows us to invoke Allison, Benkart, and Gao's recognition theorem [Allison et al. 2002] and relate $im(A^{[d]})$ to an algebraic structure that is better understood. Combining their theorem with the knowledge that $im(A^{[d]})$ is centrally closed, we conclude that, up to isomorphism, $\mathfrak{im}(A^{[d]})$ is the universal covering algebra of $\mathfrak{so}_{2r+1}(\mathfrak{a}, \eta, C, \chi)$. The algebra $so_{2r+1}(\mathfrak{a}, \eta, C, \chi)$ is like the usual matrix model $so_{2r+1}(\mathbb{C})$ of a finitedimensional Lie algebra of type B_r , except that we now replace the field \mathbb{C} with an associative algebra a, which possesses an involution (that is, period two antiautomorphism) η , and we involve a right \mathfrak{a} -module C that has a hermitian form $\chi: C \times C \to \mathfrak{a}$. The defining relations of the generalized intersection matrix algebra and, hence, the intersection matrix algebra, in concert with the existence of a central, graded, surjective Lie algebra homomorphism ψ from $\mathfrak{im}(A^{[d]})$ to $\mathfrak{so}_{2r+1}(\mathfrak{a}, \eta, C, \chi)$ allow us to understand each of \mathfrak{a} , η , C, and χ . For example, we get (i) two generators of a, namely x and x^{-1} , for every long root of the form $\pm(\epsilon_i + \epsilon_{i+1})$, and (ii) four generators of a, namely y, y^{-1} , z, and z^{-1} , for every other type of long root that we adjoin. We are also able to study the relations among the generators, determine the action of the involution η , and discover that C = 0 and $\chi = 0$. Through constructing a surjective Lie algebra homomorphism $\varphi : \mathfrak{gim}(A^{[d]}) \to \mathfrak{so}_{2r+1}(\mathfrak{a}, \eta, C, \chi)$ we verify that we indeed have a complete description of the "coordinate algebra" a.

Our work continues the line of research initiated by Berman, Moody, Benkart and Zelmanov. Berman and Moody [1992] were the first to find realizations of intersection matrix algebras over Lie algebras graded by root systems of types A_r $(r \ge 2)$, D_r , E_6 , E_7 , and E_8 . Benkart and Zelmanov [1996] found realizations of intersection matrix algebras over Lie algebras graded by root systems of types A_1 , B_r , C_r , F_4 , and G_2 . In this paper, we find realizations of intersection matrix algebras over Lie algebras graded by root systems of types BC_r with grading subalgebras of type B_r $(r \ge 3)$.

2. Multiply affinizing Cartan matrices

In this paper, we focus on generalized intersection matrix algebras that arise from multiply affinizing a Cartan matrix of type B_r , where $r \ge 3$, with long roots in the root system Δ_{B_r} .

Consider a root system of type B_r . Up to isomorphism, Δ_{B_r} may be described as

$$\Delta_{B_r} = \{\pm \epsilon_i \pm \epsilon_j : 1 \le i \ne j \le r\} \cup \{\pm \epsilon_i : i = 1, \dots, r\}.$$

Once we fix an ordering of the simple roots $\alpha_1, \ldots, \alpha_r$ in a base Π , the Cartan matrix *A* is described by

$$A_{ij} = \frac{2(\alpha_i, \alpha_j)_{\text{Killing}}}{(\alpha_i, \alpha_i)_{\text{Killing}}} \quad \text{for } 1 \le i, j \le r.$$

Choose any *d* long roots in Δ_{B_r} , say $\alpha_{r+1}, \ldots, \alpha_{r+d}$, and consider the r+d by r+d matrix $A^{[d]}$ given by

$$A_{ij}^{[d]} = \frac{2(\alpha_i, \alpha_j)_{\text{Killing}}}{(\alpha_i, \alpha_i)_{\text{Killing}}} \quad \text{for } 1 \le i, j \le r+d,$$

with respect to the ordering $(\alpha_1, \ldots, \alpha_r, \alpha_{r+1}, \ldots, \alpha_{r+d})$ of the *r* roots in the base Π plus the *d* "adjoined" roots. The axioms of a root system tell us that all the entries of $A^{[d]}$ are integers. Moreover, since the Killing form is symmetric, we have $A_{ji}^{[d]} = 0$ if $A_{ij}^{[d]} = 0$, or if $A_{ij}^{[d]}$ and $A_{ji}^{[d]}$ are nonzero, then they share the same sign. In other words, $A^{[d]}$ is a generalized intersection matrix.

Since the "d-affinized" Cartan matrix $A^{[d]}$ is a generalized intersection matrix, $gim(A^{[d]})$ is a generalized intersection matrix algebra.

Note that if we affinize the Cartan matrix A of type B_r with the negative of the highest long root of Δ_{B_r} then the resulting generalized intersection matrix algebra gim $(A^{[1]})$ is the affine Kac–Moody Lie algebra of type $B_r^{(1)}$.

3. Intersection matrix algebras

Fix a Cartan matrix A of type B_r $(r \ge 3)$ with, say, $\alpha_1, \alpha_2, \ldots, \alpha_r$ being the simple roots in a base of Δ_{B_r} that were used to form A. Let

- Ω = set of all long roots of the form $\pm(\epsilon_i + \epsilon_{i+1})$ that we adjoin,
- Θ = set of all remaining long roots that are adjoined,
- N_{μ} = the number of copies of the long root μ we have adjoined, and

•
$$d = \sum_{\mu \in \Omega \cup \Theta} N_{\mu}$$
.

Let $A^{[d]}$ be the resulting generalized intersection matrix and $\mathfrak{gim}(A^{[d]})$ the corresponding generalized intersection matrix algebra.

We begin a move towards a quotient algebra of $\mathfrak{gim}(A^{[d]})$ using a slight generalization of the work done by Benkart and Zelmanov [1996]. Let Γ be the integer lattice generated by the Δ , where

$$\Delta = \{\pm \epsilon_i \pm \epsilon_j : 1 \le i \ne j \le r\} \cup \{\pm \epsilon_i, \pm 2\epsilon_i : i = 1, \dots, r\}$$

is a root system of type BC_r .

We define a Γ -grading on $gim(A^{[d]})$ as follows:

$$\deg e_i = \alpha_i = -\deg f_i, \quad \deg h_i = 0$$

for i = 1, ..., r, and

$$\deg e_{\mu,i} = \mu = -\deg f_{\mu,i}, \quad \deg h_{\mu,i} = 0$$

for $\mu \in \Omega \cup \Theta$ and $i = 1, \ldots, N_{\mu}$.

Next, we define the *radical* \mathfrak{r} of $\mathfrak{gim}(A^{[d]})$ to be the ideal generated by the root spaces $\mathfrak{gim}(A^{[d]})_{\gamma}$ where $\gamma \notin \Delta \cup \{0\}$. Since the ideal \mathfrak{r} is homogeneous, the resulting quotient algebra

$$\mathfrak{im}(A^{[d]}) = \mathfrak{gim}(A^{[d]})/\mathfrak{r}$$

is also Γ-graded. Moreover,

$$\mathfrak{im}(A^{\lfloor d \rfloor})_{\gamma} = 0 \quad \text{if } \gamma \notin \Delta \cup \{0\}.$$

We call $im(A^{[d]})$ the *intersection matrix* (im) *algebra* corresponding to the generalized intersection matrix algebra $gim(A^{[d]})$.

3.1. $im(A^{[d]})$ is *BC_r-graded*. Allison, Benkart, and Gao gave the following definition of a Lie algebra graded by a root system of type *BC*.

Definition 3 [Allison et al. 2002]. Let *r* be a positive integer greater than or equal to 3. A Lie algebra *L* over \mathbb{C} is graded by the root system BC_r or is BC_r -graded with a grading subalgebra of type B_r if

- (i) *L* contains, as a subalgebra, a finite-dimensional simple Lie algebra \mathfrak{g} whose root system relative to a Cartan subalgebra $\mathfrak{h} = \mathfrak{g}_0$ is Δ_{B_r} ,
- (ii) $L = \bigoplus_{\mu \in \Delta \cup \{0\}} L_{\mu}$, where $L_{\mu} = \{x \in L \mid [h, x] = \mu(h)x \text{ for all } h \in \mathfrak{h}\}$ for $\mu \in \Delta \cup \{0\}$, and Δ is the root system of type BC_r , and
- (iii) $L_0 = \sum_{\mu \in \Delta} [L_{\mu}, L_{-\mu}].$

Proposition 4. The algebra $im(A^{[d]})$ is BC_r -graded with a grading subalgebra of type B_r .

Proof. The subalgebra in $im(A^{[d]})$ generated by $e_1 + \mathfrak{r}, \ldots, h_r + \mathfrak{r}$, due to the relations on these elements induced by the relations on their preimages in $gim(A^{[d]})$, is isomorphic to a finite-dimensional simple Lie algebra \mathfrak{g} of type B_r . We have already shown in Section 3.1 that $im(A^{[d]})$ is Γ -graded with $im(A^{[d]})_{\gamma} = 0$ if $\gamma \notin \Delta \cup \{0\}$. That is,

$$\mathfrak{im}(A^{[d]}) = \bigoplus_{\mu \in \Delta \cup \{0\}} \mathfrak{im}(A^{[d]})_{\mu}.$$

Finally, our initial degree assignments for the generators of $\mathfrak{gim}(A^{[d]})$, the \mathfrak{gim} algebra relations like $h_i = [e_i, f_i]$ and $h_\mu = [e_\mu, f_\mu]$, and the fact that movement into the 0 root space can only occur by bracketing an element from an $\mathfrak{im}(A^{[d]})_\mu$ space with one from the $\mathfrak{im}(A^{[d]})_{-\mu}$ space all combine to lead us to the conclusion that

$$\operatorname{im}(A^{[d]})_0 = \sum_{\mu \in \Delta} [\operatorname{im}(A^{[d]})_{\mu}, \operatorname{im}(A^{[d]})_{-\mu}].$$

3.2. $\operatorname{im}(A^{[d]})$ *is centrally closed.* Recall that a Lie algebra *L* is said to be perfect if it equals its derived algebra, that is, L = [L, L]. Furthermore, if *L* is perfect and is its own universal covering then we say that *L* is centrally closed [Moody and Pianzola 1995].

Proposition 5. The algebra $gim(A^{[d]})$ is a perfect Lie algebra.

Proof. Being a Lie algebra, $\mathfrak{gim}(A^{[d]})$ is closed under the operation of taking brackets; hence $\left[\mathfrak{gim}(A^{[d]}), \mathfrak{gim}(A^{[d]})\right] \subset \mathfrak{gim}(A^{[d]})$. To show the reverse inclusion, it suffices to show that all of the generators of $\mathfrak{gim}(A^{[d]})$ lie in $\left[\mathfrak{gim}(A^{[d]}), \mathfrak{gim}(A^{[d]})\right]$. But this is indeed the case because the generators e_i , f_i , h_i (for $1 \le i \le r$) and the $e_{\mu,i}$, $f_{\mu,i}$, $h_{\mu,i}$, which arise from adjoining the *i*-th copy of a long root μ , satisfy the relations (R1) of Definition 2.

Our next theorem is Proposition 1.6 in [Benkart and Zelmanov 1996] adapted to our context.

Theorem 6. The algebra $im(A^{[d]})$ is centrally closed.

Proof. Let (\tilde{U}, ϕ) be the universal covering algebra of $\operatorname{im}(A^{[d]})$. Let \mathfrak{g} be the simple finite dimensional subalgebra of type *B* contained in $\operatorname{im}(A^{[d]})$ with Cartan subalgebra \mathfrak{h} whose root space decomposition induces a *BC*-gradation on $\operatorname{im}(A^{[d]})$. The preimage $\phi^{-1}(\mathfrak{h})$ of \mathfrak{h} contains ker ϕ . Since ϕ is a central map, ker ϕ lies in the center of \tilde{U} . So

$$\mathfrak{h}' = \phi^{-1}(\mathfrak{h}) / \ker \phi$$

acts on \widetilde{U} via the adjoint action. If $h' \in \mathfrak{h}'$, $\phi(h') = h \in \mathfrak{h}$, and $\mu(t) \in \mathbb{C}[t]$ is the minimal polynomial of $\operatorname{ad}_{U_L}(h)$, then

$$\mu\left(\operatorname{ad}_{\widetilde{U}}(h')\right)\left(\widetilde{U}\right)\subset \ker\phi.$$

So $\operatorname{ad}_{\widetilde{U}}(h')$ satisfies the polynomial $t\mu(t)$. Therefore \widetilde{U} is a sum of root spaces with respect to $\operatorname{ad}_{\widetilde{U}}\mathfrak{h}'$, and $\widetilde{U}_{\gamma} \neq (0)$ if and only if $\gamma \in \Delta \cup \{0\}$. So ϕ induces an isomorphism between the nonzero root spaces of \widetilde{U} and those of $\operatorname{im}(A^{[d]})$. Moreover,

$$\widetilde{U}_0 = \sum_{\gamma \in \Delta} \left[\widetilde{U}_{-\gamma}, \widetilde{U}_{\gamma} \right] + \ker \phi \quad \text{implies that} \quad \left[\widetilde{U}_0, \widetilde{U}_0 \right] \subset \sum_{\gamma \in \Delta} \left[\widetilde{U}_{-\gamma}, \widetilde{U}_{\gamma} \right].$$

Since $\widetilde{U} = [\widetilde{U}, \widetilde{U}]$, it follows that

$$\widetilde{U}_{0} = \left[\widetilde{U}_{0}, \widetilde{U}_{0}\right] + \sum_{\gamma \in \Delta} \left[\widetilde{U}_{-\gamma}, \widetilde{U}_{\gamma}\right] = \sum_{\gamma \in \Delta} \left[\widetilde{U}_{-\gamma}, \widetilde{U}_{\gamma}\right]$$

Consequently, ϕ is an isomorphism.

4. Recognition theorem

The following construction, given in Example 1.23 of [Allison et al. 2002], is a more general version of the classical construction of so_{2r+1} (\mathbb{C}), the simple Lie algebra of type B_r .

Let *r* be a positive integer, a be a unital associative algebra over \mathbb{C} with an involution (that is, period two antiautomorphism) η , *C* be a right a-module with a hermitian form $\chi : C \times C \rightarrow \mathfrak{a}$, that is a biadditive map $\chi : C \times C \rightarrow \mathfrak{a}$ satisfying

$$\chi(c,c'\cdot a) = \chi(c,c') \cdot a, \quad \chi(c \cdot a,c') = \eta(a) \cdot \chi(c,c'), \quad \chi(c,c') = \eta\bigl(\chi(c',c)\bigr),$$

for $c, c' \in C$, $a \in \mathfrak{a}$, and G be the $(2r + 1) \times (2r + 1)$ matrix

$$G = \begin{bmatrix} 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & \cdots & 1 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 1 & \cdots & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 \end{bmatrix}.$$

Also, given any $c \in C$, define $\chi_c \in C^*$ by $\chi_c(c') := \chi(c, c')$, for any $c' \in C$, and given any

$$\underline{c} = \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} \in C^{2r+1}, \quad \text{define } \chi_{\underline{c}} := \begin{bmatrix} \chi_{c_1} \\ \vdots \\ \chi_{c_n} \end{bmatrix} \in (C^*)^{2r+1}.$$

Now set

$$\mathfrak{A}(\chi) := \left\{ N \in \operatorname{End}_{\mathfrak{a}}(C) : \chi(Nc, c') + \chi(c, Nc') = 0 \text{ for all } c, c' \in C \right\},\\ \mathfrak{A} := \left\{ \begin{bmatrix} M & \chi_{\underline{c}} \\ \underline{c}^{t}G & N \end{bmatrix} : M \in \operatorname{M}_{2r+1}(\mathfrak{a}), (M^{\eta})^{t}G + GM = 0, \underline{c} \in C^{2r+1}, N \in \mathfrak{A}(\chi) \right\}.$$

It can be checked that \mathfrak{A} is a Lie algebra that contains a simple Lie algebra

$$\mathfrak{g} = \left\{ \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} : M \in \mathcal{M}_{2r+1}(\mathbb{C}), \, M^t G + G M = 0 \right\},\$$

of type B_r . If \mathfrak{h} denotes the Cartan subalgebra of diagonal matrices in \mathfrak{g} , then the adjoint action of \mathfrak{h} on \mathfrak{A} induces a root space decomposition

$$\mathfrak{A} = \bigoplus_{\mu \in \Delta \cup \{0\}} \mathfrak{A}_{\mu}, \quad \text{where } \mathfrak{A}_{\mu} = \big\{ T \in \mathfrak{A} : [h, T] = \mu(h) \ T \text{ for all } h \in \mathfrak{h} \big\}.$$

The following abbreviated notation helps describe these root spaces:

for
$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_{2r+1} \end{bmatrix} \in \mathbb{C}^{2r+1}$$
 and $c \in C$, let $vc := \begin{bmatrix} v_1c \\ \vdots \\ v_{2r+1}c \end{bmatrix} \in C^{2r+1}$.

Then $C^{2r+1} = \bigoplus_{i=1}^{2r+1} e_i C$, where e_1, \ldots, e_{2r+1} is the standard basis for \mathbb{C}^{2r+1} . Letting *B* denote the set of skew-symmetric elements of a relative to the involution η , we have

$$\begin{split} \mathfrak{A}_{\epsilon_{i}-\epsilon_{j}} &= \left\{ E_{i,j}(a) + E_{2r+2-j,2r+2-i}(-\eta(a)) : a \in \mathfrak{a} \right\}, &1 \leq i \neq j \leq r, \\ \mathfrak{A}_{\epsilon_{i}+\epsilon_{j}} &= \left\{ E_{i,2r+2-j}(a) + E_{j,2r+2-i}(-\eta(a)) : a \in \mathfrak{a} \right\}, &1 \leq i, j \leq r, \\ \mathfrak{A}_{-\epsilon_{i}-\epsilon_{j}} &= \left\{ E_{2r+2-i,j}(a) + E_{2r+2-j,i}(-\eta(a)) : a \in \mathfrak{a} \right\}, &1 \leq i, j \leq r, \\ \mathfrak{A}_{\epsilon_{i}} &= \left\{ \begin{bmatrix} 0\chi_{e_{i}c} \\ (e_{2r+2-i}c)^{t} & 0 \end{bmatrix} : c \in C \right\} \\ &+ \left\{ E_{i,r+1}(a) + E_{r+1,2r+2-i}(-\eta(a)) : a \in \mathfrak{a} \right\}, &1 \leq i \leq r, \\ \mathfrak{A}_{-\epsilon_{i}} &= \left\{ \begin{bmatrix} 0 & \chi_{e_{2r+2-i}c} \\ (e_{i}c)^{t} & 0 \end{bmatrix} : c \in C \right\} \\ &+ \left\{ E_{r+1,i}(a) + E_{2r+2-i,r+1}(-\eta(a)) : a \in \mathfrak{a} \right\}, &1 \leq i \leq r, \\ \mathfrak{A}_{0} &= \left\{ \sum_{i=1}^{r} E_{ii}(a) + E_{2r+2-i,2r+2-i}(-\eta(a)) : a \in \mathfrak{a} \right\} + \left\{ \begin{bmatrix} 0 & 0 \\ 0 & N \end{bmatrix} : N \in \mathfrak{A}(\chi) \right\} \\ &+ \left\{ E_{r+1,r+1}(b) : b \in B \right\} + \left\{ \begin{bmatrix} 0 & \chi_{e_{r+1}c} \\ (e_{r+1}c)^{t} & 0 \end{bmatrix} : c \in C \right\}. \end{split}$$

The subalgebra

$$\operatorname{so}_{2r+1}(\mathfrak{a},\eta,C,\chi) := \sum_{\mu \in \Delta} \mathfrak{A}_{\mu} + \sum_{\mu \in \Delta} [\mathfrak{A}_{\mu},\mathfrak{A}_{-\mu}]$$

of \mathfrak{A} has the root spaces

$$so_{2r+1}(\mathfrak{a}, \eta, C, \chi)_0 = so_{2r+1}(\mathfrak{a}, \eta, C, \chi) \cap \mathfrak{A}_0,$$

$$so_{2r+1}(\mathfrak{a}, \eta, C, \chi)_\mu = \mathfrak{A}_\mu \quad \text{for } \mu \in \Delta.$$

In particular,

$$\operatorname{so}_{2r+1}(\mathfrak{a},\eta,C,\chi)_0 = \sum_{\mu \in \Delta} \left[\operatorname{so}_{2r+1}(\mathfrak{a},\eta,C,\chi)_{\mu}, \operatorname{so}_{2r+1}(\mathfrak{a},\eta,C,\chi)_{-\mu} \right].$$

Remark. In [Allison et al. 2002] the notation *L* is used to refer to the Lie algebra that we are calling $so_{2r+1}(\mathfrak{a}, \eta, C, \chi)$.

To shorten the description of elements in $so_{2r+1}(\mathfrak{a}, \eta, C, \chi)$, we use the following notation: Given any $1 \le k \le r$ and $a \in \mathfrak{a}$, let

$$E_{k,r+1}^{[\square]}(a) := E_{k,r+1}(a) + E_{r+1,2r+2-k}(-\eta(a)),$$

$$E_{r+1,k}^{[\square]}(a) := E_{r+1,k}(a) + E_{2r+2-k,r+1}(-\eta(a)),$$

and for any $1 \le p, q \le r$ and $a \in \mathfrak{a}$, let

$$E_{p,q}^{\square}(a) := E_{p,q}(a) + E_{2r+2-q,2r+2-p}(-\eta(a)),$$

$$E_{p,2r+2-q}^{\square}(a) := E_{p,2r+2-q}(a) + E_{q,2r+2-p}(-\eta(a)),$$

$$E_{2r+2-p,q}^{\square}(a) := E_{2r+2-p,q}(a) + E_{2r+2-q,p}(-\eta(a)).$$

We often also denote the involution η on \mathfrak{a} by $\overline{\cdot}$. So, for example, we would write

$$E_{2r+2-p,q}^{[b]}(a)$$
 (above) as $E_{2r+2-p,q}(a) + E_{2r+2-q,p}(-\bar{a})$.

Allison, Benkart, and Gao's classification results on BC_r -graded Lie algebras [Allison et al. 2002] say the following in our setting:

Theorem 7 [Allison et al. 2002, Theorem 3.10]. Let $r \ge 3$. Then L is BC_r -graded with grading subalgebra \mathfrak{g} of type B_r if and only if there exists an associative algebra \mathfrak{a} with involution η , and an \mathfrak{a} -module C with a hermitian form χ such that L is centrally isogenous to the BC_r -graded Lie algebra $\mathfrak{so}_{2r+1}(\mathfrak{a}, \eta, C, \chi)$.

Since $im(A^{[d]})$ is BC_r -graded with a grading subalgebra of type B_r and is centrally closed, we have the following result.

Corollary 8. The intersection matrix algebra $im(A^{[d]})$ is isomorphic to the universal covering algebra of the Lie algebra $so_{2r+1}(\mathfrak{a}, \eta, C, \chi)$. In particular, there exists a graded, central, surjective Lie algebra homomorphism

$$\psi:\mathfrak{im}(A^{[d]})\to \mathrm{so}_{2r+1}(\mathfrak{a},\eta,C,\chi).$$

5. Arriving at a "minimal" understanding of a, η, C , and χ

The graded nature of the map $\psi : im(A^{[d]}) \to so_{2r+1}(\mathfrak{a}, \eta, C, \chi)$ along with the relations among the generating elements of $im(A^{[d]})$ allow us to study each of components \mathfrak{a}, η, C , and χ involved in $so_{2r+1}(\mathfrak{a}, \eta, C, \chi)$.

Since the elements $e_1 + \mathfrak{r}, \ldots, e_r + \mathfrak{r}, f_1 + \mathfrak{r}, \ldots, f_r + \mathfrak{r}, h_1 + \mathfrak{r}, \ldots, h_r + \mathfrak{r}$ in $\mathfrak{im}(A^{[d]})$ generate a subalgebra isomorphic to a simple Lie algebra of type B_r , and since ψ is a graded homomorphism, we may assume without loss of generality that (after relabeling the $e_i + \mathfrak{r}, f_i + \mathfrak{r}$, and $h_i + \mathfrak{r}$ as e_i, f_i , and h_i , respectively)

$$\begin{split} \psi(e_i) &= E_{i,i+1}^{\square}(1) \quad \text{for } 1 \le i \le r-1, \qquad \psi(e_r) = E_{r,r+1}^{\square}(\sqrt{2}), \\ \psi(f_i) &= E_{i+1,i}^{\square}(1) \quad \text{for } 1 \le i \le r-1, \qquad \psi(f_r) = E_{r+1,r}^{\square}(\sqrt{2}), \\ \psi(h_i) &= E_{i,i}^{\square}(1) + E_{i+1,i+1}^{\square}(-1) \text{ for } 1 \le i \le r-1, \quad \psi(h_r) = E_{r,r}^{\square}(2). \end{split}$$

Remark. Here we are using the notation established in Section 4. The generators of $\operatorname{im}(A^{[d]})$ coming from a simple root $\alpha_j \in \Pi$ are denoted by e_j , f_j , and h_j , while the generators coming from an *i*-th copy of an adjoined root $\alpha \in \Delta_{B_r}$ are denoted by $e_{\alpha,i}$, $f_{\alpha,i}$, and $h_{\alpha,si}$.

5.1. Understanding the invertibility of some coordinates of a.

Proposition 9. (i) Let $e_{\epsilon_p-\epsilon_q,i}$, $f_{\epsilon_p-\epsilon_q,i}$, $h_{\epsilon_p-\epsilon_q,i}$ be the generators of $\operatorname{im}(A^{[d]})$ that result from adjoining the *i*-th copy of a long root $\epsilon_p - \epsilon_q$ $(1 \le p, q \le r, p \ne q)$. If

$$\psi(e_{\epsilon_p-\epsilon_q,i}) = E_{p,q}^{\Box}(a)$$

for some $a \in \mathfrak{a}$, then a is an invertible element and

$$\psi(f_{\epsilon_p-\epsilon_q,i}) = E_{q,p}^{\Box}(a^{-1}).$$

(ii) Let $e_{\epsilon_p+\epsilon_q,i}$, $f_{\epsilon_p+\epsilon_q,i}$, $h_{\epsilon_p+\epsilon_q,i}$ be the generators of $\operatorname{im}(A^{[d]})$ that result from adjoining the *i*-th copy of a long root $\epsilon_p + \epsilon_q$ $(1 \le p, q \le r, p \ne q)$. If

$$\psi(e_{\epsilon_p+\epsilon_q,i}) = E_{p,2r+2-q}^{\square}(b)$$

for some $b \in \mathfrak{a}$, then b is an invertible element and

$$\psi(f_{\epsilon_p+\epsilon_q,i}) = E_{2r+2-q,p}^{\square}(b^{-1}).$$

(iii) Let $e_{-\epsilon_p-\epsilon_q,i}$, $f_{-\epsilon_p-\epsilon_q,i}$, $h_{-\epsilon_p-\epsilon_q,i}$ be the generators of $im(A^{[d]})$ that result from adjoining the *i*-th copy of a long root $-\epsilon_p - \epsilon_q$ ($1 \le p, q \le r, p \ne q$). If

$$\psi(e_{-\epsilon_p-\epsilon_q,i}) = E_{2r+2-p,q}(c)$$

for some $c \in \mathfrak{a}$, then c is an invertible element and

$$\psi(f_{-\epsilon_p-\epsilon_q,i}) = E_{q,2r+2-p}^{\square}(c^{-1}).$$

Proof. (i) Since ψ is a graded homomorphism,

$$\psi(e_{\epsilon_p-\epsilon_q,i}) = E_{p,q}^{\square}(a) \text{ and } \psi(f_{\epsilon_p-\epsilon_q,i}) = E_{q,p}^{\square}(a')$$

for some $a, a' \in \mathfrak{a}$. Without loss of generality, assume that p < q. Then

$$\begin{split} \left[\left[\psi(e_{\epsilon_{p}-\epsilon_{q},i}), \psi(f_{\epsilon_{p}-\epsilon_{q},i}) \right], \psi(e_{q}) \right] \\ &= \begin{cases} \left[\left[E_{p,q}^{\Box}(a), E_{q,p}^{\Box}(a') \right], E_{q,q+1}^{\Box}(1) \right] & \text{if } q < r, \\ \left[\left[E_{p,q}^{\Box}(a), E_{q,p}^{\Box}(a') \right], E_{r,r+1}^{\Box}(\sqrt{2}) \right] & \text{if } q = r, \end{cases} \\ &= \begin{cases} E_{q,q+1}^{\Box}(-a'a) & \text{if } q < r, \\ E_{r,r+1}^{\Box}(-\sqrt{2}a'a) & \text{if } q = r. \end{cases} \end{split}$$

But since

$$\begin{bmatrix} e_{\epsilon_p - \epsilon_q, i}, f_{\epsilon_p - \epsilon_q, i} \end{bmatrix} = h_{\epsilon_p - \epsilon_q, i} = \begin{bmatrix} h_{\epsilon_p - \epsilon_q, i}, e_q \end{bmatrix} = \begin{cases} A_{\epsilon_p - \epsilon_q, \epsilon_q - \epsilon_{q+1}} e_q & \text{if } q < r, \\ A_{\epsilon_p - \epsilon_q, \epsilon_r} e_q & \text{if } q = r \end{cases} = -e_q$$

and ψ is a homomorphism,

$$\begin{bmatrix} \begin{bmatrix} \psi(e_{\epsilon_p-\epsilon_q,i}), \psi(f_{\epsilon_p-\epsilon_q,i}) \end{bmatrix}, \psi(e_q) \end{bmatrix} = -\psi(e_q) = \begin{cases} E_{q,q+1}^{\Box}(-1) & \text{if } q < r, \\ E_{r,r+1}^{\Box}(-\sqrt{2}) & \text{if } q = r, \end{cases}$$

So whether q < r or q = r, we have

We show that aa' also equals 1. Indeed,

$$\begin{bmatrix} \begin{bmatrix} \psi(e_{\epsilon_p-\epsilon_q,i}), \psi(f_{\epsilon_p-\epsilon_q,i}) \end{bmatrix}, \psi(e_p) \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} E_{p,q}(a), E_{q,p}(a') \end{bmatrix}, E_{p,p+1}(1) \end{bmatrix}$$
$$= \begin{cases} E_{p,p+1}(aa') & \text{if } q-p \ge 2, \\ E_{p,p+1}(aa'+a'a) & \text{if } q=p+1. \end{cases}$$

But because

$$\begin{bmatrix} [e_{\epsilon_p-\epsilon_q,i}, f_{\epsilon_p-\epsilon_q,i}], e_p \end{bmatrix} = A_{\epsilon_p-\epsilon_q,\epsilon_p-\epsilon_{p+1}}e_p$$
$$= (1+\delta_{q,p+1})e_p = \begin{cases} e_p & \text{if } q \ge p+2, \\ 2e_p & \text{if } q = p+1, \end{cases}$$

we have

$$\begin{bmatrix} \psi(e_{\epsilon_p-\epsilon_q,i}), \psi(f_{\epsilon_p-\epsilon_q,i}) \end{bmatrix}, \psi(e_p) = \begin{cases} E_{p,p+1}^{\square}(1) & \text{if } q \ge p+2, \\ E_{p,p+1}^{\square}(2) & \text{if } q = p+1. \end{cases}$$

So if $q \ge p+2$, then aa' = 1. If q = p+1, then aa' + a'a = 2. But, by (1), a'a = 1. Hence, in either case, aa' = 1.

(ii) Since ψ is a graded homomorphism,

$$\psi(e_{\epsilon_p+\epsilon_q,i}) = E_{p,2r+2-q}^{\square}(b)$$
 and $\psi(f_{\epsilon_p+\epsilon_q,i}) = E_{2r+2-q,p}^{\square}(b')$

for some $b, b' \in \mathfrak{a}$. Again without loss of generality, we may assume that p < q. Then $\left[\left[\psi(e_{\epsilon_p + \epsilon_q, i}), \psi(f_{\epsilon_p + \epsilon_q, i}) \right], \psi(e_q) \right]$ equals

$$\begin{cases} \left[\left[E_{p,2r+2-q}^{\Box}(b), E_{2r+2-q,p}^{\Box}(b') \right], E_{q,q+1}^{\Box}(1) \right] & \text{if } q < r, \\ \left[\left[E_{p,2r+2-q}^{\Box}(b), E_{2r+2-q,p}^{\Box}(b') \right], E_{r,r+1}^{\Box}(\sqrt{2}) \right] & \text{if } q = r, \\ \\ = \begin{cases} E_{q,q+1}^{\Box} \left(\eta(b) \eta(b') \right) & \text{if } q < r, \\ E_{r,r+1}^{\Box} \left(\sqrt{2}\eta(b) \eta(b') \right) & \text{if } q = r. \end{cases} \end{cases}$$

But it also equals

$$\psi\left(\left[\left[e_{\epsilon_{p}+\epsilon_{q},i}, f_{\epsilon_{p}+\epsilon_{q},i}\right], e_{q}\right]\right) = \psi(e_{q}) = \begin{cases} E_{q,q+1}^{\Box}(1) & \text{if } q < r, \\ E_{r,r+1}^{\Box}(\sqrt{2}) & \text{if } q = r, \end{cases}$$

whence $\eta(b) \eta(b') = 1$. Applying (the antiautomorphism) η to both sides, we get that

$$b'b = 1.$$

To show that bb' = 1, we first compute that

$$\begin{bmatrix} \left[\psi(e_{\epsilon_p + \epsilon_q, i}), \psi(f_{\epsilon_p + \epsilon_q, i}) \right], \psi(e_p) \end{bmatrix} = \begin{cases} E_{p, p+1}^{\bullet}(bb') & \text{if } q \ge p+2, \\ E_{p, p+1}^{\bullet}(bb' - \eta(b) \eta(b')) & \text{if } q = p+1. \end{cases}$$

Since

$$\begin{bmatrix} e_{\epsilon_p+\epsilon_q,i}, f_{\epsilon_p+\epsilon_q,i} \end{bmatrix}, e_p \end{bmatrix} = A_{\epsilon_p+\epsilon_q,\epsilon_p-\epsilon_{p+1}} e_p = (1-\delta_{q,p+1})e_p = \begin{cases} e_p & \text{if } q \ge p+2, \\ 0 & \text{if } q = p+1, \end{cases}$$

we also have

$$\left[\left[\psi(e_{\epsilon_p+\epsilon_q,i}),\psi(f_{\epsilon_p+\epsilon_q,i})\right],\psi(e_p)\right] = \begin{cases} E_{p,p+1}^{\square}(1) & \text{if } q-p \ge 2, \\ 0 & \text{if } q=p+1. \end{cases}$$

So if $q \ge p+2$, then bb' = 1. If q = p+1, then $bb' - \eta(b) \eta(b') = 0$, which implies, using (2), that

$$bb' = \eta(b) \eta(b') = \eta(b'b) = \eta(1) = 1.$$

In either case, bb' = 1.

(iii) The proof follows using similar calculations as above.

5.2. Understanding the involution η on α .

Proposition 10. (i) *If*

$$\psi(e_{\epsilon_p+\epsilon_{p+1},i}) = E_{p,2r+2-(p+1)}^{\square}(a)$$

for some $1 \le p \le r - 1$ and $a \in \mathfrak{a}$, then $\eta(a) = a$.

(ii) If

$$\psi(e_{-\epsilon_p-\epsilon_{p+1},i}) = E_{2r+2-p,p+1}(b)$$

for some $1 \le p \le r - 1$ and $b \in \mathfrak{a}$, then $\eta(b) = b$.

Proof. We prove (i). The proof of (ii) is similar. Observe that

$$\left[\psi(e_{\epsilon_p+\epsilon_{p+1},i}),\psi(e_p)\right] = \left[E_{p,2r+2-(p+1)}(a), E_{p,p+1}(1)\right] = E_{p,2r+2-p}(a).$$

But $A_{\epsilon_p+\epsilon_{p+1},\epsilon_p-\epsilon_{p+1}} = 0$ implies that $(\operatorname{ad} e_{\epsilon_p+\epsilon_{p+1},i})^{-0+1}e_p = [e_{\epsilon_p+\epsilon_{p+1},i}, e_p] = 0$, which, in turn, implies that $[\psi(e_{\epsilon_p+\epsilon_{p+1},i}), \psi(e_p)] = 0$. So

$$E_{p,2r+2-p}^{\Box}(a) = E_{p,2r+2-p}(a-\eta(a)) = 0$$

and thus

$$\eta(a) = a.$$

5.3. Understanding the relations on generators of a.

Proposition 11. If, as a consequence of adjoining an *i*-th copy of the long root $\epsilon_p - \epsilon_q$ and a *j*-th copy of the long root $\epsilon_p + \epsilon_q$, where $1 \le p, q \le r$ with $p \ne q$,

$$\psi(e_{\epsilon_p-\epsilon_q,i}) = E_{p,q}^{\square}(s) \text{ and } \psi(e_{\epsilon_p+\epsilon_q,j}) = E_{p,2r+2-q}^{\square}(t),$$

for some $s, t \in \mathfrak{a}$, then:

(a) If |p-q| = 1, the elements s, t, and $\eta(s)$ in a satisfy the relation

$$s \cdot t = t \cdot \eta(s).$$

(b) If $|p-q| \ge 2$, the elements s, t, $\eta(s)$, and $\eta(t)$ in a satisfy the relation

$$s \cdot \eta(t) = t \cdot \eta(s).$$

Proof. Observe that

$$\begin{split} \left[\psi(e_{\epsilon_p - \epsilon_q, i}), \psi(e_{\epsilon_p + \epsilon_q, j}) \right] &= \left[E_{p,q}^{\square}(s), E_{p,2r+2-q}^{\square}(t) \right] \\ &= E_{p,2r+2-p}^{\square}(-s \cdot \eta(t)) \\ &= E_{p,2r+2-p}(-s \cdot \eta(t) + t \cdot \eta(s)) \\ &= \begin{cases} E_{p,2r+2-p}(-s \cdot t + t \cdot \eta(s)) & \text{if } |p-q| = 1, \\ E_{p,2r+2-p}(-s \cdot \eta(t) + t \cdot \eta(s)) & \text{if } |p-q| \ge 2. \end{cases} \end{split}$$

(The division into two cases in the last step follows from the use of (3).) But since $A_{\epsilon_p-\epsilon_q,\epsilon_p+\epsilon_q} = 0$, the generalized intersection matrix algebra relations tell us that

$$\left(\operatorname{ad} e_{\epsilon_p-\epsilon_q,i}\right)^{-0+1}e_{\epsilon_p+\epsilon_q,j}=0.$$

That is, $[e_{\epsilon_p-\epsilon_q,i}, e_{\epsilon_p+\epsilon_q,j}] = 0$. So we must have $[\psi(e_{\epsilon_p-\epsilon_q,i}), \psi(e_{\epsilon_p+\epsilon_q,j})] = 0$. This implies that

$$-s \cdot t + t \cdot \eta(s) = 0 \quad \text{if } |p-q| = 1,$$

$$-s \cdot \eta(t) + t \cdot \eta(s) = 0, \quad \text{if } |p-q| \ge 2.$$

Similarly:

Proposition 12. If, as a consequence of adjoining an *i*-th copy of the long root $\epsilon_p - \epsilon_q$ and a *j*-th copy of the long root $-\epsilon_p - \epsilon_q$, where $1 \le p, q \le r$ with $p \ne q$,

$$\psi(e_{\epsilon_p-\epsilon_q,i}) = E_{p,q}^{\square}(s) \quad and \quad \psi(e_{-\epsilon_p-\epsilon_q,j}) = E_{2r+2-p,q}^{\square}(t),$$

for some $s, t \in \mathfrak{a}$, then:

(a) If |p - q| = 1, the elements s, t, and $\eta(s)$ in a satisfy the relation

$$\eta(s) \cdot t = t \cdot s.$$

(b) If $|p-q| \ge 2$, the elements s, t, $\eta(s)$, and $\eta(t)$ in a satisfy the relation

$$\eta(s) \cdot t = \eta(t) \cdot s.$$

5.4. A description of the module C. Since ψ is a graded, surjective homomorphism from $\operatorname{im}(A^{[d]})$ to $\operatorname{so}_{2r+1}(\mathfrak{a}, \eta, C, \chi)$ and we are only adjoining long roots, we can examine the image of $\operatorname{im}(A^{[d]})$ under ψ to help us understand C.

Proposition 13. The module C is zero.

Proof. The generators of $\psi(\operatorname{im}(A^{[d]}))$ all have the form $\begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix}$, $M \in M_{2r+1}(\mathfrak{a})$. Since the matrices of this form in $\operatorname{so}_{2r+1}(\mathfrak{a}, \eta, C, \chi)$ form a subalgebra and since ψ is surjective, we have C = 0.

6. Achieving a "sufficient" understanding of a, η, C , and χ

In the previous section we used the homomorphism

$$\psi:\mathfrak{im}(A^{[d]})\to \mathrm{so}_{2r+1}(\mathfrak{a},\eta,C,\chi),$$

given by the recognition theorem of [Allison et al. 2002] to get a sense (i) of what the generators of a ought to be, (ii) of what the involution η on a ought to be, (iii) of what the relations on the generators of a ought to be, and (iv) that C = 0 and $\chi = 0$.

In this section, we show that the understanding we have arrived at is complete. We do so as follows:

1. Take the 4-tuple of associative algebra, involution, module, and hermitian form as we presently understand it. That is:

(i) Let Ω = the set of all long roots of the form ±(ϵ_i+ϵ_{i+1}) that we have adjoined,
 Θ = the set of all long roots in Δ_B which we have adjoined but that are not in Ω,

and let

$$\begin{aligned} X_e &= \bigcup_{\omega \in \Omega} \left\{ x_{\omega,1}, \dots, x_{\omega,N_\omega} \right\}, \quad X_f &= \bigcup_{\omega \in \Omega} \left\{ x_{\omega,1}^{-1}, \dots, x_{\omega,N_\omega}^{-1} \right\}, \\ Y_e &= \bigcup_{\theta \in \Theta} \left\{ y_{\theta,1}, \dots, y_{\theta,N_\theta} \right\}, \qquad Y_f &= \bigcup_{\theta \in \Theta} \left\{ y_{\theta,1}^{-1}, \dots, y_{\theta,N_\theta}^{-1} \right\}, \\ Z_e &= \bigcup_{\theta \in \Theta} \left\{ z_{\theta,1}, \dots, z_{\theta,N_\theta} \right\}, \qquad Z_f &= \bigcup_{\theta \in \Theta} \left\{ z_{\theta,1}^{-1}, \dots, z_{\theta,N_\theta}^{-1} \right\}, \end{aligned}$$

denote collections of indeterminates indexed by the sets Ω and Θ . Let \mathfrak{b} be the unital associative \mathbb{C} -algebra generated by the indeterminates in

$$X_e \cup X_f \cup Y_e \cup Y_f \cup Z_e \cup Z_f,$$

subject to the relations

$$y_{\epsilon_{p}-\epsilon_{q},i}x_{\epsilon_{p}+\epsilon_{q},j} = x_{\epsilon_{p}+\epsilon_{q},j}z_{\epsilon_{p}-\epsilon_{q},i},$$

$$y_{\epsilon_{p}-\epsilon_{q},i}z_{\epsilon_{p}+\epsilon_{q},j} = y_{\epsilon_{p}+\epsilon_{q},j}z_{\epsilon_{p}-\epsilon_{q},i},$$

$$z_{\epsilon_{p}-\epsilon_{q},i}x_{-\epsilon_{p}-\epsilon_{q},k} = x_{-\epsilon_{p}-\epsilon_{q},k}y_{\epsilon_{p}-\epsilon_{q},i},$$

$$z_{\epsilon_{p}-\epsilon_{q},i}y_{-\epsilon_{p}-\epsilon_{q},k} = z_{-\epsilon_{p}-\epsilon_{q},k}y_{\epsilon_{p}-\epsilon_{q},i},$$

where $i = 1, ..., N_{\epsilon_p - \epsilon_q}$ for $\epsilon_p - \epsilon_q \in \Theta$, $j = 1, ..., N_{\epsilon_p + \epsilon_q}$ for $\epsilon_p + \epsilon_q \in \Omega \cup \Theta$, and $k = 1, ..., N_{-\epsilon_p - \epsilon_q}$ for $-\epsilon_p - \epsilon_q \in \Omega \cup \Theta$.

(ii) Define an involution, which we also call η and sometimes denote by $\overline{\cdot}$, on b by

$$\begin{aligned} \eta(x_{\omega,i}) &= x_{\omega,i} & \text{if } \omega \in \Omega \text{ and } 1 \le i \le N_{\omega}, \\ \eta(y_{\theta,i}) &= z_{\theta,i} & \text{if } \theta \in \Theta \text{ and } 1 \le i \le N_{\theta}, \\ \eta(z_{\theta,i}) &= y_{\theta,i} & \text{if } \theta \in \Theta \text{ and } 1 \le i \le N_{\theta}. \end{aligned}$$

- (iii) Let C = 0 be the trivial b-module.
- (iv) Let $\chi = 0$ be a hermitian form on *C*.
- **Remarks.** (a) The indeterminates in $X_e \cup X_f \cup \cdots \cup Z_f$ are intended to capture the elements of the form a, a', b, b', c, and c' of a that we studied in Section 5, which arose from the images of the map ψ .
- (b) In the relations listed above, our use of the indeterminates $x_{\epsilon_p+\epsilon_q,j}$ and $x_{-\epsilon_p-\epsilon_q,j}$ signals that we are working with roots in Ω and, hence, |p-q| = 1 in this setting. Likewise, our use of the indeterminates $y_{\epsilon_p+\epsilon_q,j}$, $z_{\epsilon_p+\epsilon_q,j}$, $y_{-\epsilon_p-\epsilon_q,j}$, and $z_{-\epsilon_p-\epsilon_q,j}$ signals that we are working with roots in Θ and p, q such that $|p-q| \ge 2$.
- 2. Construct a map

$$\varphi:\mathfrak{gim}(A^{[d]}) \to \mathrm{so}_{2r+1}(\mathfrak{b},\eta,C,\chi)$$

sending the generators

$$e_{1}, \dots, e_{r}, \quad \bigcup_{\omega \in \Omega} \{e_{\omega,1}, \dots, e_{\omega,N_{\omega}}\}, \quad \bigcup_{\theta \in \Theta} \{e_{\theta,1}, \dots, e_{\theta,N_{\theta}}\},$$

$$f_{1}, \dots, f_{r}, \quad \bigcup_{\omega \in \Omega} \{f_{\omega,1}, \dots, f_{\omega,N_{\omega}}\}, \quad \bigcup_{\theta \in \Theta} \{f_{\theta,1}, \dots, f_{\theta,N_{\theta}}\},$$

$$h_{1}, \dots, h_{r}, \quad \bigcup_{\omega \in \Omega} \{h_{\omega,1}, \dots, h_{\omega,N_{\omega}}\}, \quad \bigcup_{\theta \in \Theta} \{h_{\theta,1}, \dots, h_{\theta,N_{\theta}}\},$$

of $gim(A^{[d]})$ to

$$\begin{split} \tilde{e}_{1}, \dots, \tilde{e}_{r}, & \bigcup_{\omega \in \Omega} \{ \tilde{e}_{\omega,1}, \dots, \tilde{e}_{\omega,N_{\omega}} \}, & \bigcup_{\theta \in \Theta} \{ \tilde{e}_{\theta,1}, \dots, \tilde{e}_{\theta,N_{\theta}} \} \\ \tilde{f}_{1}, \dots, \tilde{f}_{r}, & \bigcup_{\omega \in \Omega} \{ \tilde{f}_{\omega,1}, \dots, \tilde{f}_{\omega,N_{\omega}} \}, & \bigcup_{\theta \in \Theta} \{ \tilde{f}_{\theta,1}, \dots, \tilde{f}_{\theta,N_{\theta}} \}, \\ \tilde{h}_{1}, \dots, \tilde{h}_{r}, & \bigcup_{\omega \in \Omega} \{ \tilde{h}_{\omega,1}, \dots, \tilde{h}_{\omega,N_{\omega}} \}, & \bigcup_{\theta \in \Theta} \{ \tilde{h}_{\theta,1}, \dots, \tilde{h}_{\theta,N_{\theta}} \}, \end{split}$$

respectively, where

$$\begin{split} \tilde{e}_{i} &:= E_{i,i+1}^{\square}(1), \quad 1 \leq i \leq r-1, \\ \tilde{e}_{r} &:= E_{r,r+1}^{\square}(\sqrt{2}), \\ \tilde{e}_{\omega,i} &:= \begin{cases} E_{p,2r+2-(p+1)}^{\square}(x_{\omega,i}) & \text{if } \omega = \epsilon_{p} + \epsilon_{p+1}, \\ E_{2r+2-p,p+1}^{\square}(x_{\omega,i}) & \text{if } \omega = -\epsilon_{p} - \epsilon_{p+1}, \end{cases} \text{ for } \omega \in \Omega \text{ and } 1 \leq i \leq N_{\omega}, \\ \tilde{e}_{\theta,i} &:= \begin{cases} E_{p,q}^{\square}(y_{\theta,i}) & \text{if } \theta = \epsilon_{p} - \epsilon_{q}, \\ E_{p,2r+2-q}^{\square}(y_{\theta,i}) & \text{if } \theta = \epsilon_{p} + \epsilon_{q}, \\ E_{2r+2-p,q}^{\square}(y_{\theta,i}) & \text{if } \theta = -\epsilon_{p} - \epsilon_{q}, \end{cases} \end{split}$$

$$\begin{split} \tilde{f}_{i} &:= E_{i+1,i}^{-}(1), \quad 1 \leq i \leq r-1, \\ \tilde{f}_{r} &:= E_{r+1,r}^{-}(\sqrt{2}), \\ \tilde{f}_{\omega,i} &:= \begin{cases} E_{2r+2-(p+1),p}^{-}(x_{\omega,i}^{-1}) & \text{if } \omega = \epsilon_{p} + \epsilon_{p+1}, \\ E_{p+1,2r+2-p}^{-}(x_{\omega,i}^{-1}) & \text{if } \omega = -\epsilon_{p} - \epsilon_{p+1}, \end{cases} \text{ for } \omega \in \Omega \text{ and } 1 \leq i \leq N_{\omega}, \\ \tilde{f}_{\theta,i} &:= \begin{cases} E_{q,p}^{-}(y_{\theta,i}^{-1}) & \text{if } \theta = \epsilon_{p} - \epsilon_{q}, \\ E_{2r+2-q,p}^{-}(y_{\theta,i}^{-1}) & \text{if } \theta = -\epsilon_{p} - \epsilon_{q}, \\ E_{q,2r+2-p}^{-}(y_{\theta,i}^{-1}) & \text{if } \theta = -\epsilon_{p} - \epsilon_{q}, \end{cases} \\ \tilde{h}_{i} &:= E_{i,i}^{-}(1) + E_{p+1,p+1}^{-}(1) & \text{if } \omega = \epsilon_{p} + \epsilon_{p+1} \\ E_{p,p}^{-}(-1) + E_{p+1,p+1}^{-}(-1) & \text{if } \omega = -\epsilon_{p} - \epsilon_{p+1} \end{cases} \text{ for } \omega \in \Omega, 1 \leq i \leq N_{\omega}, \\ \tilde{h}_{\theta,i} &:= \begin{cases} E_{p,p}^{-}(1) + E_{p+1,p+1}^{-}(-1) & \text{if } \theta = -\epsilon_{p} - \epsilon_{p+1} \\ E_{p,p}^{-}(-1) + E_{p+1,p+1}^{-}(-1) & \text{if } \theta = -\epsilon_{p} - \epsilon_{p+1} \end{cases} \text{ for } \omega \in \Omega, 1 \leq i \leq N_{\omega}, \\ \tilde{h}_{\theta,i} &:= \begin{cases} E_{p,p}^{-}(1) + E_{q,q}^{-}(-1) & \text{if } \theta = \epsilon_{p} - \epsilon_{q}, \\ E_{p,p}^{-}(-1) + E_{q,q}^{-}(-1) & \text{if } \theta = \epsilon_{p} - \epsilon_{q}, \end{cases} \end{cases}$$

- 3. We show that φ is
- (a) a Lie algebra homomorphism (Theorem 14),
- (b) that is surjective (Proposition 15), and
- (c) graded (Proposition 16).

4. We show that the radical \mathfrak{r} of $\mathfrak{gim}(A^{[d]})$ lies in the kernel of this map φ (see just before Proposition 17), hence inducing a surjective, graded, Lie algebra homomorphism

$$\phi: \mathfrak{im}(A^{\lfloor d \rfloor}) \to \operatorname{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi).$$

5. Finally, we show that ϕ is a central map and that $\mathfrak{b} \cong \mathfrak{a}$ (Proposition 18).

Theorem 14. The map $\varphi : \mathfrak{gim}(A^{[d]}) \to \mathfrak{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$ is a Lie algebra homomorphism.

Proof. We show that the images in $so_{2r+1}(\mathfrak{b}, \eta, C, \chi)$ of the generators of $\mathfrak{gim}(A^{[d]})$, under the map φ , satisfy the relations (R1)–(R3) of Definition 2 with respect to

the same $(r + d) \times (r + d)$ generalized intersection matrix $A^{[d]}$ as used in the construction of the algebra $gim(A^{[d]})$.

While working with the various long roots in our proof, we use labels like *u* or *v* to denote the indeterminates $x_{\omega,i}$ or $y_{\theta,i}$.

The reason that we can substitute u or v for the actual indeterminates is that the result of taking a bracket like

$$\left[\tilde{e}_{-\epsilon_p-\epsilon_q,i}, \tilde{e}_{-\epsilon_k-\epsilon_l,j}\right] = \left[E_{2r+2-p,q}^{\square}(y_{-\epsilon_p-\epsilon_q,i}), E_{2r+2-k,l}^{\square}(y_{-\epsilon_k-\epsilon_l,j})\right]$$

depends primarily on the indices p, q, k, and l rather than on the particular elements of the algebra b being housed at these sites.

If we agree on this convention of using substitute variables like u, then we must recognize that

$$\overline{u} = \begin{cases} x_{\omega,i} & \text{if } u = x_{\omega,i}, \\ z_{\theta,i} & \text{if } u = y_{\theta,j}. \end{cases}$$

That is, the involution $\overline{\cdot}$ applied to *u* depends on whether *u* is substituting for a variable associated to a root in Ω or a root in Θ .

We show the computations for the interactions between the generators corresponding to the long roots $\epsilon_p - \epsilon_q$ and $\epsilon_k - \epsilon_l$. The remaining computations are similar.

Let $1 \leq p, q, k, l \leq r$ with $p \neq q$ and $k \neq l, u, v \in \{x_{\omega,i}, x_{\omega,j}, y_{\theta,i}y_{\theta,j}\}$ and $u^{-1}, v^{-1} \in \{x_{\omega,i}^{-1}, x_{\omega,j}^{-1}, y_{\theta,i}^{-1}y_{\theta,j}^{-1}\}$, where $\omega \in \Omega, \theta \in \Theta$, and $1 \leq i, j \leq N_{\omega}$ or $1 \leq i, j \leq N_{\theta}$.

Using the definition of $A_{\epsilon_p-\epsilon_q,\epsilon_k-\epsilon_l}$, we see that

$$A_{\epsilon_p-\epsilon_q,\epsilon_k-\epsilon_l} = \delta_{p,k} - \delta_{p,l} - \delta_{q,k} + \delta_{q,l} = \begin{cases} 0 & \text{if } p, q \notin \{k, l\}, \\ 1 & \text{if } p = k \text{ but } q \neq l, \\ -1 & \text{if } p = l \text{ but } q \neq k, \\ -1 & \text{if } p \neq l \text{ but } q = k, \\ 1 & \text{if } p \neq k \text{ but } q = l, \\ 2 & \text{if } p = k \text{ and } q = l, \\ -2 & \text{if } p = l \text{ and } q = k. \end{cases}$$

A.
$$\begin{bmatrix} \tilde{e}_{\epsilon_p - \epsilon_q, i}, \tilde{e}_{\epsilon_k - \epsilon_l, j} \end{bmatrix} = \begin{bmatrix} E_{p,q}^{\bullet}(u), E_{k,l}^{\bullet}(v) \end{bmatrix} = \delta_{q,k} E_{p,l}^{\bullet}(uv) + \delta_{l,p} E_{k,q}^{\bullet}(-vu)$$
$$= \begin{cases} E_{k,q}^{\bullet}(-vu) & \text{if } p = l \text{ but } q \neq k, \\ E_{p,l}^{\bullet}(uv) & \text{if } p \neq l \text{ but } q = k, \\ E_{p,p}^{\bullet}(uv) + E_{q,q}^{\bullet}(-vu) & \text{if } p = l \text{ and } q = k, \\ 0 & \text{otherwise.} \end{cases}$$

• If p = l but $q \neq k$, then $\left[E_{p,q}^{\square}(u), E_{k,q}^{\square}(-vu)\right] = 0$ because $q \neq k$ and $q \neq p$.

If p ≠ l but q = k, then [E[□]_{p,q}(u), E[□]_{p,l}(uv)] = 0 because q ≠ p and l ≠ p.
If p = l and q = k, then

$$\begin{bmatrix} E_{p,q}^{\Box}(u), & E_{p,p}^{\Box}(uv) + E_{q,q}^{\Box}(-vu) \end{bmatrix} = E_{p,q}^{\Box}(-uvu) + E_{p,q}^{\Box}(-uvu)$$
$$= E_{p,q}^{\Box}(-2uvu).$$

So

$$\left(\operatorname{ad} \tilde{e}_{\epsilon_p - \epsilon_q, i}\right)^{1+1} \tilde{e}_{\epsilon_k - \epsilon_l, j} = \begin{cases} E_{p,q}^{\square}(-2uvu) & \text{if } p = l \text{ and } q = k, \\ 0 & \text{otherwise.} \end{cases}$$

Since
$$[E_{p,q}(u), E_{p,q}(-2uvu)] = 0$$
, we get $(\text{ad } \tilde{e}_{\epsilon_{p}-\epsilon_{q},i})^{2+1} \tilde{e}_{\epsilon_{k}-\epsilon_{l},j} = 0$.
B. $[\tilde{f}_{\epsilon_{p}-\epsilon_{q},i}, \tilde{f}_{\epsilon_{k}-\epsilon_{l},j}] = [E_{q,p}(u^{-1}), E_{l,k}(v^{-1})]$
 $= \delta_{p,l} E_{q,k}(u^{-1}v^{-1}) + \delta_{k,q} E_{l,p}(-v^{-1}u^{-1})$
 $= \begin{cases} E_{q,k}(u^{-1}v^{-1}) & \text{if } p = l \text{ but } q \neq k, \\ E_{l,p}(-v^{-1}u^{-1}) & \text{if } p \neq l \text{ but } q = k, \\ E_{p,p}(-v^{-1}u^{-1}) + E_{q,q}(u^{-1}v^{-1}) & \text{if } p = l \text{ and } q = k, \\ 0 & \text{otherwise.} \end{cases}$

- If p = l but $q \neq k$, then $\left[E_{q,p}^{\square}(u^{-1}), E_{q,k}^{\square}(u^{-1}v^{-1})\right] = 0$ because $p \neq q$ and $k \neq q$.
- If $p \neq l$ but q = k, then $\left[E_{q,p}^{\square}(u^{-1}), E_{l,p}^{\square}(-v^{-1}u^{-1})\right] = 0$ because $p \neq l$ and $p \neq q$.
- If p = l and q = k, then

$$\begin{split} \left[E_{q,p}^{\Box}(u^{-1}), E_{p,p}^{\Box}(-v^{-1}u^{-1}) + E_{q,q}^{\Box}(u^{-1}v^{-1}) \right] \\ &= E_{q,p}^{\Box}(-u^{-1}v^{-1}u^{-1}) + E_{q,p}^{\Box}(-u^{-1}v^{-1}u^{-1}) \\ &= E_{q,p}^{\Box}\left(-2u^{-1}v^{-1}u^{-1}\right). \end{split}$$

So

$$\left(\operatorname{ad} \tilde{f}_{\epsilon_p-\epsilon_q,i}\right)^{1+1} \tilde{f}_{\epsilon_k-\epsilon_l,j} = \begin{cases} E_{q,p}^{\bullet}(-2u^{-1}v^{-1}u^{-1}) & \text{if } p = l \text{ and } q = k, \\ 0 & \text{otherwise.} \end{cases}$$

Since $\left[E_{q,p}^{\square}(u^{-1}), E_{q,p}^{\square}(-2u^{-1}v^{-1}u^{-1})\right] = 0$, we get that (ad $\tilde{f}_{\epsilon_p-\epsilon_q,i})^{2+1}\tilde{f}_{\epsilon_k-\epsilon_l,j} = 0$.

$$\begin{aligned} \mathbf{C}. \quad & \left[\tilde{h}_{\epsilon_{p}-\epsilon_{q},i}, \tilde{h}_{\epsilon_{k}-\epsilon_{l},j}\right] \\ &= \left[E_{p,p}^{\Box}(1) + E_{q,q}^{\Box}(-1), E_{k,k}^{\Box}(1) + E_{l,l}^{\Box}(-1)\right] \\ &= \delta_{p,k} E_{p,p}^{\Box}([1,1]) + \delta_{p,l} E_{p,p}^{\Box}([1,-1]) + \delta_{q,k} E_{q,q}^{\Box}([-1,1]) + \delta_{q,l} E_{q,q}^{\Box}([-1,-1]) \\ &= \delta_{p,k} E_{p,p}^{\Box}(0) + \delta_{p,l} E_{p,p}^{\Box}(0) + \delta_{q,k} E_{q,q}^{\Box}(0) + \delta_{q,l} E_{q,q}^{\Box}(0) \\ &= 0 \end{aligned}$$

D.
$$\left[\tilde{e}_{\epsilon_{p}-\epsilon_{q},i}, \tilde{f}_{\epsilon_{k}-\epsilon_{l},j}\right] = \left[E_{p,q}^{\Box}(u), E_{l,k}^{\Box}(v^{-1})\right]$$

 $= \delta_{q,l}E_{p,k}^{\Box}\left(u\,v^{-1}\right) + \delta_{k,p}E_{l,q}^{\Box}\left(-v^{-1}\,u\right)$
 $= \begin{cases}E_{l,q}^{\Box}\left(-v^{-1}u\right) & \text{if } p = k \text{ but } q \neq l, \\E_{p,k}^{\Box}\left(uv^{-1}\right) & \text{if } p \neq k \text{ but } q = l, \\E_{p,p}^{\Box}\left(uv^{-1}\right) + E_{q,q}^{\Box}\left(-v^{-1}u\right) & \text{if } p = k \text{ and } q = l, \\0 & \text{otherwise.}\end{cases}$

• If
$$p = k$$
 but $q \neq l$, then $\left[E_{p,q}^{\square}(u), E_{l,q}^{\square}(-v^{-1}u)\right] = 0$ because $q \neq l$ and $q \neq p$.

If p ≠ k but q = l, then [E[□]_{p,q}(u), E[□]_{p,k}(uv⁻¹)] = 0 because q ≠ p and k ≠ p.
If p = k and q = l, then

$$\begin{bmatrix} E_{p,q}^{\Box}(u), E_{p,p}^{\Box}(uv^{-1}) + E_{q,q}^{\Box}(-v^{-1}u) \end{bmatrix} = E_{p,q}^{\Box}(-uv^{-1}u) + E_{p,q}^{\Box}(-uv^{-1}u)$$
$$= E_{p,q}^{\Box}(-2uv^{-1}u).$$

So

$$\left(\operatorname{ad} \tilde{e}_{\epsilon_p - \epsilon_q, i}\right)^{1+1} \tilde{f}_{\epsilon_k - \epsilon_l, j} = \begin{cases} E_{p,q}^{\square}(-2uv^{-1}u) & \text{if } p = k \text{ and } q = l, \\ 0 & \text{otherwise.} \end{cases}$$

Since $[E_{p,q}(u), E_{p,q}(-2uv^{-1}u)] = 0$, we get that $(\operatorname{ad} \tilde{e}_{\epsilon_{p}-\epsilon_{q},i})^{2+1} \tilde{f}_{\epsilon_{k}-\epsilon_{l},j} = 0$. E. $[\tilde{f}_{\epsilon_{p}-\epsilon_{q},i}, \tilde{e}_{\epsilon_{k}-\epsilon_{l},j}] = [E_{q,p}(u^{-1}), E_{k,l}(v)]$ $= \delta_{p,k} E_{q,l}(u^{-1}v) + \delta_{l,q} E_{k,p}(-vu^{-1})$ $= \begin{cases} E_{q,l}(u^{-1}v) & \text{if } p = k \text{ but } q \neq l, \\ E_{k,p}(-vu^{-1}) & \text{if } p \neq k \text{ but } q = l, \\ E_{p,p}(-vu^{-1}) + E_{q,q}(u^{-1}v) & \text{if } p = k \text{ and } q = l, \\ 0 & \text{otherwise.} \end{cases}$

• If p = k but $q \neq l$, then $\left[E_{q,p}^{\square}(u^{-1}), E_{q,l}^{\square}(u^{-1}v)\right] = 0$ because $p \neq q$ and $l \neq q$.

- If $p \neq k$ but q = l, then $\left[E_{q,p}^{\bullet}(u^{-1}), E_{k,p}^{\bullet}(-v u^{-1})\right] = 0$ because $p \neq k$ and $p \neq q$.
- If p = k and q = l, then

$$\begin{bmatrix} E_{q,p}^{\Box}(u^{-1}), E_{p,p}^{\Box}(-vu^{-1}) + E_{q,q}^{\Box}(u^{-1}v) \end{bmatrix} = E_{q,p}^{\Box}(-u^{-1}vu^{-1}) + E_{q,p}^{\Box}(-u^{-1}vu^{-1})$$
$$= E_{q,p}^{\Box}(-2u^{-1}vu^{-1}).$$

So

$$\left(\operatorname{ad} \tilde{f}_{\epsilon_p-\epsilon_q,i}\right)^{1+1} \tilde{e}_{\epsilon_k-\epsilon_l,j} = \begin{cases} E_{q,p}^{\bullet}(-2u^{-1}vu^{-1}) & \text{if } p=k \text{ and } q=l, \\ 0 & \text{otherwise.} \end{cases}$$

Since
$$\left[E_{q,p}^{\bullet}(u^{-1}), E_{q,p}^{\bullet}(-2u^{-1}vu^{-1})\right] = 0$$
, we get
(ad $\tilde{f}_{\epsilon_p-\epsilon_q,i})^{2+1}\tilde{e}_{\epsilon_k-\epsilon_l,j} = 0$.

$$\begin{aligned} \mathbf{F}.\left[\tilde{h}_{\epsilon_{p}-\epsilon_{q},i},\tilde{e}_{\epsilon_{k}-\epsilon_{l},j}\right] &= \begin{bmatrix} E_{p,p}^{\Box}(1) + E_{q,q}^{\Box}(-1), E_{k,l}^{\Box}(v) \end{bmatrix} \\ &= \delta_{p,k}E_{p,l}^{\Box}(v) + \delta_{l,p}E_{k,p}^{\Box}(-v) + \delta_{q,k}E_{q,l}^{\Box}(-v) + \delta_{l,q}E_{k,q}^{\Box}(v) \\ &= \delta_{p,k}E_{k,l}^{\Box}(v) - \delta_{p,l}E_{k,l}^{\Box}(v) - \delta_{q,k}E_{k,l}^{\Box}(v) + \delta_{q,l}E_{k,l}^{\Box}(v) \\ &= (\delta_{p,k} - \delta_{p,l} - \delta_{q,k} + \delta_{q,l})E_{k,l}^{\Box}(v) \end{aligned}$$

$$\begin{aligned} \mathbf{G}. \left[\tilde{h}_{\epsilon_{p}-\epsilon_{q},i}, \tilde{f}_{\epsilon_{k}-\epsilon_{l},j} \right] \\ &= \left[E_{p,p}^{\bullet}(1) + E_{q,q}^{\bullet}(-1), E_{l,k}^{\bullet}(v^{-1}) \right] \\ &= \delta_{p,l} E_{p,k}^{\bullet}(v^{-1}) + \delta_{k,p} E_{l,p}^{\bullet}(-v^{-1}) + \delta_{q,l} E_{q,k}^{\bullet}(-v^{-1}) + \delta_{k,q} E_{l,q}^{\bullet}(v^{-1}) \\ &= \delta_{p,l} E_{l,k}^{\bullet}(v^{-1}) - \delta_{p,k} E_{l,k}^{\bullet}(v^{-1}) - \delta_{q,l} E_{l,k}^{\bullet}(v^{-1}) + \delta_{q,k} E_{l,k}^{\bullet}(v^{-1}) \\ &= -(\delta_{p,k} - \delta_{p,l} - \delta_{q,k} + \delta_{q,l}) E_{l,k}^{\bullet}(v^{-1}) \\ &= -A_{\epsilon_{p}-\epsilon_{q},\epsilon_{k}-\epsilon_{l}} \tilde{f}_{\epsilon_{k}-\epsilon_{l},j} \end{aligned}$$

Proposition 15. The map φ : $\mathfrak{gim}(A^{[d]}) \to \mathfrak{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$ is a surjective Lie algebra homomorphism.

Proof. Let $\mathfrak{B} = \text{Im}(\varphi) \subseteq L$, where $L = \text{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$. We show that $\mathfrak{B} = L$ in a sequence of steps.

#1. Let $\mathfrak{g} = \{ M \in \mathcal{M}_{2r+1}(\mathbb{C}) : M^t G = -GM \}$ and $\mathfrak{s} = \{ M \in \mathcal{M}_{2r+1}(\mathbb{C}) : M^t G = GM, \operatorname{tr}(M) = 0 \}.$ Let $A = \{a \in \mathfrak{b} : \eta a = a\}$ and $B = \{b \in \mathfrak{b} : \eta b = -b\}$. If $0 \neq a \in A$, then $a\mathfrak{g}$ is an irreducible \mathfrak{g} -submodule of \mathfrak{B} with highest weight $\epsilon_1 + \epsilon_2$. If $0 \neq b \in \mathfrak{B}$, then $b\mathfrak{s}$ is an irreducible \mathfrak{g} -submodule of \mathfrak{B} with highest weight $2\epsilon_1$. These \mathfrak{g} -modules are not isomorphic.

- #2. \mathfrak{B} is a subalgebra of L containing \mathfrak{g} , so \mathfrak{B} is a \mathfrak{g} -submodule of L.
- #3. For $1 \le p, q \le 2r + 1, p \ne q, p \ne 2r + 2 q$, let

$$I_{pq} = \left\{ x \in \mathfrak{b} : E_{pq}(x) - E_{2r+2-q,2r+2-p}(\eta x) \in \mathfrak{B} \right\}.$$

Notice that I_{pq} is a subspace of \mathfrak{B} .

#4. I_{pq} is invariant under η . Indeed, let $x \in I_{pq}$, in which case

$$X := E_{pq}(x) - E_{2r+2-q,2r+2-p}(\eta x) \in \mathfrak{B}.$$

But $X = X_1 + X_2$, where $X_1 = \frac{1}{2}(x + \eta x) (E_{pq}(1) - E_{2r+2-q,2r+2-p}(1)) \in (x + \eta x) \mathfrak{g}$ and $X_2 = \frac{1}{2}(x - \eta x) (E_{pq}(1) + E_{2r+2-q,2r+2-p}(1)) \in (x - \eta x) \mathfrak{s}$. Thus, by #1 and #2, $X_1, X_2 \in \mathfrak{B}$. So $x + \eta x, x - \eta x \in I_{pq}$, which implies that $\eta x \in I_{pq}$.

#5. By #4, $I_{pq} = I_{pq} \cap A + I_{pq} \cap B$. But by #1 and #2, $I_{pq} \cap A$ and $I_{pq} \cap B$ are independent of p, q. So $I := I_{pq}$ is independent of p, q.

#6. We have

$$\left[E_{12}(x) - E_{2r,2r+1}(\eta x), E_{23}(y) - E_{2r-1,2r}(\eta y)\right] = E_{13}(xy) - E_{2r-1,2r+1}((\eta y)(\eta x)).$$

So, by #5, *I* is a subalgebra of \mathfrak{B} , and, by #4, *I* is invariant under η .

#7. The action of φ on the generators of $\mathfrak{gim}(A^{[d]})$ tells us that *I* contains the elements $x_{\omega,i}, x_{\omega,i}^{-1}, y_{\theta,i}, y_{\theta,i}^{-1}$. So by #6, $I = \mathfrak{B}$.

#8. By #7, we have $A\mathfrak{g} + B\mathfrak{s} \subseteq \mathfrak{B}$. But since $C = \{0\}$, we have $\sum_{\alpha \in \Delta} L_{\alpha} \subseteq A\mathfrak{g} + B\mathfrak{s}$. So $\sum_{\alpha \in \Delta} L_{\alpha} \subseteq \mathfrak{B}$. Hence, since \mathfrak{B} is a subalgebra of $L, \mathfrak{B} = L$.

Continuing our plan laid out on page 273, we next show that $\varphi : \mathfrak{gim}(A^{\lfloor d \rfloor}) \rightarrow \operatorname{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$ is a graded homomorphism and that it induces a map from $\operatorname{im}(A^{\lfloor d \rfloor})$ to $\operatorname{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$.

We saw in Sections 3 and 4, respectively, that $\mathfrak{gim}(A^{[d]})$ and $\mathfrak{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$ are both Γ -graded Lie algebras, where

$$\Gamma = \bigoplus_{\mu \in \Delta} \mathbb{Z} \alpha_{\mu}.$$

The map $\varphi : \mathfrak{gim}(A^{[d]}) \to \mathrm{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$ is engineered so that, for all $\alpha \in \Gamma$,

$$\varphi(\mathfrak{gim}(A^{[d]})_{\alpha}) \subset \mathrm{so}_{2r+1}(\mathfrak{b},\eta,C,\chi)_{\alpha}.$$

That is, the following result holds by design.

Proposition 16. The map $\varphi : \mathfrak{gim}(A^{[d]}) \to \mathfrak{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$ is also a graded homomorphism.

Moreover, since $so_{2r+1}(\mathfrak{b}, \eta, C, \chi)_{\gamma} = 0$ for $\gamma \notin \Delta \cup \{0\}$, we get that the radical \mathfrak{r} of $\mathfrak{gim}(A^{[d]})$ lies in the kernel of φ .

Proposition 17. There exists a surjective, graded Lie algebra homomorphism

$$\phi : \mathfrak{im}(A^{[d]}) \to \mathfrak{so}_{2r+1}(\mathfrak{b}, \eta, C, \chi)$$

given by $\phi(u + \mathfrak{r}) = \varphi(u)$ for any $u + \mathfrak{r} \in \mathfrak{im}(A^{[d]})$, where $u \in \mathfrak{gim}(A^{[d]})$.

We now turn to centrality. Let so(\mathfrak{a}) and so(\mathfrak{b}) be shorthand for so_{2*r*+1}(\mathfrak{a} , η , *C*, χ) and so_{2*r*+1}(\mathfrak{b} , η , *C*, χ), respectively.

Since the elements of a satisfy the defining relations of b, by universality, there exists a surjective associative algebra homomorphism $g: b \to a$. In particular, $g(x_{\omega,i}) = a \in a$ if $\psi(e_{\omega,i}) = E_{pq}(a) - E_{2r+2-q,2r+2-p}(\eta a)$, and $g(y_{\theta,i}) = b \in a$ if $\psi(e_{\theta,i}) = E_{pq}(b) - E_{2r+2-q,2r+2-p}(\eta b)$. This algebra homomorphism respects the involution and induces a surjective Lie algebra homomorphism

$$\widetilde{g}$$
: so(\mathfrak{b}) \rightarrow so(\mathfrak{a}) such that $\widetilde{g}\phi = \psi$.

Hence ker $\phi \subset$ ker $\psi \subset \mathfrak{z}(\mathfrak{im}(A^{[d]}))$, where $\mathfrak{z}(\mathfrak{im}(A^{[d]}))$ denotes the center of $\mathfrak{im}(A^{[d]})$. Thus ker $\phi \subset \mathfrak{z}(\mathfrak{im}(A^{[d]}))$, implying the following result:

Proposition 18. The map $\phi : \mathfrak{im}(A^{[d]}) \to \mathfrak{so}(\mathfrak{b})$ is a central homomorphism.

We also know that $\psi : \operatorname{im}(A^{[d]}) \to \operatorname{so}(\mathfrak{a})$ is a universal central extension: so there exists a Lie algebra homomorphism $\tilde{f} : \operatorname{so}(\mathfrak{a}) \to \operatorname{so}(\mathfrak{b})$ such that $\tilde{f}\psi = \phi$. Since ψ is surjective, the generators of \mathfrak{a} are of the form $a, a^{-1}, \eta(a)$, where a is the element in \mathfrak{a} corresponding to the image $\psi(e_{\lambda,i})$ of the *i*-th copy of a long root λ in Δ_B which was adjoined.

Since $\tilde{f}\psi = \phi$, the map \tilde{f} induces an associative algebra homomorphism $f: \mathfrak{a} \to \mathfrak{b}$ given by

$$f(a) = \begin{cases} x_{\omega,i} & \text{if } a \text{ is the element in a corresponding to} \\ & \text{the image } \psi(e_{\omega,i}) = E_{pq}(a) - E_{2r+2-q,2r+2-p}(\eta a) \\ & \text{for some } 1 \le p, q \le 2r+1 \text{ with } p \ne q, p \ne 2r+2-q \\ & \text{and } \omega \in \Omega, 1 \le i \le N_{\omega}, \end{cases} \\ y_{\theta,i} & \text{if } a \text{ is the element in a corresponding to} \\ & \text{the image } \psi(e_{\theta,i}) = E_{pq}(a) - E_{2r+2-q,2r+2-p}(\eta a) \\ & \text{for some } 1 \le p, q \le 2r+1 \text{ with } p \ne q, p \ne 2r+2-q \\ & \text{and } \theta \in \Theta, 1 \le i \le N_{\theta}. \end{cases}$$

We define $f(\eta a)$ to be $\eta f(a)$ so that f preserves the involution.

But then $g \circ f = id_{\mathfrak{a}}$ and $f \circ g = id_{\mathfrak{b}}$, that is, $\mathfrak{a} \cong \mathfrak{b}$ as associative algebras.

Acknowledgments

We thank Frank Lemire and Michael Lau for their helpful discussions, and the anonymous referees for their useful comments and shortened proofs.

References

- [Allison et al. 2002] B. Allison, G. Benkart, and Y. Gao, *Lie algebras graded by the root systems BC_r*, $r \ge 2$, Mem. Amer. Math. Soc. **751**, Amer. Math. Soc., Providence, RI, 2002. MR 2003h:17038 Zbl 0998.17031
- [Benkart and Zelmanov 1996] G. Benkart and E. Zelmanov, "Lie algebras graded by finite root systems and intersection matrix algebras", *Invent. Math.* **126**:1 (1996), 1–45. MR 97k:17044 Zbl 0871.17024
- [Berman 1989] S. Berman, "On generators and relations for certain involutory subalgebras of Kac-Moody Lie algebras", *Comm. Algebra* **17**:12 (1989), 3165–3185. MR 91a:17030 Zbl 0693.17012
- [Berman and Moody 1992] S. Berman and R. V. Moody, "Lie algebras graded by finite root systems and the intersection matrix algebras of Slodowy", *Invent. Math.* **108**:2 (1992), 323–347. MR 93e:17031 Zbl 0778.17018
- [Berman et al. 2001] S. Berman, E. Jurisich, and S. Tan, "Beyond Borcherds Lie algebras and inside", *Trans. Amer. Math. Soc.* **353**:3 (2001), 1183–1219. MR 2001f:17048 Zbl 0973.17031
- [Carter 2005] R. W. Carter, *Lie algebras of finite and affine type*, Cambridge Studies in Advanced Mathematics **96**, Cambridge University Press, 2005. MR 2006i:17001 Zbl 1110.17001
- [Gabber and Kac 1981] O. Gabber and V. G. Kac, "On defining relations of certain infinite-dimensional Lie algebras", *Bull. Amer. Math. Soc.* (*N.S.*) **5**:2 (1981), 185–189. MR 84b:17011 ZbI 0474.17007
- [Gao 1996] Y. Gao, "Involutive Lie algebras graded by finite root systems and compact forms of IM algebras", *Math. Z.* **223**:4 (1996), 651–672. MR 98b:17026 Zbl 0884.17012
- [Kac 1990] V. G. Kac, *Infinite-dimensional Lie algebras*, 3rd ed., Cambridge University Press, 1990. MR 92k:17038 Zbl 0716.17022
- [Moody and Pianzola 1995] R. V. Moody and A. Pianzola, *Lie algebras with triangular decompositions*, Canadian Mathematical Society Series of Monographs and Advanced Texts, Wiley, New York, 1995. MR 96d:17025 Zbl 0874.17026
- [Neher 1996] E. Neher, "Lie algebras graded by 3-graded root systems and Jordan pairs covered by grids", *Amer. J. Math.* **118**:2 (1996), 439–491. MR 97a:17029 Zbl 0857.17019
- [Peng 2002] L. Peng, "Intersection matrix Lie algebras and Ringel–Hall Lie algebras of tilted algebras", pp. 98–108 in *Representations of algebra, I*, edited by D. Happel and Y. B. Zhang, Beijing Norm. Univ. Press, Beijing, 2002. MR 2005m:16022 Zbl 1086.16517
- [Rao et al. 1992] S. E. Rao, R. V. Moody, and T. Yokonuma, "Lie algebras and Weyl groups arising from vertex operator representations", *Nova J. Algebra Geom.* **1**:1 (1992), 15–57. MR 93h:17040 Zbl 0867.17022
- [Slodowy 1984] P. Slodowy, Singularitäten, Kac–Moody Lie-Algebren, assoziierte Gruppen und Verallgemeinerungen, Habilitationsschrift, Universität Bonn, 1984.
- [Slodowy 1986] P. Slodowy, "Beyond Kac–Moody algebras, and inside", pp. 361–371 in *Lie algebras and related topics* (Windsor, Canada, 1984), edited by D. J. Britten et al., CMS Conf. Proc. **5**, Amer. Math. Soc., Providence, RI, 1986. MR 832211 Zbl 0582.17011

Received April 22, 2012. Revised August 28, 2012.

SANDEEP BHARGAVA MATHEMATICS DEPARTMENT, SCHOOL OF LIBERAL ARTS AND SCIENCES HUMBER INSTITUTE OF TECHNOLOGY AND ADVANCED LEARNING 205 HUMBER COLLEGE BOULEVARD TORONTO, ON M9W 5L7 CANADA sandeep.bhargava@humber.ca

YUN GAO DEPARTMENT OF MATHEMATICS AND STATISTICS YORK UNIVERSITY 4700 KEELE STREET TORONTO, ON M3J 1P3 CANADA ygao@yorku.ca

STABLE FLAGS, TRIVIALIZATIONS AND REGULAR CONNECTIONS

ELIE COMPOINT AND EDUARDO COREL

We study stalkwise modifications of a holomorphic vector bundle endowed with a meromorphic connection on a compact Riemann surface. We introduce the notion of Birkhoff–Grothendieck trivialization, in the case of the Riemann sphere, and show that its computation corresponds to shortest paths in some local affine Bruhat–Tits building. We use this to compute how the type of a bundle changes under stalk modifications, and give several corresponding algorithmic procedures. We finally deduce from these results some applications to the Riemann–Hilbert problem.

Introduction

The motivation of this article originates in the Riemann–Hilbert problem on a compact Riemann surface, and the present work follows it as a guideline. The results presented herein are however not directly related to this problem. The reader who is exclusively interested in new advances on the Riemann–Hilbert problem will nevertheless find a couple of improvements on already known results. The real interest of the present paper in the eyes of its authors consists in translating this classical problem in a "new" setting (the setting of Bruhat–Tits buildings). In this new context, the Riemann–Hilbert problem reveals new geometric objects (such as *Birkhoff–Grothendieck trivializations*), whose study appears to be interesting by itself, and seems also promising for the original problem.

The Riemann–Hilbert problem (RHP) has a long and distinguished history, not even devoid of suspense, for it has been solved several times, using different tools, in a seemingly complete and positive way. It was finally A. A. Bolibrukh, in a celebrated series of papers at the beginning of the 1990s, who clarified the situation, by rigorously defining (and exhibiting a counterexample to) the strongest version of the RHP, thereby showing that people before him had either committed a mistake, or solved in reality a weaker problem.

MSC2010: primary 51N30, 34M03; secondary 34M50.

Keywords: meromorphic connection, vector bundle, Birkhoff–Grothendieck theorem, Bruhat–Tits building, Riemann–Hilbert problem.

The modern approach to the RHP was initiated by H. Röhrl in the 1950s, who used the theory of vector bundles in a way that has been conserved since. First, one constructs a vector bundle \mathscr{C} outside the singular points, whose cocycle mimics the monodromy. We call this the *topological* RH problem, since the monodromy is so much encoded in the topology of the constructed bundle, that construction of the required connection becomes essentially trivial. The second step consists of extending the bundle (and the connection) to the singular points by means of a *local solution* to the inverse monodromy problem. It has been exposed in great generality in [Deligne 1970] how to extend a holomorphic vector bundle \mathscr{C} , defined over the complement of a divisor \mathfrak{D} and endowed with a holomorphic connection ∇ having a prescribed monodromy about \mathfrak{D} , into a logarithmic connection ($\overline{\mathscr{C}}, \overline{\nabla}$) with singularities on the divisor.

In this way, we get all logarithmic extensions of \mathscr{C} with *nonresonant* residue (the *Deligne lattices*). These two steps are sufficient to solve positively the weak Riemann–Hilbert problem (i.e., with regular singularities). Note, however, that in this second level, two different types of problems have been mixed. The connection constructed is essentially unique up to *meromorphic* equivalence whereas the holomorphic vector bundle setting already introduces much finer *holomorphic* equivalence problems. This fact can contribute to explain some of the confusions that have surrounded the precise formulation of the RHP.

The *strong Riemann–Hilbert problem* asks for a logarithmic bundle (with the prescribed monodromy) which is moreover *trivial*. So, to solve the Riemann–Hilbert problem in this way, one must modify the constructed Deligne bundle, over the support of the singular divisor exclusively (to keep the singular set invariant), while conserving its logarithmic character, until a trivial bundle is eventually found. Until Bolibrukh's celebrated counterexample [1990], it was widely acknowledged that this was possible, and it is indeed so in several "generic" instances, although some mistakes in the seemingly general solution by Plemelj had already been pointed out (e.g., in [Treibich Kohn 1983]).

The counterexample found by Bolibrukh to the strong Riemann–Hilbert problem requires the knowledge of *all* the logarithmic extensions of a regular connection, in order to prove that none is trivial. Despite the production of both counterexamples and sufficient conditions for a positive answer, no general necessary and sufficient conditions for the solubility of the strong Riemann–Hilbert problem have been given in terms of the monodromy representation only, except for the remarkable case of an *irreducible* representation [Bolibrukh 1990; Kostov 1992].

As already stated, the strong Riemann–Hilbert problem admits a solution if and only if the stalks of the Deligne bundle over the singular set can be replaced by logarithmic lattices in such a way that the resulting bundle is trivial. To tackle this problem, one should be able to

- (a) determine the set of all logarithmic lattices above a given point, and
- (b) get a criterion for the triviality of the modified bundle.

Problem (a) can be considered classical, in the sense that it has been repeatedly solved, under different guises, going back as far as [Gantmacher 1959], and including [Levelt 1961; Babbitt and Varadarajan 1983; Bolibrukh 1990; Sabbah 2002]. In this paper, we give a complete description of the logarithmic lattices in terms of flags stabilized under the action of the residue of the connection; in a way, our contribution is to make the solution to problem (a) given in [Sabbah 2002] completely explicit and effective. We also give a partial answer to problem (b). In the case of $\mathbb{P}^1(\mathbb{C})$, the type of a vector bundle gives such a triviality criterion. In our approach, starting with the Deligne bundle \mathfrak{D}^{mod} . The question is then to compute the type of the modified bundle \mathfrak{D}^{mod} . Generalizing a result by Gabber and Sabbah (Proposition 28), we show how to determine the type of the Deligne bundle. In a second step, we show that this problem in turn is reduced to the well-known problem of *connection matrices*.

With these problems in mind, we introduce and study in this paper the notion of *Birkhoff–Grothendieck trivialization* of a bundle \mathcal{E} , which is a pointwise modification of \mathcal{E} such that

- (i) the resulting bundle \mathcal{F} is holomorphically trivial, and
- (ii) the relative elementary divisors of the stalks give the *type* of the bundle \mathcal{E} .

The paper is organized as follows. In a first section, we define the category in which we will work, and what we precisely mean by "modifying a bundle over one or several points". In a second part, we describe the geometry on the local lattices involved. We describe this geometry in terms of the affine Bruhat–Tits building of SL_n .

The third part contains the main results of the paper. We use the previously introduced setting to give an effective method to compute how the type of an arbitrary bundle E is modified under certain pointwise modifications. This algorithm can also be applied to compute the type of the bundle E. This third section concludes with a generalization of an essential result originally due to Bolibrukh, the *permutation lemma*, for which we provide an interesting geometric interpretation. This result allows us to give a quite complete insight into the structure of Birkhoff–Grothendieck trivializations, which we sum up as follows.

Theorem 1. Let Λ^0 be the set of pointwise modifications over $x \in \mathbb{P}^1(\mathbb{C})$ of a bundle \mathscr{C} that give a trivial bundle. Then $M \in \Lambda^0$ is a Birkhoff–Grothendieck trivialization of \mathscr{C} if and only if M realizes the minimum of the canonical metric between \mathscr{C}_x and Λ^0 in the local affine Bruhat–Tits building of SL_n at x.

The fourth section gives the complete description of the set of logarithmic lattices in terms of flags which are stable under the action of the residue of the connection on the Deligne bundle¹.

In the last part, we use these tools for study of the Riemann–Hilbert problem. After recalling the construction of the classical Röhrl–Deligne bundle, we give a very short proof of Plemelj's theorem on the Riemann–Hilbert solubility. This well-known result becomes an immediate consequence of the geometrical interpretations of the permutation lemma and the set of logarithmic lattices. We describe all trivializations of the Deligne bundle over an arbitrary point, and we give a concise proof of the Bolibrukh–Kostov theorem on the solvability of Riemann–Hilbert in the irreducible case. Finally, we give algorithmically effective procedures that allow to search the space of weak solutions.

1. Holomorphic vector bundles

Let X be a compact Riemann surface and let $\pi : E \to X$ be a holomorphic vector bundle of rank n. The sheaf \mathscr{C} of holomorphic sections of E is a locally free sheaf of \mathbb{O}_X -modules of the same rank n, where \mathbb{O}_X denotes as usual the sheaf of holomorphic functions on X. There is a well-known equivalence between these two categories. However, this equivalence does not hold for subobjects of the same rank. Therefore we will privilege the sheaf-theoretical formulation.

Meromorphic connections. Let $\mathfrak{D} = \sum_{i=1}^{p} m_i x_i$ be a positive divisor on X. Let $\mathbb{O}_{\mathfrak{D}}$ be the sheaf of meromorphic functions on X having pole orders bounded by \mathfrak{D} , and $|\mathfrak{D}| = \{x_1, \ldots, x_p\}$ be the *support* of \mathfrak{D} . For any finite set $S = \{y_1, \ldots, y_t\}$, let $(S) = y_1 + \cdots + y_t$.

Let $\nabla : \mathscr{C} \to \mathscr{C} \otimes_{\mathbb{O}_X} \Omega^1_X(\mathfrak{D})$ be a meromorphic connection with singular divisor \mathfrak{D} on a vector bundle \mathscr{C} of rank n. Sometimes for simplicity we'll just say "connection" for the pair (\mathscr{C}, ∇) . If we assume that \mathfrak{D} is the smallest possible (as we will do), then the *Poincaré rank* of ∇ at $x \in X$ is the integer $\mathfrak{p}_X(\nabla) = \max(0, m_X - 1)$. We will omit ∇ whenever possible. If $\mathfrak{p}_X = 0$, the sheaf \mathscr{C} is said to be *logarithmic at* x. Let $\mathscr{G} = |\mathfrak{D}|$ be the *singular*, and $\mathscr{G}_{\log} = \{x \in \mathscr{G} \mid \mathfrak{p}_X = 0\}$ the *logarithmic singular sets* of ∇ . If $\mathscr{G}_{\log} \neq \emptyset$, then one can define the residue map Res $\nabla \in \operatorname{End}(\mathscr{C}/\mathscr{C}_{-(\mathscr{G}_{\log})})$. We will specify in parentheses the bundle if necessary.

The meromorphic bundle. Let \mathcal{M}_X be the sheaf of meromorphic functions on X and let $\mathcal{V} = \mathscr{C} \otimes_{\mathcal{O}_X} \mathcal{M}_X$ be the sheaf of meromorphic sections of E. A meromorphic connection ∇ on \mathscr{C} has a canonical extension to \mathcal{V} . Since the sheaf \mathscr{C} can be

¹Independently, P. Boalch [2011] has taken a similar view on local logarithmic lattices, in terms of stable filtrations and Bruhat–Tits buildings, albeit on the more general setting of a complex reductive group. Restricting to GL_n and SL_n enables us to give however more explicit formulæ; see Sections 3.2 and 3.3.

embedded into \mathcal{V} , we consider from now on the set

$$\mathcal{H} = \left\{ \mathcal{F} \subset \mathcal{V} \mid \mathcal{F} \stackrel{\text{loc.}}{\simeq} \mathbb{O}_X^n \right\}$$

1 - -

of *holomorphic vector bundles* of \mathcal{V} . Each such bundle $\mathcal{F} \in \mathcal{H}$ is automatically endowed with a meromorphic connection induced by ∇ . For simplicity, we won't make any notational difference between all these connections.

We say that $\mathcal{F} \in \mathcal{H}$ is *trivial* if $\mathcal{F} \simeq \mathbb{O}_X^n$, or, equivalently, if \mathcal{F} is generated by its global sections, and *quasitrivial* if there exists a line bundle \mathcal{L} such that $\mathcal{F} \otimes_{\mathbb{O}_X} \mathcal{L}$ is trivial. Let $\mathcal{H}_0 \subset \mathcal{H}$ be the subset of trivial holomorphic bundles in \mathcal{V} . The following result is well known (e.g., [Sabbah 2002, p. 134]).

Lemma 2. Let $\mathcal{F} \in \mathcal{H}_0$ be a trivial holomorphic vector bundle in (\mathcal{V}, ∇) . The space $\mathcal{F}(X) = \Gamma(X, \mathcal{F})$ of global sections is a \mathbb{C} -vector space of dimension n. For any logarithmic singularity $s \in \mathcal{G}_{\log}(\mathcal{F})$, the residue $\operatorname{Res}_s^{\mathcal{F}} \nabla$ induces a well-defined endomorphism $\psi_s \in \operatorname{End}_{\mathbb{C}}(\mathcal{F}(X))$.

Stalks and lattices. For $\mathcal{F} \in \mathcal{H}$, the stalk \mathcal{F}_x at any $x \in X$ is a free $(\mathbb{O}_X)_x$ -submodule of rank *n* (or *lattice*) of \mathcal{V}_x , which is a vector space of dimension *n* over $(\mathcal{M}_X)_x$. Let $\mathbb{O}_x = (\widehat{\mathbb{O}_X})_x$ be the formal completion of $(\mathbb{O}_X)_x$, and $K_x = \operatorname{Frac}(\mathbb{O}_x)$ its field of fractions. Similarly, the formal completion $\widehat{\mathcal{F}_x} = \mathcal{F}_x \otimes_{(\mathbb{O}_X)_x} \mathbb{O}_x$ is a lattice in the vector space $V_x = \mathcal{V}_x \otimes_{(\mathcal{M}_X)_x} K_x$. This operation is harmless, as $\Lambda \mapsto \widehat{\Lambda}$ is a bijection between the sets of lattices in \mathcal{V}_x and V_x (cf. [Malgrange 1996]). We define an equivalence relation \sim_x on \mathcal{H} as

$$\mathscr{F} \sim_x \widetilde{\mathscr{F}}$$
 if and only if $\mathscr{F}_{|X \setminus \{x\}} = \widetilde{\mathscr{F}}_{|X \setminus \{x\}}$.

For simplicity, we will drop the index x as soon as no ambiguity can arise. Let Λ_x be the set of lattices in V_x . Any coset $[\mathcal{F}]$ of \mathcal{H}/\sim_x can be identified with the set Λ_x , by identifying $\mathcal{F}' \in [\mathcal{F}]$ with its formal stalk $\widehat{\mathcal{F}'_x} \in \Lambda_x$ at x. Since X is compact, two vector bundles $\mathcal{E}, \mathcal{F} \in \mathcal{H}$ have equal stalks outside a finite set.

Lemma 3. Let $\mathcal{C} \in \mathcal{H}$ be a holomorphic vector bundle. For any family of formal lattices $M_x \in \Lambda_x$ for x in a discrete set \mathcal{G} , there exists a unique vector bundle $\mathcal{C}^M \in \mathcal{H}$ such that

$$(\mathscr{E}^{M})_{x} = \begin{cases} \mathscr{E}_{x} & \text{if } x \notin \mathcal{G}, \\ M_{x} \cap \mathscr{V}_{x} & \text{if } x \in \mathcal{G}. \end{cases}$$

Conversely, for any $\mathcal{F} \in \mathcal{H}$, there exists a discrete set \mathcal{S} and a family $(M_x \in \Lambda_x)_{x \in \mathcal{G}}$ of lattices such that $\mathcal{F} = \mathcal{E}^M$. If \mathcal{E} is endowed with a meromorphic connection ∇ , there is a canonical extension ∇^M of $\nabla_{|X \setminus \mathcal{G}}$ as a meromorphic connection on \mathcal{E}^M . In particular, a bundle $\mathcal{E} \in \mathcal{H}$ is completely determined by its coset $[\mathcal{E}] \in \mathcal{H} / \sim_x$ and the lattice $\widehat{\mathcal{E}_x} \in \Lambda_x$. The sheaf \mathcal{V} is trivial (as a meromorphic bundle), and the group G of (meromorphic) automorphisms of the space $\Gamma(X, \mathcal{V})$ is isomorphic to $\operatorname{GL}_n(\mathbb{C}(X))$. Let $\Lambda^0_x = \Lambda_x \cap \mathcal{H}_0$ be the set of trivial bundles in the coset Λ_x . The subgroup $G_x \subset G$ of automorphisms of $\Gamma(X, \mathcal{V})$ that leave Λ^0_x globally invariant is called the *group of monopole gauge transforms at x*. Each element of G_x sends a trivial sheaf \mathcal{F} to a trivial sheaf \mathcal{F} such that $\mathcal{F}_{|X\setminus\{x\}} = \mathcal{F}_{|X\setminus\{x\}}$. An element of G_x modifies at most the stalk \mathcal{F}_x .

2. Lattices and the affine building of SL_n

2.1. *Flags and filtrations.* Let *V* be a vector space over a field K of characteristic 0. Given a flag *F* of vector spaces $0 = F_0 \subset F_1 \subset \cdots \subset F_s = V$, with *length* |F| = s and *signature* $\sigma(F) = (n_1, \ldots, n_s)$ where $n_i = \dim_{\mathbb{K}}(F_i/F_{i-1})$, a map $u \in \operatorname{End}_{\mathbb{K}}(V)$ *stabilizes* the flag *F* if $u(F_i) \subset F_i$ for all $0 \leq i \leq |F|$. Let $\mathfrak{Fl}(V)$ denote the set of flags of *V*, and $\mathfrak{Fl}_u(V)$ the subset of flags that are stabilized by *u*. A flag *F'* is *transversal* to *F* if |F'| = |F| = s and $F'_i \oplus F_{s-i} = V$ for $1 \leq i \leq s$. In this case, the signature of *F'* is equal to $\sigma^- = (n_s, \ldots, n_1)$. The left action of the permutation group S_n on a sequence $a = (a_1, \ldots, a_n)$ is given by $\tau a = (a_{\tau(1)}, \ldots, a_{\tau(n)})$ for $\tau \in S_n$. Let s_i be the transposition (i, i + 1) exchanging *i* and i + 1. The set $S = \{s_1, \ldots, s_{n-1}\}$ makes (S_n, S) into a Coxeter group of type A_{n-1} .

Let (e) be a basis of V, and let $\sigma = (n_1, \dots, n_s)$ be a signature. Let $F^{\sigma}(e)$ be the flag with elements $F_i^{\sigma}(e) = \langle e_1, \dots, e_{\nu_i} \rangle$ where $\nu_i = n_1 + \dots + n_i$. The basis (e) is said to be *adapted* to F if any element of the flag is spanned by a subfamily of (e), strictly adapted if $\mathbf{F} = F^{\sigma}(e)$, and transversal to \mathbf{F} if $F_i \oplus \langle e_{\nu_i+1}, \ldots, e_n \rangle = V$ for all *i*. The *parabolic subgroup* stabilizing a flag F (for the action on a strictly adapted basis) is the subgroup $W_F = \langle s_i | i \in I \rangle$ of S_n generated by the generators corresponding to the missing dimensions $I = [n] \setminus \{\dim_{\mathbb{K}} F_i \mid 1 \leq i \leq n\}$. These properties depend in fact only on the K-vector subspaces spanned by the vectors of (e). The opposite flag $F^{-(e)}$ is then defined as the unique flag transversal to F for which (e) is adapted. This last notion does not even depend on the order in which the vectors of (e) are taken. A flag F' is transversal to F if and only if there is a basis (e) of V strictly adapted to F such that $F' = F^{-(e)}$. A K-frame is an unordered set $\Phi = \{L_1, \dots, L_n\}$ of one-dimensional K-vector subspaces of V such that $L_1 + \cdots + L_n = V$. The notions defined in the previous paragraph make sense for a frame (with a fixed order on the lines for some of them). The relative position $\rho(\mathbf{F}, \mathbf{F}') \in S_n$ of two flags \mathbf{F} and \mathbf{F}' is the² permutation $\tau \in S_n$ such that

²Strictly speaking, ρ is only unique when both flags are complete. Otherwise, we shall a bit imprecisely consider ρ either as its double coset in $W_{F'}\rho W_F$ modulo the parabolic subgroups attached to F and F', or to the unique minimal length representative of this coset (or possibly even to *any* such representative).

there exists a basis (e) strictly adapted to F for which $(\tau e) = (e_{\tau(1)}, \dots, e_{\tau(n)})$ is strictly adapted to F'.

Similarly, a *flag* S of length s in $[n] = \{1, ..., n\}$ is an increasing sequence $S : S_0 = \emptyset \subset S_1 \subset \cdots \subset S_s = [n]$ of subsets of [n], whose *signature* is the sequence $\sigma(S) = (n_1, ..., n_s)$ where $n_i = |S_i| - |S_{i-1}|$. For a given signature $\sigma = (n_1, ..., n_s)$, the *standard ascending flag* $S^{\nearrow}(\sigma)$ of signature σ is the flag composed of initial segments of [n] of lengths $n_1 + \cdots + n_i$ for $1 \le i \le s$.

Given a sequence $D = (d_1, \ldots, d_n) \in \mathbb{Z}^n$, let D^{\nearrow} be the sequence of the elements of D arranged in increasing order. Define the *ascendent flag* $S^{\nearrow}(D)$ of D as the sequence of subsets of indices corresponding to blocks of equal elements of D^{\nearrow} . Let also $(\mathbb{Z}^n)^{\nearrow}$ denote the set of nondecreasing integer sequences. Finally, let

(1)
$$D \cong D'$$
 if $D^{\nearrow} = (D')^{\nearrow}$.

We denote with a \searrow symbol all similarly defined *descending* quantities. Note that $D^{\searrow} = w_0(D^{\nearrow})$ where the permutation $w_0 = (n, n-1, ..., 1)$ is the largest element of S_n in the Bruhat order.

We further define $D + \ell = (d_1 + \ell, \dots, d_n + \ell)$, $D_0 = D - \min D \in \mathbb{N}^n$ and $D^0 = \max D - D \in \mathbb{N}^n$. Let $\operatorname{Tr} D = \sum_{i=1}^n d_i$ and $\Delta D = \max D - \min D$, and finally

$$||D|| = \sum_{j=1}^{n} \left(d_j - \frac{\operatorname{Tr} D}{n} \right)^2$$
 and $i(D) = \sum_{j=1}^{n} (\max D - d_j) \in \mathbb{N}.$

We list some useful and obvious properties in the following lemma.

Lemma 4. For $D \in \mathbb{Z}^n$ and $\ell \in \mathbb{Z}$, we have

- (i) $\Delta(D+\ell) = \Delta D$ and $i(D+\ell) = i(D)$,
- (ii) $\Delta(-D) = \Delta D$ and $i(-D) = i(D^0) = \sum_{j=1}^n (d_j \min D)$,
- (iii) $i(D) + i(D^0) = n\Delta D$.

An *F*-admissible sequence is an integer sequence whose ascending flag is equal to the standard ascendent flag of signature $\sigma(F) = (n_1, \ldots, n_s)$; in more concrete terms, an integer sequence

$$D = (\underbrace{d_1, \dots, d_1}_{n_1 \text{ times}}, \underbrace{d_2, \dots, d_2}_{n_2 \text{ times}}, \dots, \underbrace{d_s, \dots, d_s}_{n_s \text{ times}}) \text{ with } d_1 < \dots < d_s.$$

Let $\mathbb{Z}^n(F)$ be the set of integer *F*-admissible sequences, and let

$$\Xi(V) = \{ (F, D) \mid F \in \mathfrak{Fl}(V) \text{ and } D \in \mathbb{Z}^n(F) \}$$

be the set of F-filtrations of V.

2.2. Lattices. In the remainder of the section, we fix a point $x \in X$ and a coset $\Lambda \in H/\sim_x$. We drop the index x for simplicity. The field $K = (\mathcal{M}_X)_x$ is local, and endowed with the discrete valuation $v = \operatorname{ord}_x$, whose valuation ring and maximal ideal we denote by \mathbb{O} and \mathfrak{m} . Let V be the K-vector space $\mathcal{V}_X \otimes_{(\mathcal{M}_X)_x} K$ of dimension n.

Let $\mathcal{L}(u)$ denote the free \mathbb{O} -module spanned by a family (u) of vectors in V. An \mathbb{O} -module $M \in V$ is a lattice if there exists a K-basis (e) of V such that $M = \mathcal{L}(e)$. Let

$$v_{\Lambda}(x) = \max\{k \in \mathbb{Z} \mid x \in \mathfrak{m}^{k}\Lambda\}$$

be the natural valuation of V induced by Λ . For any lattices $M \subset \Lambda$ in V, we define the *interval* $[M, \Lambda]$ as

$$[M, \Lambda] = \{ N \in \Lambda \mid M \subset N \subset \Lambda \}.$$

Let $\pi_{\Lambda} : \Lambda \to \overline{\Lambda} = \Lambda/\mathfrak{m}\Lambda \simeq \mathbb{C}^n$ denote the canonical surjection on the quotient module.

Elementary divisors. Let z be a uniformizing parameter of K. For any two lattices Λ and M in V, there exists a unique increasing sequence of integers $d_1 \leq \cdots \leq d_n$ (the elementary divisors of M in Λ) and an \mathbb{O} -basis (e_1, \ldots, e_n) of Λ such that $(z^{d_1}e_1, \ldots, z^{d_n}e_n)$ is a basis of M. Such a basis (e) is called a *Smith basis of* Λ for M. We will write them $d_i^{\Lambda}(M)$ if we want to specify the respective lattices, and we put

$$\mathbf{ED}_{\Lambda}(M) = \left(d_1^{\Lambda}(M), \dots, d_n^{\Lambda}(M) \right).$$

Note that $d_1^{\Lambda}(M) = v_{\Lambda}(M)$ and $\mathbf{ED}_{\Lambda}(z^k M) = \mathbf{ED}_{\Lambda}(M) + k$. Let also

$$M_{\Lambda} = z^{-\nu_{\Lambda}(M)}M$$
 and $\mu_{\Lambda}(M) = \dim_{\mathbb{C}} M_{\Lambda}/\mathfrak{m}\Lambda$.

If $P \in GL_n(K)$ is a basis change from Λ to M, the sequence $ED_{\Lambda}(M) = (d_1, \ldots, d_n)$ can be computed in the following manner. For a subset $I \subset [n]$ of cardinality |I|, let $S_I = \sum_{i \in I} s_i$, and for |I| = |J| = k, let $P_{I,J}$ denote the (I, J)-submatrix of P. The sequence (d_1, \ldots, d_n) satisfies

(2)
$$d_k = e_k - e_{k-1}$$
 where $e_k = \min_{|I|=|J|=k} \{v(\det P_{I,J})\}$ and $e_0 = 0$.

It is convenient to be a bit more lax in the definition, and allow the elementary divisors to appear in another order. To avoid ambiguities, we will specify that (e) is an *ascending Smith basis* of Λ for M, if the vectors in (e) are ordered according to $\mathbf{ED}_{\Lambda}(M)^{\nearrow}$.

We say that a matrix *P* is *D*-parabolic if $P_{ij} \neq 0 \Rightarrow d_i \leq d_j$ for $1 \leq i, j \leq n$. For any commutative ring *R*, put $G_D(R)$ for the group of *D*-parabolic matrices of
$GL_n(R)$. The subgroup of $GL_n(\mathbb{O})$ that acts on the set of Smith bases of Λ for M is the subgroup of \mathcal{G}_D -parabolic matrices

$$\mathscr{G}_D = \{ P \in \operatorname{GL}_n(\mathbb{O}) \mid v(P_{ij}) \ge d_i - d_j \}.$$

Being \mathcal{G}_D -parabolic is stronger than *D*-parahoric, which only stabilizes the induced *D*-flag in $\Lambda/\mathfrak{m}\Lambda$, and weaker than *D*-parabolic, which stabilizes a *K*-flag of signature $\sigma(D)$. For any sequence $D = (d_1, \ldots, d_n)$, let z^D be the diagonal matrix diag $(z^{d_1}, \cdots, z^{d_n})$. We will frequently use the following type of diagram:

$$\begin{array}{c} \Lambda:(e) \xrightarrow{z^{D}} M:(\varepsilon) \\ \downarrow^{P} & \downarrow^{\tilde{P}} \\ \Lambda:(e') \xrightarrow{z^{D}} M:(\varepsilon') \end{array}$$

which means that (e) and (e') are two Smith bases of Λ for M. Note that in this case $P \in \operatorname{GL}_n(\mathbb{O})$ is \mathcal{G}_D -parabolic. Since $G_D(\mathbb{C}) \subset \mathcal{G}_D$ holds, the obvious factorization $P = P_0 P_I$ of a lattice gauge $P \in \operatorname{GL}_n(\mathbb{O})$ into a constant term $P_0 = P(0) \in \operatorname{GL}_n(\mathbb{C})$ and a term $P_I = I + zU$ with $U \in \mathfrak{gl}_n(\mathbb{O})$ tangent to I satisfies the property

$$P \in \mathcal{G}_D \iff (P \in G_D(\mathbb{C}) \text{ and } P_I \in \mathcal{G}_D).$$

We can therefore usually assume that P is tangent to I. Note that this also holds for a right factorization $P = P'_I P_0$.

Sometimes, we will find it more convenient to consider the elementary divisors with their multiplicities. In this case, we will put d_1, \ldots, d_s for the distinct elementary divisors of M in Λ and let n_i be their respective multiplicities. The set $[n] = \{1, \ldots, n\}$ of indices of ordinary (*simple*) elementary divisors is partitioned into the subsets I_i corresponding to a single value of the elementary divisors:

$$I_j = \{1 \leq \ell \leq n \mid d_\ell = \boldsymbol{d}_j\} \text{ for } 1 \leq j \leq s.$$

2.3. *Relative flag of a lattice.* Any lattice M induces a natural flag in $\overline{\Lambda} = \Lambda/\mathfrak{m}\Lambda$. For any $k \in \mathbb{Z}$, let

$$M_k = (\mathfrak{m}^{-k} M \cap \Lambda) + \mathfrak{m} \Lambda \in [\mathfrak{m} \Lambda, \Lambda].$$

Lemma 5. Let Λ , M be lattices in V. Let d_1, \ldots, d_s be the distinct elementary divisors of M in Λ . The flag in $\overline{\Lambda}$

$$F^{\Lambda}(M): 0 \subset \pi_{\Lambda}(M_{d_1}) \subset \cdots \subset \pi_{\Lambda}(M_{d_s}) \subset \overline{\Lambda}$$

has signature $\sigma(\mathbf{ED}_{\Lambda}(M))$.

Proof. Let (e) be a basis of elementary divisors of Λ for M, and $I = \{1 \le i \le n | d_i \le k\}$. Then M_k admits (u) as basis where $u_i = e_i$ if $i \in I$ and $u_i = ze_i$ if $i \notin I$.

The spaces M_k are thus embedded lattices, all belonging to the interval $[m\Lambda, \Lambda]$, so they take at most n + 1 different values. Their images $\overline{M}_k = \pi_{\Lambda}(M_k)$ in the quotient space $\overline{\Lambda}$ form a flag $F^{\Lambda}(M)$, and it is clear that $\overline{M}_{k-1} \subsetneq \overline{M}_k$ if and only if k is an elementary divisor of M. If d_1, \ldots, d_s are the distinct elementary divisors of M in Λ , with multiplicities n_i , the subset of indices corresponding to d_j can be written as

$$I_j = [n_1 + \dots + n_j, n_1 + \dots + n_{j+1} - 1].$$

The lattices M_k , M_ℓ coincide if and only if there exists *i* such that $d_i \leq k, \ell < d_{i+1}$ (with the conventions $d_0 = -\infty$ and $d_{s+1} = +\infty$). Therefore the flag $F^{\Lambda}(M)$ has exactly length *s*, and its signature is equal to the sequence (n_1, \ldots, n_s) .

The components of $F^{\Lambda}(M)$ can be indexed either as \overline{M}_{d_i} , by the value of the elementary divisor d_i to which it is attached (if known), or as \overline{M}_i , by its index in the flag (here *i*). In this latter case, we will also use the notation $F_i^{\Lambda}(M)$. It will hopefully be always clear which convention we are using.

Lemma 6. Let Λ , M, N be lattices in V. Let $D = \mathbf{ED}_{\Lambda}(M)^{\nearrow}$ and $D' = \mathbf{ED}_{\Lambda}(N)^{\nearrow}$. If either

- (i) there exists a common Smith basis for Λ , M and N, or
- (ii) the flags $F^{\Lambda}(N)$ and $F^{\Lambda}(M)$ are transverse,

then we have

$$\mathbf{ED}_{\boldsymbol{M}}(N) \cong D' - \sigma D,$$

where $\sigma = \rho(\mathbf{F}^{\Lambda}(N), \mathbf{F}^{\Lambda}(M))$ is the relative position of the induced flags in $\Lambda/\mathfrak{m}\Lambda$.

A similar formula holds for the descending sequences D^{\searrow} , but with $w_0 \sigma w_0$ instead of σ .

Proof. We summarize the setting by means of the following scheme:

$$\begin{array}{c} \Lambda : (e) \xrightarrow{z^{S^{\prime}}} M : (\varepsilon) \\ \downarrow^{P} & \downarrow^{\tilde{P}} \\ \Lambda : (e') \xrightarrow{z^{T^{\prime}}} M : (\varepsilon') \end{array}$$

If there is an apartment containing Λ , M, N, then one can assume that P is a permutation matrix, namely $\Pi_{\sigma^{-1}}$. In the second case, it is possible to choose $P = \Pi_{w_0}(I+zU)$ with $U \in \mathfrak{gl}(\mathbb{O})$. Let Q = I+zU. Then we have $\tilde{P}_{ij} = Q_{ij}z^{t_j-s_{n+1-i}}$. According to formula (2), we have $d_1 = \min_{1 \le i,j \le n}(v(P_{ij}) + t_j - s_{n+1-i})$. The minimum of $t_j - s_{n+1-i}$ is attained for (i, j) = (1, 1), and by assumption we have $v(P_{11}) = 0$. Therefore, we get $d_1 = t_1 - s_n$. Let us prove by induction that $d_i = t_i - s_{n+1-i}$ for all *i*. Assume that $e_i = T_{\{1,\dots,i\}} - S_{\{n+1-i,\dots,n\}}$ for $i \le k-1$.

Formula (2) yields $e_k = \min_{|I|=|J|=k} \{v(\det P_{I,J}) + T_J - (w_0 S)_I\}$. The minimum for $T_J - (w_0 S)_I$ is attained for $I = J = \{1, \dots, k\}$, while by assumption the principal minor $P_{I,I}$ has valuation 0. Hence $e_k = T_{[k]} - (w_0 S)_{[k]}$ and thus $d_k = e_k - e_{k-1} = t_k - s_{n+1-k}$ for all k.

For a flag $F = (F_i)_{1 \le i \le s}$ of *K*-vector spaces, and a lattice Λ , let $F_{\Lambda} = (F \cap \Lambda)/\mathfrak{m}\Lambda$ be the flag $((F_i \cap \Lambda)/\mathfrak{m}\Lambda)_{1 \le i \le s}$ induced in $\Lambda/\mathfrak{m}\Lambda$. The following result is easily established.

Lemma 7. Let F and F' be two K-flags of V. For any basis (e), adapted for F and F', one has

$$\rho(\mathbf{F}, \mathbf{F}') = \rho(\mathbf{F}_{\Lambda}, \mathbf{F}'_{\Lambda}), \text{ where } \Lambda = \mathcal{L}(e).$$

2.4. *The affine building of* SL(V). For this section, good references are [Garrett 1997; Ronan 1989], and especially [Abramenko and Brown 2008]. The affine building B_n naturally attached to SL(V) is the following (n - 1)-dimensional simplicial complex. Two lattices Λ and M are *homothetic* if there exists $\alpha \in K^*$ such that $M = \alpha \Lambda$. Let [Λ] be the homothety class of the lattice Λ in V. Two classes L and L' are *adjacent* if and only if there exist representatives Λ of L and M of L' such that $m\Lambda \subset M \subset \Lambda$. Consider a graph whose vertices are the homothety classes of lattices in V, and edges connect all adjacent vertices. The affine building B_n is the *flag simplicial complex* associated with this graph, or in other terms, its *clique complex*. A simplex A is a set $\{L_1, \ldots, L_k\}$ of mutually adjacent vertices, and the *face relation* $B \leq A$ is defined by the inclusion $B \subset A$ of these sets.

Lemma 8. Let L, L', L'' be vertices in B_n . If L' and L'' are adjacent, then for any representatives $\Lambda \in L$, $M \in L'$ and $N \in L''$, the flags $F^{\Lambda}(M)$ and $F^{\Lambda}(N)$ are compatible (in the sense that their components are all pairwise comparable). In particular, a maximal simplex C induces a complete flag in $\Lambda/\mathfrak{m}\Lambda$.

Proof. One can find representatives $\Lambda \in L$, $M \in L'$ and $N \in L''$, and a basis (e) of Λ such that $M = \mathcal{L}(z^{D}e)$ and $N = \mathcal{L}(z^{D'}e)$, and that $D' - D \in \{0, 1\}^n$. Indeed, let (u) be a basis of $M/\mathfrak{m}M$ which is adapted to both flags $F^M(\Lambda)$ and $F^M(N)$. There exists a lifting (ε) of (u) in M which is a Smith basis of M for Λ . Since M and N are adjacent, any lifting of (u) is a Smith basis for N. Then $(e) = (z^{ED_M(\Lambda)}\varepsilon)$ is a common Smith basis of Λ for M and N. By definition of adjacency, one can ensure that the representatives satisfy $\mathfrak{m}M \subset N \subset M$; that is, $D' - D \in \{0, 1\}^n$. Let now $I_k = \{1 \leq i \leq n \mid d_i \leq k\}$ and $I'_k = \{1 \leq i \leq n \mid d_i \leq k\}$. Then we have $I_{k-1} \subset I'_k \subset I_k$ for all i. Since $F^{\Lambda}(M)$ is spanned by the subfamilies $(\overline{e_i})_{i \in I_k}$ for $k \in \mathbb{Z}$, the claim follows.

Simplices and chambers. A maximal simplex, or chamber in B_n , is an *n*-chain C of vertices L_0, \ldots, L_{n-1} with representatives Λ_i for $0 \le i \le n-1$ satisfying

$$\mathfrak{m}\Lambda_0\subset\Lambda_1\subset\cdots\subset\Lambda_{n-1}\subset\Lambda_0.$$

A basis (e) is a fundamental basis for C at L_i if

$$L_{(i+j) \mod n} = [\mathscr{L}(e_1, \dots, e_j, ze_{j+1}, \dots, ze_n)] \quad \text{for } 0 \le j \le n-1.$$

Lemma 9. The set of chambers which contain a given vertex L is in bijection with the set of complete flags in $\Lambda/\mathfrak{m}\Lambda$ with $\Lambda \in L$. A basis (e) is fundamental for Cat L if and only if its image in $\Lambda/\mathfrak{m}\Lambda$ is strictly adapted to the flag $F^{\Lambda}(C)$.

A (partial) flag F in $\Lambda/\mathfrak{m}\Lambda$ can be lifted (by $\pi_{[\Lambda]}^{-1}$) to a uniquely defined simplex in B_n containing the vertex $L = [\Lambda]$.

Definition 10. The graph-theoretic distance, canonical metric and index on B_n are defined, respectively, by

$$d(L, L') = \Delta \left(\mathbf{E} \mathbf{D}_{\Lambda}(\Lambda') \right), \quad d(L, L') = \left\| \mathbf{E} \mathbf{D}_{\Lambda}(\Lambda') \right\|, \quad [L:L'] = i \left(\mathbf{E} \mathbf{D}_{\Lambda}(\Lambda')^{\mathbf{0}} \right)$$

for any representatives Λ , Λ' of L, L'.

Note that $d(L, L') = -v_{\Lambda}(\Lambda') - v_{\Lambda'}(\Lambda)$ also holds, and that these three maps are indeed symmetric. The *d* metric makes the geometric realization of B_n into a CAT(0) space [Abramenko and Brown 2008, Theorem 11.16, p. 555].

Apartments. Let $\Phi = \{d_1, \dots, d_n\}$ be a *K*-frame of *V*. The set

$$\mathcal{A}_{\Phi} = \{\Lambda = \ell_1 + \dots + \ell_n \mid \ell_i \text{ is a lattice in } d_i\}$$

of lattices spanned over multiples of the vectors in Φ induces a simplicial subcomplex in the affine building B_n called the *apartment* spanned by Φ . For any lattice $\Lambda \in \Lambda$, a Λ -basis of the apartment \mathcal{A}_{Φ} is a collection $(e) = (e_1, \ldots, e_n)$ of vectors such that e_i spans d_i and $v_{\Lambda}(e_i) = 0$. Such a family is unique up to permutation and to multiplication of each e_i by a scalar $\lambda_i \in \mathbb{O}^*$. The lattice is an element of the apartment \mathcal{A}_{Φ} if and only if the family $(e) = (e_1, \ldots, e_n)$ is actually a basis of the lattice Λ . Equivalently, and without reference to a basis, this means that

$$\Lambda = \bigoplus_{i=1}^n \Lambda \cap d_i.$$

In the general case, the lattice $\Lambda_{\Phi} = \bigoplus_{i=1}^{n} \Lambda \cap d_i$ is the *largest sublattice* of Λ in the apartment \mathcal{A}_{Φ} . The homothety class $L_{\Phi} = [\Lambda_{\Phi}]$ is therefore the closest point projection of $L = [\Lambda]$ on \mathcal{A}_{Φ} , and the map

$$\rho_{\Phi}: B_n \to \mathcal{A}_{\Phi}$$
$$\Lambda \mapsto \Lambda_{\Phi}$$

is the *retract* onto \mathcal{A}_{Φ} .

Galleries and type. Two chambers *C* and *C'* are *adjacent* if they share n-1 vertices in common (a so-called *panel*). A *gallery* Γ from *C* to *C'* of length $\ell(\Gamma) = \ell$ is a sequence of chambers $\Gamma = (C_0 = C, C_1, \dots, C_{\ell} = C')$ such that C_i is adjacent to C_{i+1} . The *gallery distance* is defined as

$$\ell(C, D) = \min{\{\ell(\Gamma) \mid \Gamma \text{ is a gallery from } C \text{ to } D\}}$$

The vertices in the building B_n can be labelled by $\mathbb{Z}/n\mathbb{Z}$, which we assume given from now on. Let $S = \{s_0, \ldots, s_{n-1}\}$. The affine Weyl group W of the building then has a presentation

$$W = \langle S | s_i^2 = 1 \text{ and } (s_i s_{i+1})^3 = 1 \rangle$$

where the indices are understood modulo *n*. Every chamber has exactly one vertex labelled *i*, that we denote $M_i(C)$, for every $i \in \mathbb{Z}/n\mathbb{Z}$. Two chambers are called *i*-adjacent when their common panel precisely does not contain those vertices labelled *i*. Two *i*-adjacent chambers *C* and *C'* are then said to have *W*-distance $\delta(C, C') = s_i$. If $\Gamma = (C, C_1, \ldots, C_{t-1}, C')$ is any gallery from *C* to *C'* lying inside an apartment, the *W*-distance $\delta(C, C')$ is defined as the product

$$\delta(C, C') = \delta(C, C_1) \cdots \delta(C_{t-1}, C').$$

There is a bijection between minimal galleries between *C* and *C'* and minimal length decompositions of the *W*-distance $\delta(C, C') \in W$ into products of generators in *S*. Let \mathbb{Z}_0^n be the set of *n*-tuples of integers summing to 0. A labelling of the vertices by $\mathbb{Z}/n\mathbb{Z}$ induces an isomorphism

$$W \stackrel{\phi}{\simeq} S_n \rtimes \mathbb{Z}_0^n \quad \text{given by} \begin{cases} \phi(s_i) = ((i, i+1), 0) & \text{for } 1 \le i \le n, \\ \phi(s_0) = ((1, n), (1, 0, \dots, 0, -1)). \end{cases}$$

Let $\delta = \delta(C, C')$ be the *W*-distance between two chambers *C* and *C'*. Let $L = M_0(C)$ and $L' = M_0(C')$ be the respective unique vertices of type 0 of *C* and *C'*. Then $\phi(\delta) = (\sigma, K)$ is the unique couple such that there exists a fundamental basis (e) of *C* at *L* for which $(z^K e_{\sigma}) = (z^{k_1} e_{\sigma(1)}, \dots, z^{k_n} e_{\sigma(n)})$ is a fundamental basis of *C'* at *L'*.

Walls. Let (*e*) be a basis of *V*, and let \mathcal{A} be the apartment spanned by (*e*). In the basis (*e*), any lattice *L* in \mathcal{A} can be represented by a unique (up to an integer multiple of (1, ..., 1)) *n*-tuple $(x_1(L), ..., x_n(L)) \in \mathbb{Z}^n$. The set \mathcal{H} of *walls* of *V* for *W* is the set of hyperplanes

(3)
$$H_{i,j}^{(k)} = \{x \in \mathbb{R}^n \mid x_i - x_j = k\} \text{ for } 1 \le i < j \le n \text{ and } k \in \mathbb{Z}.$$

Define the corresponding half-spaces

$$H_{i,j}^{(k)+} = \{ x \in \mathbb{R}^n \mid x_i - x_j \ge k \},\$$

$$H_{i,j}^{(k)-} = \{ x \in \mathbb{R}^n \mid x_i - x_j \le k \}.$$

For a hyperplane $H \in \mathcal{H}$ and a simplex A, let

$$\sigma_H(A) = \begin{cases} + & \text{if } A \subset H^+ \backslash H^-, \\ - & \text{if } A \subset H^- \backslash H^+, \\ 0 & \text{if } A \subset H^+ \cap H^-. \end{cases}$$

A wall *H* separates the simplices *A* and *B* if $\sigma_H(A)\sigma_H(B) = -$ for the usual sign product rule. Any simplex *A* is completely defined by its sign sequence $(\sigma_H(A))_{H \in \mathcal{H}}$. However, there is only a finite number of them which are relevant (i.e., whose defining equations or inequalities are not redundant). If we define the fundamental chamber as $C_0 = (L_0, \ldots, L_{n-1})$, where $L_i = [\mathcal{L}(e_1, \ldots, e_i, ze_{i+1}, \ldots, ze_n)]$, the (essential) walls of C_0 are the hyperplanes $H_{i,(i \mod n)+1}^{(0)}$ for $1 \le i \le n$, and its sign sequence is

$$\sigma_{H_{i,(i \bmod n)+1}^{(k)}}(C_0) = \begin{cases} + & \text{for } 1 \leq i \leq n-1 \text{ and } k \leq 0 \text{ or } i = n \text{ and } k \geq 0, \\ - & \text{for } 1 \leq i \leq n-1 \text{ and } k \geq 0 \text{ or } i = n \text{ and } k \leq 0. \end{cases}$$

Definition 11. Let *A*, *B* be simplices in B_n . The gallery distance $\ell(A, B)$ is defined as $\min_{A \leq C, B \leq D} \ell(C, D)$.

Lemma 12. For two simplices A, B we have

$$\ell(A,B) = |\{H \in \mathcal{H} \mid \sigma_H(A)\sigma_H(B) = -\}|.$$

If $A = [\Lambda]$ and B = [M] are vertices, we have

$$\ell(A, B) = \sum_{1 \le i < j \le n} \max(0, d_j - d_i - 1), \quad \text{where } (d_1, \dots, d_n) = \mathbf{ED}_{\Lambda}(M)^{\nearrow}.$$

Proof. The first assertion is known (see [Abramenko and Brown 2008, p. 32]). Take a basis (e) of V where A has coordinates a and B coordinates b. The description (3) of walls shows that A and B are separated by $H_{i,j}^{(k)}$ if and only if $a_i - a_j \le k - 1$ and $b_i - b_j \ge k + 1$. Taking for (e) a Smith basis of Λ for M gives a = 0 and $b = (d_1, \ldots, d_n)$. The result is then a straightforward count.

2.5. Forms. For $L \in B_n$, let $\overline{L} = \Lambda/\mathfrak{m}\Lambda \simeq \mathbb{C}^n$ for $\Lambda \in L$. This definition is independent of the choice of Λ as there is a canonical isomorphism between $\Lambda/\mathfrak{m}\Lambda$ and $\mathfrak{m}^k \Lambda/\mathfrak{m}^{k+1}\Lambda$ for any k. For $L' \in B_n$, let N be the unique representative of L' such that $v_{\Lambda}(N) = 0$, and define

(4)
$$\pi_L(L') = (N + \mathfrak{m}\Lambda)/\mathfrak{m}\Lambda \in L.$$

The map π_L induces an isomorphism of simplicial complexes between the link lk(L) of $L = [\Lambda]$ and the set *E* of chains of linear subspaces of \overline{L} . A *form* in a lattice Λ is

296

a \mathbb{C} -vector subspace Y of Λ spanned by an \mathbb{O} -basis (*e*) of Λ . For any $x \in \Lambda$, there is a unique representative x_Y of the coset $x + \mathfrak{m}\Lambda$ in Y. This induces a well-defined isomorphism

$$\phi_Y: Y \xrightarrow{\simeq} \Lambda/\mathfrak{m}\Lambda.$$

We will say that $Y \subset \Lambda$ is a *Smith form* for *M* if there is a basis of *Y* which is a Smith basis of Λ for *M*.

Lemma 13. Let Λ be a lattice in V and Y be a form in Λ .

- (i) For any basis (e) of the lattice Λ, there exists a unique C-basis (e_Y) of the form Y whose image in Λ/mΛ coincides with the image of (e). We call (e_Y) the Y-basis of (e).
- (ii) Given a filtration (\mathbf{F}, D) in $\Xi(\Lambda/\mathfrak{m}\Lambda)$, there exists a unique lattice $M = \mathscr{L}_Y(\mathbf{F}, D)$ such that $\mathbf{F}^{\Lambda}(M) = \mathbf{F}$ and $\mathbf{ED}_{\Lambda}(M) = D$ that admits a Smith basis in Y.
- (iii) For any lattice M such that $d(\Lambda, M) = 1$, let (\overline{e}) a basis of $\Lambda/\mathfrak{m}\Lambda$ respecting $\pi_{\Lambda}(M) = M/\mathfrak{m}\Lambda$. Then the Y-basis (e_Y) is a Smith basis for M.

Proof. The basis (ε) obtained by putting $\varepsilon_i = (e_i)_Y = \phi_Y^{-1}(\pi_\Lambda(e_i))$ obviously satisfies the conditions of (i). For any \mathbb{C} -basis (e) of Y which respects the flag F, put $M = \bigoplus_{i=1}^n z^{d_i} e_i$. Let (\tilde{e}) be another basis of Y and $\tilde{M} = \bigoplus_{i=1}^n z^{d_i} \tilde{e}_i$. The matrix of the change of basis from $(z^D e)$ to $(z^D \tilde{e})$ is equal to $P = z^D C z^{-D}$, where $C \in \operatorname{GL}_n(\mathbb{C})$ is the matrix of the change of basis from (e) to (\tilde{e}). By definition of the parabolic subgroup P_F , one has $z^D C z^{-D} \in \operatorname{GL}_n(\mathbb{O}) \iff C \in P_F$; hence $M = \tilde{M}$ if and only if (e) and (\tilde{e}) both respect the flag F. Note that the gauge from the basis (e) to its Y-basis is always of the form $P = I + zU \in \operatorname{GL}_n(\mathbb{O})$. Let $W = \pi_\Lambda(M)$ and let

$$T = \begin{pmatrix} 0_r & 0\\ 0 & I_{n-r} \end{pmatrix}$$

be the elementary divisors of $M_{\Lambda} = z^{-\nu_{\Lambda}(M)}M$ with respect to Λ . Assume that (e) satisfies the assumptions of (iii). Then the Y-basis (e_Y) is obtained by a gauge $P = I + zU \in GL_n(\mathbb{O})$. Putting

$$U = \begin{pmatrix} U_{11} & U_{12} \\ U_{21} & U_{22} \end{pmatrix},$$

we have

j

since

$$\tilde{U} = z^{-T} (I + zU) z^{T} = \begin{pmatrix} I_{r} + zU_{11} & z^{2}U_{12} \\ U_{21} & I_{n-r} + zU_{22} \end{pmatrix}.$$

The basis (e_Y) is therefore indeed a Smith basis of Λ for M.

For any two forms Y and \tilde{Y} in Λ , the set of gauges between bases of Y and \tilde{Y} is an element of the double coset $\operatorname{GL}_n(\mathbb{C})\backslash\operatorname{GL}_n(\mathbb{C})$. Let z be a uniformizing parameter. With the convention that $\operatorname{deg}_z P = \infty$ if $P \in \operatorname{GL}_n(\mathbb{O})\backslash \mathfrak{gl}(\mathbb{C}[z])$, the following definition makes sense.

Definition 14. Let Y, \tilde{Y} be two forms in Λ . Define the *z*-distance as

$$\delta_z(Y, \widetilde{Y}) = \min(\deg_z P, \deg_z P^{-1}) \in \mathbb{N} \cup \{\infty\}$$

for any gauge P from a basis of Y to a basis of \tilde{Y} .

Lemma 15. If $t = d(\Lambda, M)$, then for any form Y of Λ , and any uniformizing parameter z, there exists a Smith form \tilde{Y} for M at a z-distance $\delta_z(Y, \tilde{Y}) \leq t - 1$.

Proof. There exists a Smith form Y' of Λ for M. Let $P = P_0 + P_1 z + \cdots \in \operatorname{GL}_n(\mathbb{O})$ be a gauge corresponding to a basis change from Y to Y'. Let $\overline{P} = P_0 + \cdots + P_k z^k$ for some $k \ge 0$, and let \widetilde{Y} be the form obtained by this gauge transformation, as explained in the following scheme:



We have $Q = z^{-D} P^{-1} \overline{P} z^D = (\widetilde{P}_{ij} z^{d_j - d_i})$, where $\widetilde{P} = P^{-1} \overline{P}$. By construction, we have $\widetilde{P} = P^{-1}(P - (P - \overline{P})) = I + z^{k+1}U$ with $U \in \mathfrak{gl}(\mathbb{O})$. As soon as $k \ge t-1$, we have $Q \in \operatorname{GL}_n(\mathbb{O})$; hence the form \widetilde{Y} is a Smith form for M. \Box

2.6. Shortest paths and elementary splittings. A shortest path is a path Γ in B_n such that for any vertices $L, L' \in \Gamma$, the length of the path between L and L' induced by Γ is equal to d(L, L'). For a path $\Gamma = (L_0, \ldots, L_t)$ in B_n , and a representative $\Lambda_0 \in L_0$, let the Λ_0 -normalized sequence of Γ be the unique sequence of lattices $\Lambda_i \in L_i$ such that $v_{\Lambda_0}(\Lambda_i) = 0$.

Lemma 16. Let Λ , $M \in \Lambda$. The Λ -normalized sequence $(\Lambda, \Lambda_1, \dots, \Lambda_t = M_\Lambda)$ of a shortest path from $[\Lambda]$ to [M] satisfies $\Lambda \supset \Lambda_1 \supset \dots \supset \Lambda_t \supset z^t \Lambda$.

Proof. Let $\Gamma = (\Lambda = \Lambda_0, \Lambda_1, \dots, \Lambda_{t-1}, M = \Lambda_t)$ be a path of minimal length *t*, and let us prove that $\Lambda \supset \Lambda_1 \supset \dots \supset \Lambda_t$ by induction on *t*. For t = 1, this is the very definition of adjacency in B_n . Assuming that the claim is established for any pair of

298

lattices at distance $\leq t-1$, we have $M \subset \Lambda \supset \Lambda_1 \supset \cdots \supset \Lambda_{t-1}$ for the normalized sequence of Γ . Since Λ_{t-1} and M are adjacent, there exists a unique $k \in \mathbb{Z}$ such that $z^k \Lambda_{t-1} \supset M \supset z^{k+1} \Lambda_{t-1}$. We know that $d(\Lambda, \Lambda_{t-1}) = t-1$, and hence we have $\Lambda \supset \Lambda_{t-1} \supset z^{t-1} \Lambda$; therefore we get

$$z^k \Lambda \supset z^k \Lambda_{t-1} \supset M \supset z^{k+1} \Lambda_{t-1} \supset z^{t+k} \Lambda.$$

If k > 0, then $v_{\Lambda}(M) \ge k > 0$, which was excluded by assumption. But if k < 0, then $d(\Lambda, M) < t$, which is also excluded. Thus we have k = 0, and the claim is proved.

The following result explains how to construct some shortest paths algebraically.

Proposition 17. Let $L, L' \in B_n$, and let t = d(L, L'). For $k \in \mathbb{N}$, let $L_k = [\Lambda' + \mathfrak{m}^k \Lambda]$, and $M_k = [\Lambda \cap \mathfrak{m}^{k-t} \Lambda']$, where $\Lambda \in L$ and $\Lambda' \in L'$ are such that $v_{\Lambda}(\Lambda') = 0$. Then $d(L_k, L_{k+1}) = 1$ for $0 \leq k \leq t-1$ and $L_t = L'$. The paths

$$\Gamma_{\min}(L, L') = (L_0, L_1, \dots, L_t)$$
 and $\Gamma_{\max}(L, L') = (M_0, M_1, \dots, M_t)$

are shortest paths from L to L', respectively called the min- and the max-shortest path from L to L'. For any Λ -normalized sequence $(N_0 = \Lambda, N_1, \dots, N_t = M)$ of a shortest path Γ' from L to L', one has

$$L_i \subset N_i \subset M_i \quad for \ 0 \leq i \leq t.$$

Proof. The existential part of the lemma is easy to verify by using Smith bases of the representatives $\Lambda \in L$ and $M \in L'$, and is left to the reader. First note that the minand max-shortest paths are in correspondence under the duality map $\Lambda \mapsto \Lambda^* = \lim_{k \to \infty} (\Lambda, \mathbb{O})$, so we only need to prove claims on one type of path. Note that the shortest path interval is symmetric. Indeed, letting $\Gamma(L', L) = (L'_0, L'_1, \dots, L'_t)$, we have $L_k = [M_{\Lambda} + z^k \Lambda]$ and $L'_{t-k} = [\Lambda_M + z^{t-k} M]$. By definition we have

$$\Lambda_M + z^{t-k}M = z^{-\upsilon_M(\Lambda)}\Lambda + z^{-\upsilon_M(\Lambda)-\upsilon_\Lambda(M)-k}M = z^{-\upsilon_M(\Lambda)-k}(z^k\Lambda + z^{-\upsilon_\Lambda(M)}M).$$

Therefore $L'_{t-k} = L_k$.

Now we turn to the proof of $\Lambda_i = M + \mathfrak{m}^i \Lambda \subset N_i \subset \Lambda \cap \mathfrak{m}^{i-t} M = M_i$ for the shortest path $\Gamma' = (N_0, \ldots, N_t)$, which we will prove by induction on the distance t = d(L, L'). For convenience, let any path $([\Lambda] = L_0, L_1, \ldots, L_{t-1}, [M] = L_t)$ be represented by its Λ -normalized sequence $(\Lambda, \Lambda_1, \ldots, \Lambda_{t-1}, \Lambda_t = M)$ of lattices $\Lambda_i \in L_i$ such that $v_{\Lambda}(\Lambda_i) = 0$. The claim is obvious for t = 1. Let us assume it holds for any shortest path between pairs of vertices in B_n whose distance is at most t-1. Suppose then that $d(\Lambda, M) = t$. Let $\Lambda = \Lambda_0 \supset \Lambda_1 \supset \cdots \supset \Lambda_{t-1} \supset M$ be the min-shortest path from Λ to M and let $\Lambda \supset \Lambda'_1 \supset \cdots \supset \Lambda'_{t-1} \supset M$ represent the path Γ of minimal length. By assumption, $d(\Lambda, \Lambda'_{t-1}) = t - 1$, therefore

we have $\Lambda \supset \Lambda'_{t-1} \supset z^{t-1}\Lambda$, and by definition, we have $\Lambda_{t-1} = M + z^{t-1}\Lambda$. Therefore, we get

$$\begin{split} \Lambda_{t-1} \cap \Lambda'_{t-1} &= (M + z^{t-1}\Lambda) \cap \Lambda'_{t-1} \\ &= M + (z^{t-1}\Lambda \cap \Lambda'_{t-1}) \quad (\text{because } M \subset \Lambda'_{t-1}) \\ &= M + z^{t-1}\Lambda \qquad (\text{since } z^{t-1}\Lambda \subset \Lambda'_{t-1}) \\ &= \Lambda_{t-1}. \end{split}$$

Thus $\Lambda_{t-1} \subset \Lambda'_{t-1}$ holds. On the other hand, $M_{t-1} = \Lambda \cap \mathfrak{m}^{-1} M$ is by construction the largest lattice containing M, contained in Λ and adjacent to M. Therefore $\Lambda_{t-1} \subset \Lambda'_{t-1} \subset M_{t-1}$. By the induction assumption, the claim is established. \Box

If $D = (d_1, ..., d_n)$ are the elementary divisors of M_{Λ} in Λ , then the lattices Λ_k on the (normalized) min-shortest path from Λ to M have elementary divisors $D_k = (\min(d_1, k), ..., \min(d_n, k))$ in Λ . The differences $\Delta_k = D_k - D_{k-1}$ are the intermediate elementary divisors as explained in the following:

(5)
$$\Lambda \xrightarrow{z^{\Delta_1}} \Lambda_1 \xrightarrow{z^{\Delta_2}} \Lambda_2 \longrightarrow \cdots \xrightarrow{z^{\Delta_{t-1}}} \Lambda_{t-1} \xrightarrow{z^{\Delta_t}} M_{\Lambda} \xrightarrow{z^{v_{\Lambda}(M)}} M.$$

Let for simplicity I = (1, ..., 1) and $\ell I = (\ell, ..., \ell)$.

Definition 18. For $D = (d_1, \ldots, d_n) \in \mathbb{Z}^n$, let $v = \min D$, let $t = \max D$ and

$$D_i = \left(\min(d_1 - v, i), \dots, \min(d_n - v, i)\right) \quad \text{for } 1 \le i \le t$$

The *elementary splitting* of *D* is the unique decomposition $D = \Delta_0 + \Delta_1 + \dots + \Delta_t$ where $\Delta_0 = vI \in \mathbb{Z}^n$ and $\Delta_k = D_k - D_{k-1} \in \{0, 1\}^n$ for $1 \le k \le t$.

We write $\Delta_i(D)$ if we need to specify the sequence D. The next straightforward lemma is left to the reader.

Lemma 19. Let $\mathbf{ED}_{\Lambda}(M) = \Delta_0 + \Delta_1 + \cdots + \Delta_t$ be the elementary splitting of $\mathbf{ED}_{\Lambda}(M)$ for Λ , $M \in \Lambda$. Then we have $\Delta_0 = v_{\Lambda}(M)I$ and the following statements hold:

- (i) $\delta_i = \operatorname{Tr} \Delta_i = \mu_M(\Lambda_i)$, where $\Lambda_i \in \Gamma_{\min}(\Lambda, M)$ as in (5).
- (ii) $\Delta_i(\mathbf{ED}_M(\Lambda)) = I \Delta_{d-i+1}(\mathbf{ED}_\Lambda(M))$ for $1 \le i \le t$. In other terms,

$$\mathbf{ED}_{M}(\Lambda) = (-t - v_{\Lambda}(M))I + (I - \Delta_{t}) + \dots + (I - \Delta_{1})$$

is the elementary splitting of $\mathbf{ED}_{M}(\Lambda)$.

(iii) Let $A = [\Lambda]$ and B = [M]. The gallery distance $\ell(A, B)$ is given by

$$\ell(A, B) = \sum_{i=1}^{t-1} (n-\delta_i)\delta_{i+1}.$$

300

Note that the sequence $\delta_1, \ldots, \delta_t$ is nonincreasing, and that we also have

$$\delta_i = |\{1 \leq j \leq n \mid d_j \geq i + v_{\Lambda}(M)\}|.$$

3. Birkhoff–Grothendieck trivializations

The central result in the theory of holomorphic vector bundles on $\mathbb{P}^1(\mathbb{C})$ is the Birkhoff–Grothendieck theorem, which states that any such bundle is isomorphic to a direct sum of line bundles. In this section, we take $X = \mathbb{P}^1(\mathbb{C})$, and we investigate what properties of the vector bundle can be retrieved by considering only the Bruhat–Tits building at a point $x \in X$.

3.1. *The Birkhoff–Grothendieck property.* According to Lemma 3, a holomorphic vector bundle $\mathscr{C} \in \mathscr{H}$ is completely described by the coset $\mathbf{\Lambda} = [\mathscr{C}] \in \mathscr{H} / \sim_x$ and the lattice $\Lambda = \mathscr{C}_x \in \mathbf{\Lambda}$. Let us take up the notations of Section 2 again. Let *V* denote the formalized meromorphic stalk $\widehat{\mathcal{V}}_x$ and let *B* be the corresponding Bruhat–Tits building. Let $B_0 = \{[M] | \mathscr{C}^M$ is trivial} be the subset of *trivializing lattices* of *B*.

Note 20. Strictly speaking, $L \in B_0$ contains exactly one lattice $M \in L$ such that the extension \mathscr{C}^M is trivial. Any other $\tilde{M} \in L$ gives a *quasitrivial* vector bundle $\mathscr{C}^{\tilde{M}}$. There is no real need to make the difference, since one gets a trivial bundle by simply tensoring with a line bundle.

Lemma 21. $B_0 \neq \emptyset$ for any $x \in X$ and $\Lambda \in \mathcal{H} / \sim_x$.

The existence of an analytic trivialization outside of an arbitrary point is a very general result on holomorphic bundles over complex analytic manifolds. This much more restrictive *algebraicity* statement follows easily from the Birkhoff–Grothendieck theorem.

Proof. Let $\mathscr{C} \in \Lambda$ be a representative of the coset Λ . Let $X^* = X \setminus \{x\}$, and Δ a neighborhood of x. The Birkhoff–Grothendieck theorem implies that there exists a sequence of integers $(a_1, \ldots, a_n) \in \mathbb{Z}^n$, and linearly independent sections $\sigma_1, \ldots, \sigma_n \in \Gamma(X^*, \mathscr{C})$ and $\tilde{\sigma}_1, \ldots, \tilde{\sigma}_n \in \Gamma(\Delta, \mathscr{C})$ such that $\tilde{\sigma}_i = z^{-a_i} \sigma_i$ for a local coordinate z at x. The section σ_i admits therefore a *meromorphic* extension at x, that we still denote σ_i . The trivial bundle $\mathcal{F} \in H_0$ spanned by $(\sigma_1 \ldots, \sigma_n)$ is a trivialization of \mathscr{C} at x as required.

We say that any such bundle \mathscr{F} is a *Birkhoff–Grothendieck trivialization of* \mathscr{E} *at* x. We'll often write BG for short. Let $BG_x(\mathscr{E}) \subset \mathscr{H}_0$ for the set of BG trivializations of \mathscr{E} at x. Note that any $\mathscr{F} \in BG_x(\mathscr{E})$ corresponds uniquely to a decomposition of \mathscr{E} as a direct sum of line subbundles (a *BG decomposition*)

(6)
$$\mathscr{E} = \bigoplus_{i=1}^{n} \mathscr{L}_{i}$$

Note 22. A line bundle on *P* is characterized by its degree. Recall that, for any decomposition (6), if the integers $a_i = \deg \mathcal{L}_i$ satisfy $a_1 \ge \cdots \ge a_n$, then the sequence $T(\mathcal{C}) = (a_1, \ldots, a_n)$ is unique and called the *type* of \mathcal{C} .

In the case where the coset $[\mathscr{E}] = \Lambda \in \mathscr{H} / \sim_x$ is fixed, we write

$$BG(\Lambda) = \{ M \in \Lambda^0 \, | \, \mathcal{E}^M \in BG_x(\mathcal{E}) \},\$$

and $T(\Lambda)$ instead of $T(\mathcal{E})$. The global sections $\mathcal{F}(X)$ of a trivial bundle $\mathcal{F} \in [\mathcal{E}]$ induce, by taking stalks at *x*, a form Y_M in the corresponding lattice $M = \mathcal{F}_x \in \Lambda$, that we call the *global form* of *M*. The link between the Birkhoff–Grothendieck theorem and the algebraic structure of the local lattices is then given by the following straightforward characterization.

Lemma 23. The trivial lattice M is in BG(Λ) if and only if there exists a Smith basis (e) of M with respect to Λ that is simultaneously a \mathbb{C} -basis of the global form Y_M of M.

In this case, the basis (e) will be called a *BG basis* for Λ . We state now the following result separately for further reference.

Proposition 24. Let $\mathscr{C} \in \mathscr{H}$ be a holomorphic vector bundle. For any BG trivialization \mathscr{F} of \mathscr{C} at any $x \in \mathbb{P}^1(\mathbb{C})$, the type $T(\mathscr{C})$ of the bundle \mathscr{C} is equal to the sequence of elementary divisors $\mathbf{ED}_{\Lambda}(M)$ of the stalk $M = \mathscr{F}_x$ with respect to $\Lambda = \mathscr{C}_x$ (viewed as lattices in \mathscr{V}_x).

The BG trivializations of a bundle \mathscr{E} are as a rule not unique, nor is any trivialization of \mathscr{E} necessarily BG. One of the aims of this section is to prove the following local characterization of BG trivializations. Recall that the Bruhat–Tits building is endowed with three distance maps: the graph-theoretical distance d, the canonical metric d and the gallery distance ℓ (which is actually a pseudodistance on the vertices).

Theorem 25. Let $\Lambda = \mathscr{C}_x$. Let \leq_{lex} be the lexicographic, and \leq_{lexrev} the reverse lexicographic orderings on \mathbb{Z}^n . Then we have:

- (i) $T(\mathscr{E}) \cong \min_{\preceq_{\operatorname{lex}}} \{ \operatorname{ED}_{\Lambda}(M)^{\searrow} \mid M \in \Lambda^0 \} = \max_{\preceq_{\operatorname{lexev}}} \{ \operatorname{ED}_{\Lambda}(M)^{\searrow} \mid M \in \Lambda^0 \}.$
- (ii) If $M \in \Lambda^0$ then the following conditions are equivalent:
 - (a) $M \in BG(\Lambda)$;
 - (b) $\mathbf{ED}_{\Lambda}(M) \cong T(\mathscr{E});$
 - (c) $d(\Lambda, M) = \min_{\tilde{M} \in \Lambda^0} d(\Lambda, \tilde{M});$
 - (d) $\ell(\Lambda, M) = \min_{\tilde{M} \in \Lambda^0} \ell(\Lambda, \tilde{M}).$

(iii) Moreover, if $M \in BG(\Lambda)$ then $d(\Lambda, M) = \min_{\tilde{M} \in \Lambda^0} d(\Lambda, \tilde{M})$.

302

In other terms, in the set of elementary divisors on Λ^0 with respect to \mathscr{C}_x , ordered by decreasing values, the minimum by direct lexicographic ordering coincides with the maximum by reverse lexicographic ordering, and this value is exactly the type $T(\mathscr{C})$. Moreover, any trivial lattice having $T(\mathscr{C})$ as elementary divisors is a BG trivialization of \mathscr{C} . Finally, the BG trivializations of a lattice Λ are exactly the geodesic projections of Λ onto Λ^0 for both metrics d and ℓ . This result will be proved in two steps: Section 3.2 is dedicated to proving the first item and the implications (a) \Rightarrow (b), (c), (d), (iii), and Section 3.3 to the converse implications.

3.1.1. Monopoles and BG trivializations. Recall from Section 1 that the group of monopole gauge transforms at x sends a trivial sheaf \mathcal{F} to a trivial sheaf $\widetilde{\mathcal{F}}$ such that $\mathcal{F}_{|X \setminus \{x\}} = \widetilde{\mathcal{F}}_{|X \setminus \{x\}}$. This group is described by the group of unimodular polynomial matrices $GL_n(\mathbb{C}[T])$, that is, matrices of the form

$$P = P_0 + P_1T + \dots + P_kT^k$$
 where $\exists \alpha \in \mathbb{C}^*$, det $P = \alpha$.

More precisely, assuming that x = 0, a matrix Π is a monopole at x if and only $\Pi = P(1/z)$ with $P \in GL_n(\mathbb{C}[X])$ and z is the standard coordinate on $\mathbb{P}^1(\mathbb{C})$.

A trivial bundle $\mathcal{F} \in \mathcal{H}_0$ is a trivialization of $\mathscr{C} \in \mathcal{H}$ at x if any basis $(\sigma) = (\sigma_1, \ldots, \sigma_n)$ of global sections of \mathcal{F} spans the stalk \mathscr{C}_y over the local ring $\mathbb{O}_y = (\mathbb{O}_X)_y$ for every $y \neq x$. Any other basis of meromorphic sections $(\tilde{\sigma})$ spans a trivialization $\tilde{\mathcal{F}}$ of \mathscr{C} at x if and only if the gauge Π from (σ) to $(\tilde{\sigma})$ is a monopole at x. In particular, since the group of units of $\mathbb{C}[X]$ is \mathbb{C}^* , a line bundle \mathscr{L} admits a unique trivialization \mathscr{L}^x at x.

Lemma 26. Two BG bases for Λ are related by $a - T(\Lambda)$ -parabolic monopole gauge.

Proof. Consider two BG trivializations M, \tilde{M} of Λ , as in the following diagram, where $D = T(\Lambda)$:



Since $v(P_{ij}) \ge 0$, the gauge Π satisfies $v(\Pi_{ij}) \ge d_j - d_i$. Therefore $\Pi_{ij} = 0$ as soon as $d_j > d_i$, which means that Π is (-D)-parabolic.

According to the previous section, \mathcal{F} is a BG trivialization of \mathcal{E} at x if there exists a basis $(\sigma) = (\sigma_1, \ldots, \sigma_n)$ of global sections of \mathcal{F} and an integer sequence $D = (d_1, \ldots, d_n)$, such that $(e) = (t^{-d_1}\sigma_1, \ldots, t^{-d_n}\sigma_n)$ spans the stalk \mathcal{E}_x over the local ring $\mathbb{O} = (\mathbb{O}_X)_x$, where t is a local coordinate at x. This coordinate t can be arbitrarily chosen, since the local behavior of \mathcal{E} only depends on the local ring \mathbb{O} . If we choose as coordinate t a *meromorphic* function on X, then the sections

 $(e) = (t^{-d_1}\sigma_1, \dots, t^{-d_n}\sigma_n)$ form a basis of global (meromorphic) sections of \mathcal{V} . The \mathbb{O}_X -module $\widetilde{\mathcal{F}}$ spanned by (e) in this case does coincide with \mathscr{E} at x, and differs from it at most on the support of the divisor of the function t.

When $X = \mathbb{P}^1(\mathbb{C})$, we can obviously find a function t with divisor (t) = x - yfor any arbitrary point $y \neq x$, that we call a *global coordinate at* x. In this case, the bundle $\tilde{\mathscr{F}}$ is a BG trivialization of \mathscr{C} at y. It is clearly independent of the global basis (σ) of \mathscr{F} , which is defined up to a (-D)-parabolic constant matrix $C \in GL_n(\mathbb{C})$, and of the specific meromorphic function t, which is only defined up to a nonzero constant. We call $t_y(\mathscr{F}) = \tilde{\mathscr{F}}$ the *transport at* y of the BG trivialization \mathscr{F} of \mathscr{C} at x.

Understood otherwise, this is the description of a nontrivial bundle \mathscr{E} by means of *two trivial* bundles \mathscr{F} and $\widetilde{\mathscr{F}}$ coinciding outside $\{x, y\}$, and glued along the cocycle $g = t^D$, where (t) = x - y.

3.1.2. *The Harder–Narasimhan flag.* The Harder–Narasimhan filtration $HN(\mathscr{E})$ of \mathscr{E} over $\mathbb{P}^1(\mathbb{C})$ can be obtained (see [Sabbah 2002, p. 65]) from a BG decomposition $\mathscr{E} = \bigoplus_{i=1}^n \mathscr{L}_i$ of \mathscr{E} as a direct sum of line bundles $\mathscr{L}_i \simeq \mathbb{O}(a_i)$ of the appropriate degree, by

$$F^k(\mathscr{E}) = \bigoplus_{i \mid a_i \ge k} \mathscr{L}_i.$$

Locally, the Harder–Narasimhan filtration can be described as follows. Let (e) be a BG basis for $\Lambda = \mathscr{C}_x$. The Harder–Narasimhan flag HN_{Λ} of $V = \mathscr{V}_x$ defined by

(7)
$$F^k = \bigoplus_{i|a_i \ge k} Ke_i$$

is independent of the BG basis for Λ . For a lattice M, define the residual HN flag

$$H^M_\Lambda = (\mathrm{HN}_\Lambda \cap M)/\mathfrak{m}M$$
 and, for simplicity, $H^\Lambda = H^\Lambda_\Lambda$.

Lemma 27. Let $\Lambda \in \Lambda$ and let HN_{Λ} be the HN flag of V.

- (i) If $M \in BG(\Lambda)$, then $H^M_{\Lambda} = F^M(\Lambda)$, hence H^{Λ} is transversal to $F^{\Lambda}(M)$.
- (ii) Conversely, for any flag F' which is transversal to H^{Λ} , there exists $M \in BG(\Lambda)$ such that $F^{\Lambda}(M) = F'$.

Proof. Let $T = T(\Lambda) = \text{diag}(a_1 I_{n_1}, \dots, a_s I_{n_s})$ with $a_i > a_{i+1}$ be the type of the lattice Λ , so that

$$\Lambda \xrightarrow{z^T} M$$

sends the basis (e) of Λ into the global basis (ε) = $z^T(e)$ of Y_M . Let $v_i = \sum_{1 \le k \le i} n_i$, and let generically \overline{x} denote the class modulo m of a vector x. The *i*-th component of the flag H^M_{Λ} is spanned by $(\overline{\varepsilon}_1, \ldots, \overline{\varepsilon}_{v_i})$, so $H^M_{\Lambda} = F^M(\Lambda)$ indeed

holds. Conversely, the (n-i+1)-th component of $F^{\Lambda}(M)$ is $F_i = \langle \bar{e}_{v_i+1}, \ldots, \bar{e}_n \rangle$ while the *i*-th component H_i of H^{Λ} satisfies $H_i = \langle \bar{e}_1, \cdots, \bar{e}_{v_i} \rangle$; hence both flags are transversal to each other. Any other BG trivialization \tilde{M} is obtained from (ε) by a monopole gauge transform Π such that $P = z^T \Pi z^{-T} \in GL_n(\mathbb{O})$. According to Lemma 26, Π is block-upper-triangular with respect to the blocks of equal elements of T, hence so is P. Let $P \in G_T(\mathbb{C})$. If z is a global coordinate, the matrix $z^{-T} P z^T$ is a monopole. The orbit of (ε) under the set of the constant T-parabolic matrices covers the set of all flags in E which are transversal to the image of $HN(\mathfrak{C})_x$ in E.

For any BG trivialization \mathcal{F} of \mathcal{C} at x, let $Y = \Gamma(X, \mathcal{F})$ be the \mathbb{C} -vector space of global sections of \mathcal{F} . The Harder–Narasimhan filtration $HN(\mathcal{E})$ also induces a canonical filtration $HN_{\Lambda}(Y)$ of \mathbb{C} -vector spaces of Y. To avoid defining new concepts, we will also refer to this filtration as the *Harder–Narasimhan filtration of* Y. Note that it depends solely on the lattice $\Lambda \in \mathbf{\Lambda}$.

3.2. *Modification of the type.* We wish to answer algebraically the following question: "what does the type of \mathscr{E} become when the stalk $\mathscr{E}_x = \Lambda$ at *x* is replaced by another lattice $\tilde{\Lambda}$?" It turns out that the question can be very explicitly answered when the lattice $\tilde{\Lambda}$ is not too far from Λ , namely at distance 1 in the graph-theoretic distance of the Bruhat–Tits building. The following proposition generalizes a result of Gabber and Sabbah.

Proposition 28. Let $\mathscr{C} \simeq \bigoplus_{i=1}^{n} \mathbb{O}(a_i)$ be a holomorphic vector bundle on $X = \mathbb{P}^1(\mathbb{C})$, with $a_1 \ge \cdots \ge a_n$, and let $x \in X$. Let $\tilde{\Lambda} \in \Lambda_x$ be a lattice such that $\mathfrak{m}_x \mathscr{C}_x \subset \tilde{\Lambda} \subset \mathscr{C}_x$. Let $E = \mathscr{C}_x/\mathfrak{m}_x \mathscr{C}_x$ be the local fiber at x, let

$$\boldsymbol{H}: H_0 = 0 \subset H_1 \subset \cdots \subset H_s = E$$

be the residual HN flag in E, and $W = \tilde{\Lambda}/\mathfrak{m}_x \mathscr{E}_x$ the image of $\tilde{\Lambda}$. Assume that the type of \mathscr{C} is written as

$$a = (\underbrace{a_1, \dots, a_1}_{n_1 \text{ times}}, \underbrace{a_2, \dots, a_2}_{n_2 \text{ times}}, \dots, \underbrace{a_s, \dots, a_s}_{n_s \text{ times}}).$$

Then the modified bundle $\mathcal{F} = \mathcal{E}^{\widetilde{\Lambda}}$ has type

$$\tilde{a} = (\underbrace{a_1, \dots, a_1}_{m_1 \text{ times}}, \underbrace{a_1 - 1, \dots, a_1 - 1}_{n_1 - m_1 \text{ times}}, \dots, \underbrace{a_s, \dots, a_s}_{m_s \text{ times}}, \underbrace{a_s - 1, \dots, a_s - 1}_{n_s - m_s \text{ times}})$$

where $m_i = \dim_{\mathbb{C}} H_i \cap W - \dim_{\mathbb{C}} H_{i-1} \cap W$.

Proof. This is explained in the following scheme. Let $\Lambda = \mathscr{C}_x$, and let t be a local coordinate at x (that we assume without loss of generality to be ∞). Let

 $D = \text{diag}(a_1, \dots, a_n)$ be the elementary divisors of the BG trivialization M in Λ (or, in this case, the type of \mathscr{C}).

(8)
$$Y_{M}:(y) \xrightarrow{t^{-D}P_{0}t^{D}} Y_{\tilde{M}}:(\tilde{y})$$

$$t^{D} \uparrow \qquad t^{D} \uparrow \qquad t^{D} \uparrow \qquad t^{D-T}$$

$$\Lambda:(e) \xrightarrow{P_{0}} (\varepsilon) \xrightarrow{t^{T}} (\tilde{e}): \tilde{\Lambda}$$

$$\downarrow^{\pi}$$

$$E = \Lambda/t\Lambda:(\bar{e}) \xrightarrow{P_{0}} (u): W$$

Let (e) be a basis of Λ , such that $(\sigma) = (t^D e)$ is a basis of the form Y_M . Under the canonical projection $\pi : \Lambda \to E = \Lambda/t\Lambda$, the HN filtration of Λ descends to a flag of \mathbb{C} -vector spaces $H: 0 = H_0 \subset \cdots \subset H_s = E$, and the quotient basis (\bar{e}) is a basis respecting this flag. Let $t\Lambda \subset \tilde{\Lambda} \subset \Lambda$ be the new lattice, and let $W \subset E$ be the subspace it is projected upon by π . Let (u) be a basis respecting both W and the flag H, and let P_0 be a change of basis from (e) to (u). Consequently, the matrix P_0 belongs to the parabolic subgroup P_H stabilizing the flag H; therefore it is block-upper-triangular, with blocks given by the equal elements among the a_i . Define now the basis (ε) of Λ as the image of (e) under the constant gauge P_0 . Here is where $d(\Lambda, \tilde{\Lambda}) \leq 1$ is important: the basis (ε) is a Smith basis of Λ (this would be not necessarily true if the lattices were further apart). Let $T = \text{diag}(t_1, \ldots, t_n)$ be the diagonal matrix such that $t_i = 0$ if $\pi(\varepsilon_i) \in W$ and $t_i = 1$ otherwise. Then $(\tilde{e}) = t^T(\varepsilon)$ is a basis of $\tilde{\Lambda}$. Let now $(\tilde{\varepsilon}) = t^D(\varepsilon)$ be the basis of \tilde{M} deduced from (ε). The matrix of the basis change from $\tilde{\Lambda}$ to \tilde{M} corresponding to the bases (σ) and ($\tilde{\varepsilon}$) is equal to $Q = t^{-D} P_0 t^{D} = (P_0)_{ii} t^{d_j - d_i}$. Now, since $P_0 \in P_H$, we have $(P_0)_{ij} = 0$ whenever $d_i - d_j < 0$. Therefore this gauge $Q = \frac{1}{t^k}Q_k + \dots + Q_0$ is a Laurent polynomial in t with only nonpositive terms, where moreover $Q_0 \in GL_n(\mathbb{C})$. Since $X = \mathbb{P}^1(\mathbb{C})$, it is possible to choose as local coordinate at ∞ a meromorphic function with divisor $(\infty) - (0)$, namely t = 1/z. Accordingly, Q is a polynomial in z, whereas det $Q = \det P_0 \in \mathbb{C}^*$. Hence $Q \in \operatorname{GL}_n(\mathbb{C}[z])$ is a monopole gauge. Since (σ) was a basis of global meromorphic sections of E, then $(\tilde{\varepsilon})$ also is. Therefore $\tilde{M} \in B_0$ is a trivializing lattice. Moreover, \widetilde{M} is a BG trivialization of both \mathscr{E} and $\mathscr{F} = \mathscr{E}^{\widetilde{\Lambda}}$, because the basis ($\widetilde{\varepsilon}$) is a Smith basis for Λ and $\tilde{\Lambda}$. Therefore, we can explicitly compute the new elementary divisors of Λ in M, which are given by the matrix D-T. Summing up, we see that the change of lattice has subtracted 1 from all the elementary divisors corresponding to the vectors of the basis (ε) whose image under π do not fall into the subspace W. We obtain the Harder-Narasimhan filtration of the modified bundle by reordering the type by decreasing values. Π

306

Proposition 28 generalizes the construction given by Sabbah based on an idea of O. Gabber in [Sabbah 2002, Proposition 4.11] (where only the case where W is one-dimensional is tackled). This result is independent of the valuation of $\tilde{\Lambda}$, and can thus be formulated as follows.

Corollary 29. Let Λ , $\tilde{\Lambda} \in \Lambda$, and let $H = (HN_{\Lambda} \cap \Lambda)/\mathfrak{m}\Lambda$ be the residual HN flag and $F = F^{\Lambda}(\tilde{\Lambda}) = (0 \subset W \subset \Lambda/\mathfrak{m}\Lambda)$ the flag induced by $\tilde{\Lambda}$ in $\Lambda/\mathfrak{m}\Lambda$. If $d(\Lambda, \tilde{\Lambda}) = 1$, then we have

$$T(\tilde{\Lambda}) \cong T(\Lambda)^{\searrow} - \sigma \left(\mathbf{ED}_{\Lambda}(\tilde{\Lambda})^{\nearrow} \right)$$

where $\sigma \in \rho(\mathbf{F}, \mathbf{H})$ represents the relative position of \mathbf{F} and \mathbf{H} . Actually, one can even say, putting $\tilde{\mathbf{H}} = (\text{HN}_{\tilde{\Lambda}} \cap \Lambda)/\mathfrak{m}\Lambda$, that

$$T(\tilde{\Lambda})^{\searrow} = \rho(\boldsymbol{H}, \boldsymbol{\tilde{H}}) T(\Lambda)^{\searrow} - \rho(\boldsymbol{F}, \boldsymbol{\tilde{H}}) \big(\mathbf{E} \mathbf{D}_{\Lambda}(\tilde{\Lambda})^{\nearrow} \big).$$

Proof. With the notations of the diagram (8), the basis (\tilde{y}) is a common BG basis for Λ and $\tilde{\Lambda}$. Therefore, the HN flags are spanned in (\tilde{y}) over K by the flags of indices $S^{\searrow}(D)$ and $S^{\searrow}(D-T)$, respectively. By applying any representative of the coset $\rho(H, \tilde{H})$ to D-T, one gets $T(\tilde{\Lambda})^{\searrow}$. Therefore, $T(\tilde{\Lambda})^{\searrow} = \rho(H, \tilde{H}) T(\Lambda)^{\searrow} - \rho(H, \tilde{H})\rho(F, H)(\text{ED}_{\Lambda}(\tilde{\Lambda})^{\checkmark})$. Take as representatives of the cosets $\rho(H, \tilde{H})$ and $\rho(F, H)$ their minimal length element (see, e.g., [Abramenko and Brown 2008, Proposition 2.23, p. 83]). Since the quotient basis $(\bar{\varepsilon})$ of $\Lambda/m\Lambda$ is adapted to the three flags H, \tilde{H} and F, we get $\rho(H, \tilde{H})\rho(F, H) = \rho(F, \tilde{H})$ [Abramenko and Brown 2008, Lemma 5.55, p. 236].

Corollary 30. *Let* $\Lambda \in \Lambda$ *. Then we have*

- (i) $BG(\Lambda) \neq \emptyset$,
- (ii) for any adjacent lattice $\tilde{\Lambda}$, we have $BG(\Lambda) \cap BG(\tilde{\Lambda}) \neq \emptyset$,
- (iii) all the elements in a chamber of B_n have a common BG trivialization.

Proof. According to Proposition 28, if a lattice Λ admits a BG trivialization, so does M for any adjacent lattice M. However, according to Proposition 17, two lattices are always connected by a path of adjacent lattices. Since a trivial lattice is its own BG trivialization, the two first results are simultaneously established. The third stems from the fact that any chamber appears as a complete flag in the quotient space of any representative. According to Bruhat's lemma, there is a basis respecting simultaneously two flags, in this case the one corresponding to the chamber and the one induced by the HN filtration.

3.2.1. An algorithm to compute a Birkhoff–Grothendieck trivialization. Let $x \in X$, and let $\Lambda = [\mathscr{C}]_x$ be the \sim_x -equivalence class of \mathscr{C} . Let $\Lambda = \mathscr{C}_x$ and $M = \mathscr{F}_x \in \Lambda$ where \mathscr{F} is an arbitrary trivialization of \mathscr{C} at x. In this local setting, we "see" the

global sections of \mathcal{F} as the *global form* $Y \subset M$. An arbitrary trivialization of \mathcal{E} is not necessarily BG. The following result is easily established.

Lemma 31. Let $M \in BG(\Lambda)$. Then we have

(i)
$$M \in BG(z^k \Lambda)$$
 for $k \in \mathbb{Z}$,

(ii) $M \in BG(\Lambda')$ for any lattice Λ' on the shortest path $\Gamma_{\min}(\Lambda, M)$.

Proposition 28 allows to construct effectively from an arbitrary trivialization M a BG one, by following shortest paths in the Bruhat–Tits building from M to \mathscr{C}_x . The following result shows how to start the construction. We can assume that $v_M(\Lambda) = 0$.

Lemma 32. Let $M \in \Lambda^0$ be a trivial lattice in B_n . For any lattice $\Lambda \in \Lambda$ such that $d(\Lambda, M) = 1$, we have $M \in BG(\Lambda)$. More precisely, let $Y \subset M$ be the global form of M. For any basis (e) respecting $W = \Lambda/\mathfrak{m}M$, the Y-basis (e_Y) is a Smith basis for Λ .

Proof. Assume that (*e*) satisfies the assumptions of the lemma. Then, according to Lemma 13, the *Y*-basis (*e_Y*) is a Smith basis of *M* for Λ . Since it is a basis of the global form of *M*, the result follows, and in particular, the Harder–Narasimhan filtration of the corresponding bundle is equal to the *Y*-lifting of the flag ($0 \subset W \subset M/\mathfrak{m}M$).

If there existed a Smith basis of M for Λ which spanned simultaneously Y_M , the lattice M would be a BG trivialization of Λ , and the sequence $\mathbf{ED}_{\Lambda}(M)$ would represent the type of \mathscr{C} .

Theorem 33. Let $\Lambda \in \mathbf{\Lambda} = [\mathcal{E}]$, and let $M \in \mathbf{\Lambda}^0$ be an arbitrary trivial lattice. Let $\mathbf{ED}_{\Lambda}(M) = \Delta_0 + \cdots + \Delta_t$ be the elementary splitting of $\mathbf{ED}_{\Lambda}(M)$. There exists a sequence of permutations $w_k \in S_n$ such that the type $T(\Lambda)$ satisfies

$$T(\Lambda) \cong \Delta_0 + w_1 \Delta_1 + w_2 \Delta_2 + \dots + w_t \Delta_t.$$

Proof. We prove the result first on the sequence $D = \mathbf{ED}_M(\Lambda) = -\mathbf{ED}_{\Lambda}(M)$. Let $\Gamma = (M = M_0, M_1, \dots, M_t = \Lambda_M)$ be (a normalized representative of) the min-shortest path from [M] to $[\Lambda]$. Let $D = \mathbf{ED}_M(\Lambda_M) = (\mathbf{k}_1 I_{n_1}, \dots, \mathbf{k}_s I_{n_s})$ where $0 = \mathbf{k}_1 < \dots < \mathbf{k}_s = t$. Consider the elementary splitting of D

(9)
$$D = \Delta_1 + \dots + \Delta_t \quad \text{where } \Delta_i = (0_{n-m_i}, I_{m_i})$$

for a nonincreasing sequence (m_i) . Recall that

$$D_k = \Delta_1 + \dots + \Delta_k = \mathbf{ED}_M(M_k).$$

Let (e) be a Smith basis of M for Λ , and let $(e^{(k)}) = z^{D_k}(e)$ for $0 \le k \le t$ (with $D_0 = 0$ by convention). The induced basis $(\bar{e}^{(k)})$ of $E_k = M_k/\mathfrak{m}M_k$ respects both flags $F^{M_k}(M)$ and $F^{M_k}(\Lambda)$. With the help of Proposition 28, we construct a BG trivialization \tilde{M}_{k+1} of the k-th element M_k of the shortest path Γ , which is simultaneously a BG trivialization of $\Lambda_M + \mathfrak{m}M_k = \Lambda_M + \mathfrak{m}(\Lambda_M + \mathfrak{m}^k M) = M_{k+1}$.

Let us describe this in more detail. By assumption, there exists a BG basis (σ) of \tilde{M}_k for M_k . Let

$$T = \mathbf{ED}_{\tilde{M}_k}(M_k)^{\nearrow}$$
 and $T' = \mathbf{ED}_{M_k}(\Lambda_M)^{\nearrow}$.

Let $(y) = z^T(\sigma)$ be the corresponding basis of M_k , and (ε) a Smith basis of M_k for Λ_M . The gauge U from (y) to (ε) can be factored as

(10)
$$U = U_0(I + tU') \quad \text{with } U_0 \in \mathrm{GL}_n(\mathbb{C}).$$

Since (σ) is a basis of the global form $Y_{\tilde{M}_k}$, the basis (\bar{y}) of $E = M_k/\mathfrak{m}M_k$ is strictly adapted to the HN flag

$$\boldsymbol{H}^{(k)} = (\mathrm{HN}_{M_k} \cap M_k) / \mathfrak{m}_{M_k}.$$

Similarly, $(\bar{\varepsilon})$ is strictly adapted to the flag $F^{(k)} = F^{M_k}(\Lambda_M)$. Let *B* be the standard Borel subgroup of $GL_n(\mathbb{C})$. By the Bruhat decomposition, the group $GL_n(\mathbb{C})$ is a disjoint union of double cosets:

$$\operatorname{GL}_n(\mathbb{C}) = \coprod_{w \in W} BwB$$

where *W* is the Weyl group $W = S_n$. The constant term U_0 of the gauge *U* belongs to only one such cell: let $w \in S_n$ be the label of the corresponding Schubert cell. We have a decomposition $U_0 = QP_wQ'^{-1}$ with $Q, Q' \in B$, where P_w is the matrix representation of the permutation *w*. Accordingly, the gauge transforms *Q* and *Q'* respect, respectively, the flags $H = H^{(k)}$ and $F = F^{(k)}$. In the quotient space *E*, we have:

$$E: (\overline{y}) \xrightarrow{Q} E: (\overline{y}')$$

$$\downarrow U_0 \qquad \qquad \downarrow P_w$$

$$E: (\overline{\varepsilon}) \xrightarrow{Q'} E: (\overline{\varepsilon}')$$

The gauge U_0 represents geometrically the change of a basis that spans the Harder– Narasimhan flag H to one that spans the flag F induced by Λ ; therefore w is a representative of the relative position $\rho(H, F)$.

Let $T' = \Delta_{k+1} + \cdots + \Delta_t$ be the elementary splitting of T'. Since $(\bar{\varepsilon})$ respects the flag F, it will in particular respect the trace of the first element $M_{k+1} = \Lambda + \mathfrak{m}M_k$ of the shortest path $\Gamma_{\min}(M_k, \Lambda)$; therefore any lifting of $(\bar{\varepsilon})$ will be a Smith basis of M_{k+1} with elementary divisors Δ_{k+1} . Put $T'' = T' - \Delta_{k+1}$. The previous

scheme gets thus lifted to the following complete picture.

$$Y_{\tilde{M}_{k}} : (\sigma) \xrightarrow{t^{T} \mathcal{Q}t^{-T}} Y_{\tilde{M}_{k+1}} : (\sigma') \xrightarrow{t^{T+w^{-1}\Delta_{k+1}}} V_{\tilde{M}_{k+1}} : (\sigma') \xrightarrow{t^{T+w^{-1}\Delta_{k+1}}} V_{\tilde{M}_{k+1}} = V_{\tilde{M}_{k}} : (\sigma') \xrightarrow{t^{T+w^{-1}\Delta_{k+1}}} M_{k+1} = V_{\tilde{M}_{k}} : (\sigma') \xrightarrow{t^{T+w^{-1}\Delta_{k+1}}}$$

As a result, the elementary divisors of M_{k+1} with respect to the common BG trivialization \tilde{M}_{k+1} of M_k and M_{k+1} are not $T + \Delta_{k+1}$ (as with respect to \tilde{M}_k), but $T + w^{-1}\Delta_{k+1}$; namely, the elements of Δ_{k+1} have been twisted according to the permutation $w^{-1} = \rho(F, H)$ indexing the Bruhat cell that contains the matrix $U_0 \in \operatorname{GL}_n(\mathbb{C})$. The resulting matrix $\tilde{T} = T + w^{-1}\Delta_{k+1}$ is not necessarily ordered by increasing values: therefore we cannot ensure that $w_{k+1} = \rho(F, H)$, since the *ordered* diagonal has the form $\sigma T + \sigma w^{-1}\Delta_{k+1}$. According to Corollary 29, however, we know that we can take

$$\sigma = \rho(\boldsymbol{H}^{M_k}, \boldsymbol{H}_{M_k}^{M_{k+1}}) = \rho(\mathrm{HN}_{M_k}, \mathrm{HN}_{M_{k+1}}).$$

Thus

$$\widetilde{T}^{\nearrow} = \rho(\mathrm{HN}_{M_k}, \mathrm{HN}_{M_{k+1}})T + \rho(H^{M_k}, H^{M_{k+1}}_{M_k})\rho(F^{(k)}, H^{M_k})\Delta_{k+1}$$

Putting $T_k = \mathbf{ED}_{\widetilde{M}_k}(M_k)^{\nearrow}$, we get

(11)
$$T_{k+1} = \rho(\text{HN}_{M_k}, \text{HN}_{M_{k+1}})T_k + \rho(F^{M_k}(\Lambda), H^{M_{k+1}}_{M_k})\Delta_{k+1}.$$

At the end of at most t steps, the lattice \tilde{M}_t is a BG trivialization of Λ_M , and thus of Λ . We have $\tilde{T} = w_1 \Delta_1 + w_2 \Delta_2 + \cdots + w_t \Delta_t$ such that

$$\tilde{M}_t \xrightarrow{z^{\tilde{T}}} \Lambda_M \xrightarrow{z^{v_M(\Lambda)}} \Lambda.$$

Since
$$T(\Lambda) = \mathbf{ED}_{\Lambda}(M_t) = -T - v_M(\Lambda)$$
, we get

$$T(\Lambda) = -\Delta_0 - w_1 \Delta_1 - w_2 \Delta_2 - \dots - w_t \Delta_t$$

$$= -\Delta_0 - tI + w_1(I - \Delta_1) + (I - w_2 \Delta_2) + \dots + (I - w_t \Delta_t)$$

$$= (-\Delta_0 - tI) + w_t(I - \Delta_t) + w_{t-1}(I - \Delta_{t-1}) + \dots + w_1(I - \Delta_1).$$

By Lemma 19, the result is established. Note that a similar relation holds for $\mathbf{ED}_{M}(\Lambda)$, and that if $T(\Lambda) = \Delta_{0} + w_{1}\Delta_{1} + w_{2}\Delta_{2} + \dots + w_{t}\Delta_{t}$, then $T(\Lambda_{k}) = \Delta_{0} + w_{k+1}\Delta_{k+1} + \dots + w_{t}\Delta_{t}$ for the normalized elements Λ_{k} of the shortest path $\Gamma_{\min}(\Lambda, M)$.

We can even specify an actual (noneffective) formula for the permutations w_i appearing in $T(\Lambda)^{\nearrow} = \Delta_0 + w_1 \Delta_1 + w_2 \Delta_2 + \cdots + w_t \Delta_t$, attached to the shortest path

$$\Lambda \xrightarrow{z^{\Delta_1}} \Lambda_1 \xrightarrow{z^{\Delta_2}} \Lambda_2 \longrightarrow \cdots \xrightarrow{z^{\Delta_{t-1}}} \Lambda_{t-1} \xrightarrow{z^{\Delta_t}} M_{\Lambda} \xrightarrow{z^{\Delta_0}} M;$$

namely,

$$w_{i} = \rho(\mathrm{HN}_{\Lambda_{1}}, \mathrm{HN}_{\Lambda}) \cdots \rho(\mathrm{HN}_{\Lambda_{i-1}}, \mathrm{HN}_{\Lambda_{i-2}}) \rho(\boldsymbol{F}^{\Lambda_{i}}(\Lambda), \boldsymbol{H}_{\Lambda_{i}}^{\Lambda_{i-1}}).$$

The formula cannot be reduced as in Corollary 29 since there is not necessarily an apartment adapted for three such consecutive flags.

3.2.2. *The abacus.* Theorem 33 has a nice combinatorial interpretation in terms of Young diagrams. For any integer sequence $D = (d_1, \ldots, d_n) \in \mathbb{Z}^n$, let $D = \Delta_0 + \cdots + \Delta_t$ be the elementary splitting of D. For $w = (w_1, w_2, \ldots, w_t) \in (S_n)^t$, let $w(D) = \Delta_0 + w_1 \Delta_1 + w_2 \Delta_2 + \cdots + w_t \Delta_t$. The *abacus* of D is the set of sequences

$$ab(D) = \{A \in \mathbb{Z}^n \mid \exists w \in (S_n)^t, A = w(D)\}.$$

The name comes from the following analogy. Let Y(D) be the Young diagram with *n* rows whose respective lengths are the elements of $D_0 = D - \min D \in \mathbb{N}^n$, assumed to be arranged in increasing order (rows of length 0 are included). Assume from now on that min D = 0. Let $d_i^* = |\{1 \le j \le n \mid d_j \ge i\}|$ be the number of boxes in the *i*-th column of Y(D). The sequence $D^* = (d_1^*, \ldots, d_t^*)$ where $t = \max D$ is the *dual* sequence of D. The elementary splitting of D is given by $\Delta_i = (0_{n-d_i^*}, I_{d_i^*})$. In view of Lemma 19, put

$$\ell(D) = \sum_{i=1}^{t-1} (n - d_i^*) d_{i+1}^*.$$

Any sequence A in the abacus ab(D) of $D \in \mathbb{Z}^n$ comes from a box diagram obtained from Y(D) by allowing to move some boxes *only vertically* inside the



Figure 1. The Young diagram Y(D) of the sequence D = (0, 1, 1, 2, 4, 4), and an element $A = (2, 2, 2, 1, 2, 3) \in ab(D)$, featuring the moved boxes (in shades of gray), so that $A^{\nearrow} = (1, 2, 2, 2, 2, 3)$. The complement (thin gray line) corresponds to the sequence $D^0 = (0, 0, 2, 3, 3, 4)$, and the bijection from Lemma 35, to reversing the arrows.

whole corresponding column of length n. The diagram thus obtained, that we call an *abacus diagram*, can have nonadjacent boxes (as shown in Figure 1). For any such diagram, the sequence (a_1, \ldots, a_n) of number of boxes contained in each of the n rows is the required element of ab(D).

Lemma 34. Let $D = \Delta_1 + \dots + \Delta_t$ be the elementary decomposition of $D \in (\mathbb{Z}^n)^{\nearrow}$, and let $A = w_1 \Delta_1 + \dots + w_t \Delta_t \in ab(D)$. There exist $w'_1, \dots, w'_t \in S_n$ such that

$$w'_1 \Delta_1 + \dots + w'_i \Delta_i = (w_1 \Delta_1 + \dots + w_i \Delta_i)^{\nearrow}$$
 for all $1 \le i \le d$.

Proof. We proceed by induction on the number *s* of columns in Y(D). Let *Y* be the Young diagram for $D = (d_1, \ldots, d_n)$ and \tilde{Y} the one obtained from *Y* by erasing the last column, i.e., corresponding to the sequence $\tilde{T} = (\Delta_1, \ldots, \Delta_{t-1})$. Let $\tilde{D} = (\tilde{d}_1, \ldots, \tilde{d}_n)$ be the associated sequence. Then we have $d_i = \tilde{d}_i$ for $1 \le i \le n - d_t^*$ and $d_i = \tilde{d}_i + 1$ for $n - d_t^* + 1 \le i \le n$. An element $A \in ab(D)$ given by the permutations $w = (w_2, \ldots, w_t)$ corresponds uniquely to the pair (\tilde{A}, w_t) where $\tilde{A} \in ab(\tilde{D})$ is given by the restriction $\tilde{w} = (w_2, \ldots, w_{t-1})$.

The claim is clear for t = 1, so assume that it is established for any \tilde{D} such that $\Delta \tilde{D} \leq t-1$. We have $D = \Delta_1 + \dots + \Delta_t = \tilde{D} + \Delta_t$ where $\tilde{D} = \Delta_1 + \dots + \Delta_{t-1}$ is the elementary decomposition of \tilde{D} . Let $A \in ab(D)$ be described by the number a_i of boxes in the *i*-th row for $1 \leq i \leq n$, and let $\tilde{A} = (\tilde{a}_1, \dots, \tilde{a}_n)$ be the restriction of A to the t-1 first columns. Let $\mathcal{J} = \{i \mid a_i = \tilde{a}_i + 1\}$. Note that $\mathcal{J} = w_t(\{n-d_t^*+1,\dots,n\})$, and thus $|\mathcal{J}| = d_t^*$. According to the assumption, we can find w'_1, \dots, w'_{t-1} such that $w'_1 \Delta_1 + \dots + w'_i \Delta_i = (w_1 \Delta_1 + \dots + w_i \Delta_i)^{\mathcal{I}}$ for all $1 \leq i \leq t-1$. In particular,

we can assume that $\tilde{a}_1 \leq \cdots \leq \tilde{a}_n$. If there exists an index *i* such that $a_i > a_{i+1}$ holds, then we have necessarily $\tilde{a}_i = \tilde{a}_{i+1} = a_{i+1}$ and $a_i = \tilde{a}_i + 1$, so $i \in \mathcal{F}$ and $i + 1 \notin \mathcal{F}$. Exchanging *i* and i + 1 will not change the resulting sequence *A* when reordered, so one can avoid the inversion by putting $w'_t = (i, i + 1)w_t$, so that $i \notin \mathcal{F}$ and $i + 1 \in \mathcal{F}$ instead. By repeating this procedure for all the inverted indices, we get the claimed result for *t*.

For an $n \times t$ rectangular matrix $M = (M_{i,j})$, define the row sum vector

$$r(M) = (r_1, \dots, r_t), \quad r_j = \sum_{i=1}^n M_{i,j},$$

and column sum vector

$$c(M) = (c_1, \dots, c_n), \quad c_i = \sum_{j=1}^t M_{i,j}.$$

An element $A \in ab(D)$ in the abacus of D can also be seen as a (0, 1) rectangular $n \times t$ matrix $\mathcal{A} = (A_{i,j})$ where $A_{i,j} = 1$ whenever $i \in w_j(\{n - d_j^* + 1, \ldots, n\})$ holds. This matrix \mathcal{A} has row sum vector $r(\mathcal{A}) = D^*$, and column sum vector $c(\mathcal{A}) = A$. Recall that for two sequences $p = (p_1, \ldots, p_n)$ and $q = (q_1, \ldots, q_t)$ having the same sum, one says that p dominates q when $\sum_{i=1}^{k} q_i \leq \sum_{i=1}^{k} p_i$ for all integers k, where one completes the missing elements with 0. The following lemma sums up the behavior of the quantities introduced in Section 2.1 under the abacus transformations.

Lemma 35. Let $D \in \mathbb{Z}^n$.

- (i) The maps max, ∆, i, ℓ and || · || all attain their maximum over ab(D) at D. The map min attains its minimum at D.
- (ii) For any sequence $A \in ab(D)$, the following hold:
 - (a) max $A = \max D \iff i(A) = i(D)$,
 - (b) $\Delta A = \Delta D \iff \max A = \max D$ and $\min A = \min D$,
 - (c) $||A|| = ||D|| \iff \ell(A) = \ell(D) \iff A^{\nearrow} = D^{\nearrow}$.

Moreover, the map $A \mapsto \max D - A$ is a bijection between $\operatorname{ab}(D_0)$ and $\operatorname{ab}(D^0)$.

Proof. We will only prove the claims about $\ell(D)$ and ||D||, since the others are easily derived from their definitions. Although the assertions are similar, the methods of proof will be different. Let us start with ℓ . Assume without loss of generality that min D = 0. According to Lemma 34, one can also assume that $A = A^{\nearrow}$ holds. Define the integers k_i by induction as $k_1 = a_1^* - d_1^*$ and $k_i = a_i^* - d_i^* + k_{i-1}$ for $i \ge 2$. Since A can be seen as a (0, 1) matrix \mathcal{A} with $r(\mathcal{A}) = D^*$ and $c(\mathcal{A}) = A$, the Gale–Ryser theorem [Krause 1996] implies that $c(\mathcal{A})^* = A^*$ dominates $r(\mathcal{A}) = D^*$.

Therefore

$$\sum_{i=1}^k d_i^* \leqslant \sum_{i=1}^k a_i^*$$

for all integers k. Hence $k_i \ge 0$ for all i. Putting $k_0 = k_t = 0$, we then have $a_i^* = d_i^* - k_{i-1} + k_i$ for $1 \le i \le d$. According to the definition, we have

$$\ell(A) = \sum_{i=1}^{t-1} (n - a_i^*) a_{i+1}^* = \sum_{i=1}^{t-1} (n - d_i^* + k_{i-1} - k_i) (d_{i+1}^* - k_i + k_{i+1})$$

=
$$\sum_{i=1}^{t-1} (n - d_i^*) d_{i+1}^* + (k_{i-1} - k_i) (d_{i+1}^* - k_i + k_{i+1}) - (k_i - k_{i+1}) (n - d_i^*)$$

=
$$\ell(D) - \sum_{i=1}^{t-1} (k_i - k_{i+1}) (n - d_i^* + d_{i+1}^* + k_{i+1} - k_i).$$

After some algebraic manipulation, we get

$$\ell(A) - \ell(D) = -k_1 \left(n - d_1^* + d_2^* - k_1 + k_2 - (d_3^* - k_2 + k_3) \right)$$

$$\cdots - \sum_{i=2}^{t-2} k_i \left(d_{i+1}^* - k_i + k_{i+1} - (d_{i+2}^* - k_{i+1} + k_{i+2}) + d_{i-1}^* - d_i^* \right)$$

$$\cdots - k_{t-1} \left(d_t^* - k_{t-1} + k_t + d_{t-2}^* - d_{t-1}^* \right)$$

$$= -k_1 \left(n - d_1^* + a_2^* - a_3^* \right) \cdots - \sum_{i=2}^{t-2} k_i \left(a_{i+1}^* - a_{i+2}^* + d_{i-1}^* - d_i^* \right)$$

$$\cdots - k_{t-1} \left(a_t^* + d_{t-2}^* - d_{t-1}^* \right).$$

Since both sequences (d_i^*) and (a_i^*) are nonincreasing, we get $\ell(A) \leq \ell(D)$. Moreover, if $\ell(A) = \ell(D)$ holds, then all these terms must be zero. Let us prove then by induction that $k_i = 0$ for all i. If $k_1 \neq 0$, then one has $n = d_1^*$ and $a_2^* = a_3^*$. Since $a_1^* \geq d_1^*$, one must therefore have $a_1^* = n$, so $k_1 = 0$ holds, and hence we get $d_1^* = a_1^*$. Assume now that $k_j = 0$ and $a_j^* = d_j^*$ hold for $j \leq i - 1$. If $k_i \neq 0$ then one has $a_i^* > d_i^*$ and $d_{i-1}^* = d_i^*$ (and $a_{i+1}^* = a_{i+2}^*$ also). But then $d_{i-1}^* = d_i^* = a_{i-1}^* \geq a_i^* \geq d_i^*$, so we have $d_{i-1}^* = d_i^* = a_{i-1}^* = a_i^*$. Thus $k_i = a_i^* - d_i^* + k_{i-1} = 0$, and so the result for ℓ is established.

Let us turn now to $\|\cdot\|$. We can assume here that $w_1 = \text{id.}$ We proceed by induction on the number *t* of columns in the Young diagram Y(D). Like in the proof of the previous result, the restricted column sequence $T' = (\Delta_1, \ldots, \Delta_{t-1})$ corresponds to the Young diagram Y(D') of $D' = (d'_1, \ldots, d'_n)$ with $d_i = d'_i$ for $1 \le i \le n - d^*_t$ and $d_i = d'_i + 1$ for $n - d^*_t + 1 \le i \le n$. By the König–Huygens

314

identity, we have

$$\begin{split} \|D\| &= \sum_{i=1}^{n} \left(d_{i} - \frac{\operatorname{Tr} D}{n} \right)^{2} = \sum_{i=1}^{n} d_{i}^{2} - \frac{(\operatorname{Tr} D)^{2}}{n} \\ &= \sum_{i=1}^{n-d_{t}^{*}} (d_{i}')^{2} + \sum_{i=n-d_{t}^{*}+1}^{n} (d_{i}'+1)^{2} - \frac{(\operatorname{Tr} D'+d_{t}^{*})^{2}}{n} \\ &= \sum_{i=1}^{n} (d_{i}')^{2} - \frac{(\operatorname{Tr} D')^{2}}{n} + d_{t}^{*} + 2 \sum_{i=n-d_{t}^{*}+1}^{n} d_{i}' - \frac{2d_{t}^{*} \operatorname{Tr} D' + (d_{t}^{*})^{2}}{n} \\ &= \|D'\| + d_{t}^{*} + 2d_{t}^{*}(t-1) - \frac{2d_{t}^{*} \operatorname{Tr} D' + (d_{t}^{*})^{2}}{n}. \end{split}$$

Let $A = w(D) \in ab(D)$ with $w = (w_2, ..., w_t)$ and $A' = w'(D') \in ab(D')$ where $w' = (w_2, ..., w_{t-1})$.

For t = 1, the claim is clear, for ||w(D)|| = ||D|| and $w(D)^{\nearrow} = D$ for any $w \in (S_t)^n$. Assume then that for any diagram Y' = Y(D') with at most t - 1 columns, we have $||A'|| \le ||D'||$ for $A' \in ab(D')$, and ||A'|| = ||D'|| if and only if $(A')^{\nearrow} = D'$. Let Y = Y(D) have t columns. Let $A \in ab(D)$ be described by the number a_i of boxes in the *i*-th row for $1 \le i \le n$, and let $A' = (a'_1, \ldots, a'_n)$ be the restriction of A to the t - 1 first columns. Let again $\mathcal{J} = \{i \mid a_i = a'_i + 1\}$. Then one has

$$\begin{split} \|A\| &= \sum_{i=1}^{n} a_{i}^{2} - \frac{(\operatorname{Tr} A)^{2}}{n} = \sum_{i \in \mathcal{I}} (a_{i}'+1)^{2} + \sum_{i \notin \mathcal{I}}^{n} (a_{i}')^{2} - \frac{(\operatorname{Tr} A'+d_{t}^{*})^{2}}{n} \\ &= \sum_{i=1}^{n} (a_{i}')^{2} - \frac{(\operatorname{Tr} A')^{2}}{n} + |\mathcal{I}| + 2\sum_{i \in \mathcal{I}}^{n} a_{i}' - \frac{2d_{t}^{*} \operatorname{Tr} A' + (d_{t}^{*})^{2}}{n} \\ &= \|A'\| + d_{t}^{*} - \frac{2d_{t}^{*} \operatorname{Tr} A' + (d_{t}^{*})^{2}}{n} + 2\sum_{i \in \mathcal{I}}^{n} a_{i}'. \end{split}$$

By construction, we have $\operatorname{Tr} A' = \operatorname{Tr} D'$, so, under the induction assumption, we get

$$\|A\| \leq \|D'\| + d_t^* - \frac{2d_t^* \operatorname{Tr} D' + (d_t^*)^2}{n} + 2\sum_{i \in \mathcal{I}}^n a_i' = \|D\| + 2\left(\sum_{i \in \mathcal{I}} a_i' - d_t^*(t-1)\right).$$

Since $a'_i \leq t-1 = \max D'$ holds by construction, and $|\mathcal{Y}| = d^*_t$, we get $||A|| \leq ||D||$. Moreover, ||A|| = ||D|| can only happen when $a'_i = t - 1$ for $i \in \mathcal{Y}$. In this case, we get

$$||A'|| = ||D|| - \left(d_t^* + 2d_t^*(t-1) - \frac{2d_t^*\operatorname{Tr} D' + (d_t^*)^2}{n}\right) = ||D'||.$$

By the induction assumption, we have $(A')^{\nearrow} = D'$, and by the definition of \mathcal{J} , we get $a_i = a'_i + 1 = d = \max D$ for $i \in \mathcal{J}$. Therefore, we get d_t^* elements in A which are equal to d, and hence $A^{\nearrow} = D$.

3.2.3. Local criteria for BG trivializations. In this section, we use the fact that the type $T(\Lambda)$ is an element of $ab(ED_{\Lambda}(M))$ for any trivial $M \in \Lambda^0$ to derive local criteria satisfied by the BG trivializations. Let $d(\Lambda, \Lambda^0) = \min_{M \in \Lambda^0} d(\Lambda, M)$.

Lemma 36. Let $\Lambda \in \Lambda$ be a lattice. For any $M \in \Lambda^0$, we have

$$d(\Lambda, M) = d(\Lambda, \Lambda^0) \iff v_{\Lambda}(M) = \max_{\tilde{M} \in \Lambda^0} v_{\Lambda}(\tilde{M}) \text{ and } v_{M}(\Lambda) = \max_{\tilde{M} \in \Lambda^0} v_{\tilde{M}}(\Lambda).$$

If
$$M \in BG(\Lambda)$$
, then $\delta(\Lambda, M) = \min_{\tilde{M} \in \Lambda^0} \delta(\Lambda, \tilde{M})$ for $\delta \in \{d, d, \ell\}$.

Proof. Let $\Lambda^0 \ni \tilde{M} \xrightarrow{z^D} \Lambda$ represent the elementary divisors of an arbitrary trivialization \tilde{M} of Λ . The BG algorithm of Section 3.2.1 can be applied to \tilde{M} to obtain $M \in BG(\Lambda)$, as in the following scheme, where $A = w(D) \in ab(D)$:

$$\Lambda^0 \ni \tilde{M} \stackrel{z^D}{\longleftrightarrow} \Lambda \stackrel{z^A}{\longrightarrow} M \in \mathrm{BG}(\Lambda) \ .$$

By definition, we have $T(\Lambda) \cong A$. By Lemma 35, for any $\tilde{M} \in \Lambda^0$, we have $v_{\Lambda}(\tilde{M}) = \min D \leq \min A = v_{\Lambda}(M)$ and $v_{\tilde{M}}(\Lambda) = -\max D \leq -\max A = v_{M}(\Lambda)$. This holds for any trivial \tilde{M} ; hence, for any BG trivialization $M \in BG(\Lambda)$, we get

$$v_{\Lambda}(M) = \max_{\tilde{M} \in \Lambda^0} v_{\Lambda}(\tilde{M})$$
 and $v_{M}(\Lambda) = \max_{\tilde{M} \in \Lambda^0} (v_{\tilde{M}}(\Lambda)).$

On the other hand, we have, for any $\tilde{M} \in \Lambda^0$,

$$d(\Lambda, \tilde{M}) = -v_{\Lambda}(\tilde{M}) - v_{\tilde{M}}(\Lambda) \ge -v_{\Lambda}(M) - v_{M}(\Lambda) = d(\Lambda, M) \text{ for } M \in BG(\Lambda).$$

The rest in a direct consequence of Lemma 35.

Proposition 37. *Let* $\Lambda \in \Lambda$ *. Then*

$$T(\Lambda) = \max_{\leq_{\text{lex}}} \{ \mathbf{ED}_{\Lambda}(M)^{\nearrow} \mid M \in \Lambda^0 \} = \min_{\leq_{\text{lexrev}}} \{ \mathbf{ED}_{\Lambda}(M)^{\nearrow} \mid M \in \Lambda^0 \}.$$

Proof. Let $M \in \Lambda^0$, and let $D = (d_1, \ldots, d_n) = \mathbf{ED}_{\Lambda}(M)^{\nearrow}$. By Theorem 33, there exists $w = (w_1, \ldots, w_t) \in (S_n)^t$ such that $w(D) \cong T(\Lambda)$, and we can, by Lemma 34, ensure that $w(D) = T(\Lambda)^{\nearrow}$. Representing D as its Young diagram Y(D), the sequence $T(\Lambda)$ is given by an abacus diagram $A \in ab(D)$ derived from Y(D) by moving boxes vertically (downwards), and such that the number

316

of boxes t_i in row *i* weakly increases with the row index (i.e., $t_{i+1} \ge t_i$). Then we have $d_1 = v_{\Lambda}(M) \le t_1 = \max_{\tilde{M} \in \Lambda^0} v_{\Lambda}(\tilde{M})$. Suppose that $d_1 = t_1$, and assume that $d_2 > t_2$ holds. Accordingly some boxes from the second row must be lowered. Therefore, we necessarily have $t_1 > d_1$. This contradicts the maximality assumption on $v_{\Lambda}(M)$. Thus $d_2 \le t_2$, and so $t_2 = \max_{M \in \Lambda^0, v_{\Lambda}(M) = t_1} d_2^{\Lambda}(M)$. Let $i = \max\{j \mid d_k = t_k \text{ for } k \le j\}$ and assume that $d_{i+1} > t_{i+1}$ holds. The same argument holds and shows that $d_{i+1} \le t_{i+1}$. Therefore we have proved that $t_i = \max_{M \in \Lambda^0, d_k^{\Lambda}(M) = t_k, 1 \le k \le i-1} d_i^{\Lambda}(M)$. A similar argument starting from $d_n = -\max v_M(\Lambda)$ proves the second relation. \Box

This result establishes the first part of Theorem 25. The proof will be complete when we prove that M is indeed a BG trivialization of Λ if and only if $M \in \Lambda^0$ and $ED_{\Lambda}(M) \cong T(\Lambda)$. This will be established in the next section.

3.3. The permutation lemma. In the previous section, we showed how to construct a lattice M whose global form Y_M contains a Smith basis for a given lattice Λ . The following result allows us to give a fairly complete geometric view of the set of BG trivializations. We recall that a *principal minor* of a matrix $A \in \mathfrak{gl}_n(\mathbb{C})$ is a minor $A_{I,I}$ where $I \subset [n]$ obtained by deleting rows and columns whose indices are not elements of I.

Proposition 38 (permutation lemma). Let $D = (d_1, \ldots, d_n) \in \mathbb{Z}^n$ be an integer sequence and $P \in GL_n(\mathbb{C}[t])$ a lattice gauge.

(1) (Bolibrukh) There exist a permutation $\tau \in S_n$ and a lattice gauge $\tilde{P} \in GL_n(\mathbb{C}[t])$ such that

$$\Pi = t^{-D} P^{-1} t^{\tau D} \tilde{P} \in \operatorname{GL}_n(\mathbb{C}[t^{-1}]),$$

where $t^{D} = \text{diag}(t^{d_{1}}, ..., t^{d_{n}})$ and $\tau D = (d_{\tau(1)}, ..., d_{\tau(n)})$.

- (2) There exists moreover a lattice gauge $Q \in GL_n(\mathbb{O})$ such that $t^D \Pi = Q t^D$.
- (3) Furthermore, one can choose $\sigma = 1$ in item (1) if and only if all principal minors of P(0) indexed by the elements of the ascending flag \mathbf{D}^{\nearrow} are nonzero.

We will give a self-contained proof of this result, following for the first item, due to Bolibrukh, the same lines as the proof of this lemma given in [Ilyashenko and Yakovenko 2007]. Item (2) is, up to our knowledge, new, as well as the necessity statement in (3) (sufficiency appears in the cited work). The proof proceeds by induction, using the following simple lemma.

Lemma 39. Let $d \leq n$ and

$$T = \begin{pmatrix} I_d & 0\\ 0 & 0_{n-d} \end{pmatrix}.$$

Let $H = \begin{pmatrix} A & B \\ C & D \end{pmatrix} \in \operatorname{GL}_n(\mathbb{C}\llbrackett\rrbracket)$ be a lattice gauge matrix, decomposed as a 2×2-block matrix according to the blocks of T. If det $A(0) \neq 0$, then there exists a monopole gauge matrix

$$\Pi = \begin{pmatrix} I_d & t^{-1} \widetilde{\Pi} \\ 0 & I_{n-d} \end{pmatrix}$$

with $\tilde{\Pi}$ a constant matrix, such that $\tilde{H} = t^{-T} H t^T \Pi$ is a lattice gauge matrix; that is, $\tilde{H} \in GL_n(\mathbb{O})$.

Proof. For simplicity, put $M_0 = M(0)$ for a holomorphic matrix M. One checks that putting $\tilde{\Pi} = -A_0^{-1}B_0$, we have

$$\widetilde{H} = t^{-T} H t^T \Pi = \begin{pmatrix} A & \widetilde{B} \\ t C & \widetilde{D} \end{pmatrix},$$

where

$$\widetilde{B} = t^{-1}(B + A\widetilde{\Pi})$$
 and $\widetilde{D} = D + C\widetilde{\Pi}$

By construction, the residue of \tilde{B} is equal to $B_0 - A_0 A_0^{-1} B_0 = 0$, and hence \tilde{B} is holomorphic; therefore \tilde{H} also is. To check that $\tilde{H} \in GL_n(\mathbb{O})$, it is sufficient to check the invertibility of

$$\widetilde{H}_0 = \begin{pmatrix} A_0 & \widetilde{B}_0 \\ 0 & D_0 - C_0 A_0^{-1} B_0 \end{pmatrix}.$$

By assumption A_0 is invertible, and it is well-known that the *Schur complement* $D - CA^{-1}B$ is invertible when $\begin{pmatrix} A & B \\ C & D \end{pmatrix} \in GL_n(\mathbb{C})$ and A both are so. \Box

Note that the upper-left block of H appears unchanged in \tilde{H} . Note also that

$$\overline{H} = t^T \Pi = \begin{pmatrix} t I_d & \widetilde{\Pi} \\ 0 & I_{n-d} \end{pmatrix}$$

Geometrically, we can summarize the construction of Lemma 39 as the following scheme.



We need two technical lemmata before giving the actual proof of the permutation lemma. Let D denote the integer sequence

$$(\underbrace{d_1,\ldots,d_1}_{n_1 \text{ times}},\ldots,\underbrace{d_s,\ldots,d_s}_{n_s \text{ times}})$$

with $d_i > d_{i+1}$. We say that a matrix H is strongly D-parabolic if it has the form

$$H = \begin{pmatrix} t^{d_1} I_{n_1} & \cdots & P_{ij} \\ & t^{d_2} I_{n_2} & \vdots \\ 0 & & t^{d_s} I_{n_s} \end{pmatrix},$$

where P_{ij} is a $n_i \times n_j$ polynomial matrix in *t* satisfying deg $P_{ij} < d_i$ and $v(P_{ij}) \ge d_j$. Lemma 40. Let *H* be strongly *D*-parabolic, and let

$$H' = \begin{pmatrix} t I_m & \Pi \\ 0 & I_{n-m} \end{pmatrix},$$

where $\widetilde{\Pi}$ is a constant matrix and $m \leq n_1$. Then the product HH' is strongly D'-parabolic, where $D' = (\underbrace{d_1 + 1, \ldots, d_1 + 1}_{d_1, \ldots, d_1}, \underbrace{d_1, \ldots, d_s}_{d_1, \ldots, d_s})$.

m times
$$n_1 - m$$
 times n_s times

Proof. Let $\overline{D} = (\underbrace{d_2, \ldots, d_2}_{n_2 \text{ times}}, \ldots, \underbrace{d_s, \ldots, d_s}_{n_s \text{ times}})$. The matrix H can be written as

$$H = \begin{pmatrix} t^{d_1} I_{n_1} & P \\ 0 & \overline{H} \end{pmatrix},$$

where \overline{H} is strongly \overline{D} -parabolic, and $P = (P_2 \cdots P_s)$ where the blocks P_i satisfy deg $P_i < d_1$ and $v(P_i) \ge d_i$. Then, if $m = n_1$, the product HH' is simply

$$HH' = \begin{pmatrix} t^{d_1+1}I_{n_1} & t^{d_1}\widetilde{\Pi} + P\\ 0 & \overline{H} \end{pmatrix}.$$

Otherwise, we split the matrices in 3×3 -blocks, as

$$HH' = \begin{pmatrix} t^m I_m & 0 & P_1 \\ 0 & t^{d_1} I_{n_1 - m} & P_2 \\ 0 & 0 & \overline{H} \end{pmatrix} \begin{pmatrix} t I_m & \widetilde{\Pi}_1 & \widetilde{\Pi}_2 \\ 0 & I_{n_1 - m} & 0 \\ 0 & 0 & I_{n - n_1} \end{pmatrix}$$
$$= \begin{pmatrix} t^{d_1 + 1} I_m & t^{d_1} \widetilde{\Pi}_1 & t^{d_1} \widetilde{\Pi}_1 + P_1 \\ 0 & t^{d_1} I_{n_1 - m} & P_2 \\ 0 & 0 & \overline{H} \end{pmatrix}.$$

In both cases, we see that the product HH' is strongly D'-parabolic as requested. \Box

In the following lemma, we prove that the factorization in Proposition 38 exists if and only if the condition on the minors of the constant term holds.

Lemma 41. Let $P, Q \in GL_n(\mathbb{O})$, and let

$$D = (\underbrace{d_1, \ldots, d_1}_{n_1 \text{ times}}, \ldots, \underbrace{d_s, \ldots, d_s}_{n_s \text{ times}})$$

with $d_i > d_{i+1}$. Assume that $\Pi = z^{-D} P z^D Q \in \mathfrak{gl}_n(\mathbb{C}[[1/z]])$. Decomposing a matrix M in blocks $M_{i,j}$ according to the multiplicities (n_1, \ldots, n_s) of D, let

$$P_{i} = \begin{pmatrix} P_{i,i} & \cdots & P_{i,s} \\ \vdots & \ddots & \vdots \\ P_{s,i} & \cdots & P_{s,s} \end{pmatrix} \quad and \quad \widetilde{\Pi}_{i} = \begin{pmatrix} \Pi_{i,1}^{(0)} & \cdots & \Pi_{i,s}^{(0)} \\ \vdots & \ddots & \vdots \\ \Pi_{s,1}^{(0)} & \cdots & \Pi_{s,s}^{(0)} \end{pmatrix}.$$

Then $\tilde{\Pi}_i$ has maximal rank if and only if $P_i(0), \ldots, P_s(0)$ are all invertible.

Proof. We prove this result by induction. We put $M = \sum_{\ell \in \mathbb{Z}} M^{(\ell)} z^{\ell}$ for a formal series matrix, with $M^{(\ell)} = 0$ when needed. Let $M_{i,(j:k)} = (M_{i,j} \cdots M_{i,k})$ for $j \leq k$, and $M_{i,\bullet} = M_{i,(1:s)}$, and put further

$$\widetilde{M}_{j} = \begin{pmatrix} M_{j,\bullet}^{(0)} \\ M_{j+1,\bullet}^{(0)} \\ \vdots \\ M_{s,\bullet}^{(0)} \end{pmatrix} \quad \text{and} \quad M_{j}^{\prime(\ell)} = \begin{pmatrix} M_{j,\bullet}^{(\ell)} \\ M_{j+1,\bullet}^{(\ell+d_{j}-d_{j+1})} \\ \vdots \\ M_{s,\bullet}^{(\ell+d_{j}-d_{s})} \end{pmatrix},$$

so that

$$M_{j}^{\prime(\ell)} = \begin{pmatrix} M_{j,\bullet}^{(\ell)} \\ M_{j+1}^{\prime(\ell+d_{j}-d_{j+1})} \end{pmatrix}.$$

Let $\check{P} = z^{-D} P z^{D}$. By construction, we have $\check{P} = (\check{P}_{i,j})$ where

$$\check{P}_{i,j} = P_{i,j} z^{d_j - d_i} = \sum_{\ell \ge 0} P_{i,j}^{(\ell)} z^{\ell + d_j - d_i} = \sum_{\ell \ge d_j - d_i} P_{i,j}^{(\ell + d_i - d_j)} z^{\ell}.$$

Since $\Pi = \check{P}Q$ holds, we have $\Pi_{i,\bullet} = \sum_{i=1}^{s} \check{P}_{i,k}Q_{k,\bullet}$, with $Q_{k,\bullet} = \sum_{\ell \ge 0} Q_{k,\bullet}^{(\ell)} z^{\ell}$. After some algebra, we get

(12)
$$\Pi_{i,\bullet} = \sum_{t \leq 0} \left(\sum_{k=1}^{s} \sum_{\ell=d_k-d_i}^{t} P_{i,k}^{(\ell+d_i-d_k)} Q_{k,\bullet}^{(t-\ell)} \right) z^t \quad \text{since } \Pi \in \mathfrak{gl}_n(\mathbb{C}\llbracket 1/z \rrbracket).$$

We will establish the main claim (MC) of the lemma by proving simultaneously the following two additional results. If the assumption of the lemma holds, then we have

- (A) \tilde{Q}_j is in the row span of $\tilde{\Pi}_j$ for $i \leq j \leq s$,
- (B) for $i \leq j \leq s$ and $\ell \leq d_{j-1} d_j 1$, there exists a matrix $X_{j\ell}$ such that $Q'_i^{(\ell)} = X_{j\ell} \tilde{Q}_j$.

Assume first that i = s. According to formula (12), we have

$$\Pi_{s,\bullet} = \sum_{t\leq 0} \left(\sum_{k=1}^{s} \sum_{\ell=d_k-d_s}^{\ell} P_{s,k}^{(\ell+d_s-d_k)} Q_{k,\bullet}^{(t-\ell)} \right) z^t$$

Since $d_k - d_s > 0$ for $k \neq s$, we have that $\prod_{s,\bullet}$ is a constant, and the formula reduces to t = 0 and k = s; thus $\prod_{s,\bullet}^{(0)} = P_{s,s}^{(0)} Q_{s,\bullet}^{(0)}$. Since rk $Q_{s,\bullet}^{(0)} = n_s$, we get

$$\operatorname{rk} \Pi_{s,\bullet}^{(0)} = n_s \iff \det P_{s,s}^{(0)} \neq 0,$$

which establishes (MC) and (A) for i = s. For $0 < t < d_{s-1} - d_s$, formula (12) reduces to

$$\Pi_{s,\bullet}^{(t)} = \sum_{\ell=0}^{t} P_{s,s}^{(\ell)} Q_{s,\bullet}^{(t-\ell)} = 0$$

since $\Pi^{(t)} = 0$ for t > 0. We can arrange all these equations into the large matrix equation

$$\begin{pmatrix} P_{s,s}^{(0)} & 0 & \cdots & 0 \\ P_{s,s}^{(1)} & P_{s,s}^{(0)} & \vdots \\ \vdots & \ddots & \vdots \\ P_{s,s}^{(d_{s-1}-d_s-1)} & \cdots & \cdots & P_{s,s}^{(0)} \end{pmatrix} \begin{pmatrix} Q_{s,\bullet}^{(0)} \\ Q_{s,\bullet}^{(1)} \\ \vdots \\ Q_{s,\bullet}^{(d_{s-1}-d_s-1)} \end{pmatrix} = \begin{pmatrix} \Pi_{s,\bullet}^{(0)} \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

It follows easily that $Q_{s,\bullet}^{(\ell)}$ is then a left multiple of $Q_{s,\bullet}^{(0)}$ for $0 \le \ell \le d_{s-1} - d_s - 1$, and also for $\ell < 0$, since $Q_{s,\bullet}^{(\ell)} = 0$, which establishes the claim (B) for i = s.

Assume now that (MC), (A) and (B) hold for $i + 1 \le j \le s$ and $\ell \le d_{j-1} - d_j - 1$. Formula (12) gives

$$\Pi_{j,\bullet}^{(t)} = \sum_{k=1}^{s} \sum_{\ell=d_k-d_j}^{t} P_{j,k}^{(\ell+d_j-d_k)} Q_{k,\bullet}^{(t-\ell)}$$
$$= \sum_{\ell \ge 0} \sum_{k=1}^{s} P_{j,k}^{(\ell)} Q_{k,\bullet}^{(t+d_j-d_k-\ell)},$$

with the convention that $Q^{(\ell)} = 0$ when $\ell < 0$. Considering exponents $t = d_i - d_j + m$ for $i \leq j \leq s$ and $0 \leq m \leq d_{i-1} - d_i - 1$, we get

$$\Pi_{j,\bullet}^{(d_i-d_j+m)} = \sum_{\ell \ge 0} \sum_{k=i}^{s} P_{j,k}^{(\ell)} \mathcal{Q}_{k,\bullet}^{(m+d_i-d_k-\ell)}$$
$$= \sum_{\ell \ge 0} P_{j,(i:s)}^{(\ell)} \mathcal{Q}_i^{\prime(m-\ell)},$$

since k < i implies $m + d_i - d_k < 0$. Therefore, we get

$$\Pi_{j,\bullet}^{(d_{i}-d_{j}+m)} = \sum_{\ell \ge 0} P_{j,i}^{(\ell)} Q_{i,\bullet}^{(m-\ell)} + P_{j,(i+1:s)}^{(\ell)} Q_{i+1}^{'(m+d_{i}-d_{i+1}-\ell)}$$

$$= \sum_{\ell=0}^{m} \left(P_{j,i}^{(\ell)} Q_{i,\bullet}^{(m-\ell)} + P_{j,(i+1:s)}^{(\ell)} Q_{i+1}^{'(m+d_{i}-d_{i+1}-\ell)} \right) + \sum_{\ell \ge m} P_{j,(i+1:s)}^{(\ell)} Q_{i+1}^{'(m+d_{i}-d_{i+1}-\ell)}$$

$$= \sum_{\ell=0}^{m} P_{j,(i:s)}^{(\ell)} Q_{i}^{'(m-\ell)} + \sum_{\ell \ge m} P_{j,(i+1:s)}^{(\ell)} Q_{i+1}^{'(m+d_{i}-d_{i+1}-\ell)}.$$

According to the induction $Q'^{(m+d_i-d_{i+1}-\ell)}_{i+1}$ is a left multiple of \tilde{Q}_{i+1} for $\ell > m$, so we have

(13)
$$\Pi_{j,\bullet}^{(d_i-d_j+m)} = \sum_{\ell=0}^{m} \left(P_{j,i}^{(\ell)} \cdots P_{j,s}^{(\ell)} \right) Q_i^{\prime(m-\ell)} + [\tilde{Q}_{i+1}] \quad \text{for } i \leq j \leq s,$$

where we let, for notational simplicity, A = B + [Q] mean "there exists a matrix X such that A - B = XQ". Assume m = 0 first. The equation for j = i is then

(14)
$$\Pi_{i,\bullet}^{(0)} = P_{i,i}^{(0)} Q_{i,\bullet}^{(0)} + \left(P_{i,i+1}^{(0)} \cdots P_{i,s}^{(0)} \right) Q_{i+1}^{\prime (d_i - d_{i+1})} + [\tilde{Q}_{i+1}].$$

The remaining equations,

$$\Pi_{j,\bullet}^{(d_i-d_j)} = P_{j,i}^{(0)} Q_{i,\bullet}^{(0)} + \left(P_{j,i+1}^{(0)} \cdots P_{j,s}^{(0)}\right) Q_{i+1}^{\prime(d_i-d_{i+1})} + [\tilde{Q}_{i+1}] = 0, \quad i+1 \le j \le s,$$

can be rewritten as

$$P_{i+1}^{(0)}Q_{i+1}^{\prime(d_i-d_{i+1})} = -\begin{pmatrix} P_{i+1,i}^{(0)} \\ \vdots \\ P_{s,i}^{(0)} \end{pmatrix} Q_{i,\bullet}^{(0)} + [\tilde{Q}_{i+1}], \text{ where } P_{i+1}^{(0)} = \begin{pmatrix} P_{i+1,i+1}^{(0)} \cdots P_{i+1,s}^{(0)} \\ \vdots & \ddots & \vdots \\ P_{s,i+1}^{(0)} \cdots & P_{s,s}^{(0)} \end{pmatrix}$$

is invertible by the induction assumption. Put

$$B = \begin{pmatrix} P_{i,i+1}^{(0)} & \cdots & P_{i,s}^{(0)} \end{pmatrix} \text{ and } C = \begin{pmatrix} P_{i+1,i}^{(0)} \\ \vdots \\ P_{s,i}^{(0)} \end{pmatrix} \text{ so that } P_i^{(0)} = \begin{pmatrix} P_{i,i}^{(0)} & B \\ C & P_{i+1}^{(0)} \end{pmatrix}.$$

Thus, we get

$$Q_{i+1}^{\prime(d_i-d_{i+1})} = -(P_{i+1}^{(0)})^{-1}CQ_{i,\bullet}^{(0)} + X\tilde{Q}_{i+1} = X^{\prime}\tilde{Q}_i;$$

hence $Q_i^{\prime(0)}$ is a left multiple of \tilde{Q}_i . Moreover, substituting in (14), we get

(15)
$$\Pi_{i,\bullet}^{(0)} = \left(P_{i,i}^{(0)} - B\left(P_{i+1}^{(0)}\right)^{-1}C\right)Q_{i,\bullet}^{(0)} + [\tilde{Q}_{i+1}]$$

By assumption \tilde{Q}_{i+1} is in the row space of

$$\begin{pmatrix} \Pi_{i+1,\bullet}^{(0)} \\ \vdots \\ \Pi_{s,\bullet}^{(0)} \end{pmatrix}.$$

Hence $\operatorname{rk} \Pi_{i,\bullet}^{(0)} = n_i$ if and only if $\operatorname{rk} (P_{i,i}^{(0)} - B(P_{i+1}^{(0)})^{-1}C)Q_{i,\bullet}^{(0)} = n_i$; that is, $\operatorname{det} (P_{i,i}^{(0)} - B(P_{i+1}^{(0)})^{-1}C) \neq 0$. This matrix is the Schur complement of

$$P_i^{(0)} = \begin{pmatrix} P_{i,i}^{(0)} & B \\ C & P_{i+1}^{(0)} \end{pmatrix},$$

which is invertible exactly when $P_{i+1}^{(0)}$ is. Therefore (MC) is established in general, and by (15), \tilde{Q}_i is in the row span of $\tilde{\Pi}_i$, so claim (A) is proved. Similarly, for a given $0 \le m \le d_{i-1} - d_i - 1$, we can stack the remaining equations, (13), corresponding to $i \le j \le s$,

$$\Pi_{j,\bullet}^{(d_i-d_j+m)} = \sum_{\ell=0}^m P_{j,(i:s)}^{(\ell)} Q_i^{\prime(m-\ell)} + [\tilde{Q}_{i+1}] = 0,$$

to get the relation

$$\begin{pmatrix} \Pi_{i,\bullet}^{(d_i-d_j+m)} \\ \vdots \\ \Pi_{s,\bullet}^{(d_i-d_j+m)} \end{pmatrix} = \sum_{\ell=0}^m P_i^{(\ell)} Q_i^{\prime(m-\ell)} = \begin{cases} \widetilde{\Pi}_i & \text{if } m = 0, \\ 0 & \text{otherwise} \end{cases}$$

Putting for notational simplicity $d = d_{i-1} - d_i$, we finally get

$$\begin{pmatrix} P_i^{(0)} & 0 & \cdots & 0 \\ P_i^{(1)} & P_i^{(0)} & \vdots \\ \vdots & & \ddots & \vdots \\ P_s^{(d-1)} & \cdots & \cdots & P_i^{(0)} \end{pmatrix} \begin{pmatrix} Q_i^{\prime(0)} \\ Q_i^{\prime(1)} \\ \vdots \\ Q_i^{\prime(d-1)} \end{pmatrix} = \begin{pmatrix} \Pi_i' \\ 0 \\ \vdots \\ 0 \end{pmatrix} + [\tilde{Q}_{i+1}] \quad \text{where } \Pi_i' = \begin{pmatrix} \tilde{\Pi}_i \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

Accordingly, $Q_i^{\prime(\ell)}$ is a left multiple of $Q_i^{\prime(0)}$ for $0 \leq \ell \leq k-1$. Since we have established that $Q_i^{\prime(0)}$ is a left multiple of \tilde{Q}_i , this ends the proof of claim (B). \Box

Proof of Proposition 38. Assume for simplicity that $D = \text{diag}(d_1 I_{n_1}, \ldots, d_s I_{n_s})$ is written by blocks, and that $d_1 > d_2 > \cdots > d_s$. Then there exist $m = d_1 - d_s$ matrices T_1, \ldots, T_m of the type

$$T_i = \begin{pmatrix} I_{b_i} & 0\\ 0 & 0_{n-b_i} \end{pmatrix},$$

where every b_i is equal to some $n_1 + \cdots + n_{t_i}$ for some decreasing sequence t_i , such that $D = T_1 + \cdots + T_m$. Secondly, assume that all left-upper square blocks of P(0) of sizes b_i are invertible. Letting $P = H_0$, according to Lemma 39, there exists a sequence of monopole matrices

$$\Pi_i = \begin{pmatrix} I_{b_i} & t^{-1} \widetilde{\Pi}_i \\ 0 & I_{n-b_i} \end{pmatrix}$$

with a constant matrix $\tilde{\Pi}_i$, and a sequence of lattice gauge transforms $H_i \in GL_n(\mathbb{O})$ such that

(16)
$$H_{i+1} = t^{-T_i} H_i t^{T_i} \Pi_i.$$

Let

$$\overline{H}_i = t^{T_i} \Pi_i = \begin{pmatrix} t I_{b_i} & \overline{\Pi}_i \\ 0 & I_{n-b_i} \end{pmatrix}.$$

It follows from Lemma 40 that $\overline{H} = \overline{H}_1 \cdots \overline{H}_m$ is strongly *D*-parabolic. It follows then, as a remarkable consequence, that the diagonal matrix t^D can be both factored from the matrix \overline{H} on the left as $\overline{H} = t^D \Pi$ with a monopole matrix Π , and simultaneously from the right as $\overline{H} = Qt^D$ with a lattice gauge $Q \in GL_n(\mathbb{O})$. Since $t^D H_{d+1} = H\overline{H}$ holds, we get on the one hand that $t^{-D} P^{-1} t^D H_{m+1} = \Pi \in$ $GL_n(\mathbb{C}[t^{-1}])$ as required for the first claim of the permutation lemma. However, and this was not stated in [Bolibrukh 1990] or [Ilyashenko and Yakovenko 2007], we also have the relation $t^D \Pi = Qt^D$, which yields the second claim.

For the third claim, the sufficiency of the minors condition has just been established. For the converse, assume that we have the following diagram, which we have completed with a monopole Π and a lattice gauge Q such that $z^{-D} P z^{D} Q = \Pi$:

$$\begin{array}{c} \Lambda:(e) \xrightarrow{z^{D}} Y_{M} \subset M \\ & \downarrow^{P} \\ \Lambda:(\varepsilon) \xrightarrow{z^{D}} \tilde{M} \xrightarrow{Q} Y_{\tilde{M}} \subset \tilde{M} \end{array}$$

Let $Q = Q_0 + Q_1 z + \cdots \in \operatorname{GL}_n(\mathbb{O})$, and $\Pi = \Pi_0 + \cdots + \Pi_t z^{-t} \in \operatorname{GL}_n(\mathbb{C}[z^{-1}])$. Recall that if Π is a monopole, then one must in particular have $\Pi_0 \in \operatorname{GL}_n(\mathbb{C})$. Then Lemma 41 implies that such a factorization only exists if the matrix P_0 satisfies the condition on the minors.

The following scheme sums up this construction.

$$\begin{array}{c|c}
\tilde{M}_{m-1} \xrightarrow{H_{m-1}} \tilde{M}_{m-1} \xrightarrow{t^{T_m}} Y_{m-1} \subset M_m^{m-1} \\
\downarrow^{T_m} & \overline{H}_m & \Pi_m \\
\tilde{M}_m \xrightarrow{H_m} \tilde{Y} \subset \tilde{M}_m
\end{array}$$

The first row corresponds to a min-shortest path $\Gamma = (\Lambda, M_1, \ldots, M)$ from Λ to a given BG trivialization M. This path Γ is included in an apartment \mathfrak{B} , namely the one spanned by a BG basis (e) of Λ corresponding to the trivialization M. By definition, the apartment \mathfrak{B} goes through the global form Y of M. The gauge H^{-1} does not map the shortest path Γ onto anything special. However, if we call $\mathcal{A} = H^{-1}(\mathfrak{B})$ the image of the apartment spanned by (e), the permutation lemma tells us how to construct a shortest path Γ' in \mathcal{A} whose end point is also a BG trivialization of Λ . Lemma 39 gives the step-by-step modification of the shortest path Γ . Row *i* of the diagram corresponds indeed to a partial shortest path $\Gamma_i =$ $(\tilde{M}_i, M_i^i, \ldots, M_m^i)$ whose end-point is a BG trivialization of the *i*-th element \tilde{M}_i of the shortest path $\Gamma' = (\Lambda, \tilde{M}_1, \ldots, \tilde{M}_m)$. Even if the end-point \tilde{M}_m is a BG trivialization of Λ , note that the apartment \mathcal{A} does not contain the global form \tilde{Y} of \tilde{M}_m , and that we still need the gauge transform H_m to obtain it.

3.3.1. Consequences of the permutation lemma. As stated in [Ilyashenko and Yakovenko 2007], one can assume that $\tau = id$ if all leading principal minors of *P* are holomorphically invertible; that is, the corresponding minors of *P*(0) are nonzero. This condition can always be ensured by a permutation of the columns of *P*. Actually, as stated in Proposition 38(3), it is sufficient that this condition

holds only for the leading principal minors of orders $n_1, n_1 + n_2, ..., n_1 + \cdots + n_{s-1}$ of P(0). We then say that *P* respects the minors condition with respect to *D* (or to $\sigma = (n_1, ..., n_s)$). Recall the following well-known result.

Lemma 42. Let $P \in GL_n(\mathbb{C})$, let (e) be the standard basis of \mathbb{C}^n and (ε) be the column vectors of P. Then P respects the minors condition with respect to a signature σ if and only if (e) is transversal to the flag $F^{\sigma}(\varepsilon)$.

Definition 43. Let \mathcal{A} be an apartment induced by the *K*-frame Φ in *V*. Consider $\Lambda \in \mathcal{A}$. Let *H* be the Harder–Narasimhan flag of Λ in $E = \Lambda/\mathfrak{m}\Lambda$. Let $W \subset S_n$ be the parabolic subgroup of S_n associated to *H*, and *W'* be the set of right cosets $W \setminus S_n$. Let (*e*) be a basis of Λ in the frame Φ whose image (\overline{e}) in *E* is transversal to *H*. The integer

$$\iota_{\Lambda}(\mathcal{A}) = |\{\tau \in W' \mid \tau(\bar{e}) \text{ transversal to } H\}|$$

is independent of (e) and is called the *transversality index* of \mathcal{A} with respect to Λ .

Theorem 44. Let \mathscr{C} be a holomorphic vector bundle over X, and let $\Lambda = \mathscr{C}_x \in \Lambda$ be its stalk at $x \in X$.

- (i) For any apartment A in the Bruhat–Tits building B at x such that [Λ] ∈ A, there exists a BG trivialization of Λ in A.
- (ii) More precisely, the number of BG trivializations of Λ in A is exactly

$$|\mathcal{A} \cap BG(\Lambda)| = \iota_{\Lambda}(\mathcal{A}).$$

Proof. Let (*e*) be a BG basis of Λ , and $M \in BG(\Lambda)$. Let (ε) be a basis of the lattice Λ which spans the apartment \mathcal{A} . Since \mathcal{A} is invariant under S_n , we can assume that the matrix $P \in GL_n(\mathbb{O})$ of the basis change from (ε) to (*e*) has invertible principal leading minors. According to the permutation lemma, there exists a matrix $\tilde{P} \in GL_n(\mathbb{O})$ such that

$$\Pi = z^{-D} P^{-1} z^D \tilde{P} \in \operatorname{GL}_n(\mathbb{C}[z^{-1}]).$$

The gauge Π sends the basis of global sections $(\sigma) = (z^D e)$ of the BG trivialization of \mathscr{C} , given at x by M, into a basis (\tilde{e}) of \tilde{M} . Since Π is a monopole, the basis (\tilde{e}) is also a global basis of sections, but spans another trivializing bundle, namely $\mathscr{F} = \mathscr{C}^{\tilde{M}}$. Therefore the arbitrary apartment \mathscr{A} spanned by (ε) indeed contains a trivial bundle. Now the matrix $\overline{H} = z^D \Pi$ admits a right factorization $\overline{H} = Qz^D$. As a consequence, if we let $(\tilde{\varepsilon})$ be the basis of Λ obtained from (e) by the matrix Q, then $z^D(\tilde{\varepsilon})$ is also a basis of $Y_{\tilde{M}}$. The following scheme sums up the situation.


Therefore the lattice M is also a BG trivialization of Λ .

The monopole gauge Π is block-upper-triangular according to *D*, and its block matrices Π_{ij} satisfy

$$d_j - d_i \leq v(\Pi_{ij}) \leq \deg \Pi_{ij} \leq 0.$$

This means that Π respects the Harder–Narasimhan filtration HN_{Λ} of V corresponding to the lattice Λ . Conversely, the lattice gauge Q has a lower D-block-triangular constant term Q_0 . Since $z^{-D} P Q z^D \in \text{GL}_n(\mathbb{O})$, the matrix PQ has an upper D-block-triangular constant term $P_0 Q_0$. Therefore, if we put P_D , P_D^- for the pair of opposite D-parabolic standard subgroups of $\text{GL}_n(\mathbb{C})$, the matrix P_0 satisfies $P_0 \in P_D^- P_D$. This means exactly that P satisfies the minors condition with respect to D. If we permute the vectors of (ε) with $\tau \in S_n$ in such a way that $D_{\tau^{-1}} \neq D$, the permutation lemma ensures that $\mathcal{L}(z^D(\varepsilon_{\tau})) \neq \tilde{M}$ is again a BG trivialization of Λ . This establishes the second claim of the theorem. \Box

Corollary 45. For any apartment $\mathcal{A} \ni \Lambda$, there exists an ordered basis (ε) of \mathcal{A} and a BG basis (e) of Λ such that the gauge P from (e) to (ε) has a lower block-D-triangular unipotent constant term P_0 , and that the following picture holds.

$$\begin{array}{c} \Lambda : (\varepsilon) \xrightarrow{z^{D}} M \\ P = P_{0} + z \widetilde{U} \uparrow & \uparrow P_{I} = I + z U \\ \Lambda : (e) \xrightarrow{z^{D}} Y_{M} \end{array}$$

The number of BG trivializations in \mathcal{A} can hence be computed from a matrix $P_0 \in GL_n(\mathbb{C})$ with a simple structure:

$$P_0 = \begin{pmatrix} I_{n_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ X_{ij} & \cdots & I_{n_s} \end{pmatrix}$$

as the number of permutations not leaving D invariant whose action on the columns of P_0 preserve the minors condition. The following result completes the proof of Theorem 25.

Corollary 46. Let $\Lambda \in \Lambda$ and $T = T(\Lambda)$.

- (i) For any basis (ε) of Λ , let $\tilde{M} = \mathcal{L}(z^T(\varepsilon))$. Then either $\tilde{M} \in BG(\Lambda)$ or $\tilde{M} \notin \Lambda^0$.
- (ii) For any flag \mathbf{F}' in $\Lambda/\mathfrak{m}\Lambda$ transversal to $\operatorname{HN}_{\Lambda}$, and any form Y in Λ , the lattice $\tilde{M} = \mathscr{L}_{Y}(\mathbf{F}', T)$ is a BG trivialization of Λ .

In particular, if $M \in BG(\Lambda)$, then $\tilde{M} = \mathscr{L}_Y(F^{\Lambda}(M), T) \in BG(\Lambda)$ for any form Y in Λ .

Proof. For (i), if \tilde{M} is trivial, then, by Lemma 41, the gauge from (e) to (ε) satisfies the minors condition; hence \tilde{M} is BG. According to Lemma 27, there exists $M \in BG(\Lambda)$ such that $F' = F^{\Lambda}(M)$. Let (e) be a Smith basis of Λ for M. The lattice \tilde{M} is spanned by $z^{T}(e_{Y})$ where (e_{Y}) is the Y-basis of (e). The gauge P from (e) to (e_{Y}) has invertible principal minors, since it is tangent to I. Item (ii) follows.

The permutation lemma is a sort of converse to the Birkhoff–Grothendieck theorem. It can also be seen as a lifting and factorization lemma. For a standard parabolic subgroup $P \subset GL_n(\mathbb{C})$, let

$$\mathbb{Z}_P^n = \{ D \in \mathbb{Z}^n \mid d_i < d_j \Rightarrow A_{ij} = 0 \text{ for all } A \in P \}.$$

Corollary 47. Let $A \in GL_n(\mathbb{C})$. Let P, P^- be a pair of opposed standard parabolic subgroups of $GL_n(\mathbb{C})$. Then we have, for $w \in S_n$,

$$A_0 \in PwP^- \Leftrightarrow \begin{cases} \forall D \in \mathbb{Z}_P^n, \ \exists \Pi \in \operatorname{GL}_n(\mathbb{C}[z^{-1}]), \ \exists B \in \operatorname{GL}_n(\mathbb{O}) \\ such \ that \ A = z^D \Pi B z^{-wD}. \end{cases}$$

Moreover, if this holds, then $B \in \mathcal{G}_{-D}$ and $\Pi \in \hat{P}$ hold, where \hat{P} stands here for the subgroup of upper-D-triangular matrices of $\operatorname{GL}_n(\mathbb{C}[z^{-1}])$.

This result means that there is a cell decomposition of $GL_n(\mathbb{O})$ whose cells (defined by the right-hand side of the previous relation) are mapped surjectively (by the canonical surjection $GL_n(\mathbb{O}) \rightarrow GL_n(\mathbb{C})$) on the Schubert–Bruhat cells of $GL_n(\mathbb{C})$. Finally, a last consequence of Proposition 38 is that the lattices involved in the Birkhoff–Grothendieck algorithm given in Section 3.2.1 can be taken inside a single apartment.

Corollary 48. Let $\Lambda \in \Lambda$. For any $M \in \Lambda^0$, let $\Gamma(M, \Lambda) = (\Lambda_0 = M, \dots, \Lambda_t = \Lambda)$ be the min-shortest path from M to Λ . For any apartment $\mathcal{A} \ni \Lambda$, M, there exist $\tilde{M}_1, \dots, \tilde{M}_t \in \mathcal{A} \cap \Lambda^0$ such that $\tilde{M}_i \in BG(\Lambda_i)$ for $1 \leq i \leq t$.

Proof. Since $\mathcal{A} \ni \Lambda$, M, we have $\Lambda_i \in \mathcal{A}$ for $1 \le i \le t$. By Theorem 44, any lattice in \mathcal{A} admits a BG trivialization in \mathcal{A} .

4. Local meromorphic connections

Let $\text{Der}_{\mathbb{C}}(K)$ be the *K*-vector space of dimension one of \mathbb{C} -derivations of *K* and $\Omega = \Omega^{1}_{\mathbb{C}}(K)$ the dual composed of differentials of *K*. The valuation *v* extends naturally to these spaces by the formulæ $v(\vartheta) = v(f)$ and $v(\omega) = v(g)$ if $\vartheta = fd/dz$ and $\omega = g dz$ for any uniformizing parameter *z* of *K*. The space Ω is naturally filtered by the rank-one free \mathbb{O} -modules $\Omega(k) = \{\omega \in \Omega \mid v(\omega) \ge -k\}$.

Let *V* be a *K*-vector space of finite dimension *n* and let $\Omega(V) = V \otimes_K \Omega^1_{\mathbb{C}}(K)$. We fix a *meromorphic connection* ∇ on *V*. This is an additive map $\nabla : V \to \Omega(V)$ satisfying the Leibniz rule

$$\nabla(fv) = v \otimes df + f \nabla v$$
 for all $f \in K$ and all $v \in V$.

For any basis $(e) = (e_1, \ldots, e_n)$ of V, the matrix $Mat(\nabla, (e))$ of the connection ∇ in the basis (e) is the matrix $A = (A_{ij}) \in M_n(\Omega)$ such that

$$\nabla e_j = -\sum_{i=1}^n e_i \otimes A_{ij}$$
 for all $j = 1, \dots, n$.

If the matrix $P = \text{Mat}(\text{id}_V, (\varepsilon), (e)) \in \text{GL}_n(K)$ is the basis change from (e) to any other basis (ε), then the matrix of ∇ in (ε) is given by the *gauge transform* of A:

(17)
$$A_{[P]} = P^{-1}AP - P^{-1}dP.$$

For any derivation $\tau \in \text{Der}(K/\mathbb{C})$, the contraction of ∇ with τ induces a differential operator ∇_{θ} on V. The connection ∇ is *regular* whenever the set of *logarithmic* lattices

$$\mathbf{\Lambda}_{\log} = \{ \Lambda \in \mathbf{\Lambda} \mid \nabla(\Lambda) \subset \Lambda \otimes_{\mathbb{O}} \Omega(1) \}$$

is nonempty. For any logarithmic lattice $\Lambda \in \Lambda_{\log}$, the connection ∇ induces a well-defined residue endomorphism $\operatorname{Res}_{\Lambda} \nabla \in \operatorname{End}_{\mathbb{C}}(\Lambda/\mathfrak{m}\Lambda)$. Note that, since the set Λ_{\log} is closed under homothety and module sums and intersections [Corel 2004, Lemma 2.5], it induces a path-convex subset of the Bruhat–Tits building: if $L, L' \in \Lambda_{\log}$, then every shortest path between L and L' is a subset of Λ_{\log} . This applies in particular to both $\Gamma_{\max}(L, L')$ and $\Gamma_{\min}(L, L')$.

4.1. *The Deligne lattice.* As is well known, the choice of a matrix logarithm of the monodromy corresponds to fixing a special lattice in the space V. More precisely, let $V^{\nabla} \subset V \otimes_K H$ be the \mathbb{C} -vector space of horizontal sections on any Picard–Vessiot extension H of K. Let $g = g_s g_u \in \text{End}(V^{\nabla})$ be the multiplicative Jordan decomposition of the corresponding local monodromy map. Then the logarithm of the unipotent part g_u is canonically defined (by the Taylor expansion formula for $\log(1 + x)$), but there are several ways to define the logarithm of the semisimple

part g_s . Namely, one must fix a branch of the complex logarithm for every *distinct* eigenvalue of g_s .

A well known result (variously attributed to Deligne, Manin, ...) says that this choice uniquely defines a lattice in V. In Deligne's terms, for any section σ of $\mathbb{C} \to \mathbb{C}/\mathbb{Z}$, there is a unique logarithmic lattice Δ_{σ} such that the eigenvalues of the residue map $\operatorname{Res}_{\Delta_{\sigma}} \nabla$ are in the image $\operatorname{Im} \sigma$ of σ . As a habit, one usually takes $\operatorname{Re}(\operatorname{Im} \sigma) \subset [0, 1[$. In fact, such a habit is not as arbitrary as it seems.

Proposition 49. Assume that the connection ∇ admits an apparent singularity (i.e., the monodromy map is trivial). Then the matrix $Mat(\nabla, (e))$ is holomorphic if and only if the lattice spanned by (e) is equal to the Deligne lattice Δ attached to $Re(Im \sigma) \subset [0, 1]$.

Proof. Since the monodromy map is trivial, its normalized logarithm with respect to Δ is 0. Hence, there is a basis of Δ where the connection has matrix 0. In any other basis (e) of Δ , the connection has matrix $A = P^{-1}dP$ with $P \in GL_n(\mathbb{O})$, which is holomorphic. Let M be another lattice, and let (e) be a Smith basis of Δ for M. Then the matrix in a basis of M is given by the gauge equation

$$\widetilde{A} = z^{-D}Az^D - z^{-D}d(z^D) = (A_{ij}z^{k_j-k_i}) - D\frac{dz}{z}.$$

The nonzero diagonal terms of the matrix D of elementary divisors of M give necessarily rise to a pole of order 1 in \tilde{A} . Therefore, Δ is the only lattice where the connection has a holomorphic matrix.

As a result, we will call Δ *the* Deligne lattice of V.

4.1.1. *Birkhoff forms.* According to a very classical result (see, e.g., [Gantmacher 1959, p. 150]), if

$$\Omega = \operatorname{Mat}(\nabla, (e)) = \sum_{k \ge 0} A_k z^k \frac{dz}{z}$$

is the series expansion in z of the matrix of ∇ in a basis (e) of Δ , the gauge $P = \sum_{k \ge 0} P_k z^k \in GL_n(\mathbb{O})$ defined recursively by

(18)
$$\begin{cases} P_0 = I, \\ P_k = \Phi_{A_0, A_0 - kI}^{-1}(Q_k), & \text{where } Q_k = \sum_{i=1}^k A_i P_{k-i}, \end{cases}$$

transforms Ω into $A_0 dz/z$. Here we put $\Phi_{U,V}(X) = XU - VX$. Recall that the map $\Phi_{U,V}$ is an automorphism of $\mathfrak{gl}(\mathbb{C})$ when the spectra of U and V are disjoint. The gauge P thus defined is uniquely determined; moreover, the set of bases where ∇ has matrix $L\frac{dz}{z}$ where $L \in M_n(\mathbb{C})$ is a constant matrix spans a form Υ_z of Δ , that we call the *Birkhoff form* of the Deligne lattice Δ . The gauge transform P sends in fact the basis (e) to its Υ_z -basis, that we denote here for simplicity by (e_z) .

As it results from the proof of Proposition 49, when the singularity is apparent, the Birkhoff form is uniquely defined. Otherwise, however, the form Υ_z depends on the choice of the local coordinate z. Two Birkhoff forms are nevertheless canonically isomorphic.

Lemma 50. Let z, t be two local coordinates, and let $\alpha \in \mathbb{O}^*$ such that $z = \alpha t$. Let P_z and P_t be the gauge transforms that send (e) to (e_z) and (e_t) , respectively. There is a unique gauge transform \tilde{P} that sends (e_z) to (e_t) .

Proof. One has $\frac{dz}{z} = u\frac{dt}{t}$ with $u = 1 + \frac{\theta_t \alpha}{\alpha}$ where $\theta_t = t\frac{d}{dt}$. Put $u = \sum_{i=0}^{\infty} u_i t^i$. Accordingly, the matrix of the connection in (e_z) satisfies

$$\operatorname{Mat}(\nabla, (e_z)) = A_0 \frac{dz}{z} = A_0 \left(\sum_{i=0}^{\infty} u_i t^i \right) \frac{dt}{t}.$$

There exists therefore a uniquely defined gauge transform $\tilde{P} = \sum_{i=0}^{\infty} \tilde{P}_i t^i$ that transforms the expression $A_0 dz/z$ into $A_0 dt/t$, as explained in the following scheme.



The matrix series \tilde{P} is determined recursively by the equations (18) applied to the series $\sum_{i=0}^{\infty} A_0 u_i t^i$. The coefficients \tilde{P}_i are even *polynomials in* A_0 , defined by the induction rule

$$\widetilde{P}_0 = I, \quad \widetilde{P}_k = \frac{1}{k} \sum_{i=1}^k u_i A_0 \widetilde{P}_{k-i}.$$

4.2. Logarithmic lattices and stable flags. When two lattices Λ , M are adjacent, all the relevant information on M can be retrieved from the quotient $M/\mathfrak{m}\Lambda$. This is also true in presence of a connection.

Lemma 51. Let $\Lambda \in \Lambda_{\log}$ be a logarithmic lattice. For any adjacent lattice $M \in [\mathfrak{m}\Lambda, \Lambda]$, we have $M \in \Lambda_{\log}$ if and only if $M/\mathfrak{m}\Lambda$ is $\operatorname{Res}_{\Lambda}\nabla$ -stable.

Proof. In any basis (e) of Λ such that the images of the first $m = \dim W$ vectors span $W = M/\mathfrak{m}\Lambda$, the connection matrix $\Omega = \operatorname{Mat}(\nabla, (e))$ has a residue of the form $\begin{pmatrix} A & B \\ 0 & C \end{pmatrix} \in \operatorname{M}_n(\mathbb{C})$, where $A \in \operatorname{M}_m(\mathbb{C})$. Putting $T = \operatorname{diag}(0_m, I_{n-m})$, the basis $(\varepsilon) = z^T(e)$ spans M. It is then straightforward that the matrix $z^{-T}\Omega z^T - T\frac{dz}{z}$ of ∇ in (ε) has a simple pole.

When the lattices are further apart, this correspondence fails. However, there is also a complete description of the logarithmic lattices as follows. Let Δ be

the Deligne lattice, and let $\delta_{\Delta} = \text{Res}_{\Delta}\nabla$ be the residue \mathbb{C} -endomorphism on $\mathbf{D} = \Delta/\mathbb{m}\Delta$. Let Υ be the Birkhoff form of Δ attached to a uniformizing parameter *z*. Logarithmic lattices can then be characterized as stable flags (as already remarked in [Sabbah 2002, Theorem III.1.1]).

Proposition 52. The set Λ_{\log} of logarithmic lattices is in bijection with the subset $\Xi_0(\Upsilon)$ of filtrations of Υ which are stable under the residue, namely

$$\Xi_0(\Upsilon) = \{ (F, D) \in \Xi(\Upsilon) \mid F \in \mathfrak{Fl}_{\delta_{\Lambda}} \}.$$

Proof. According to a classical, although not so well known, result (which can be found for instance in [Babbitt and Varadarajan 1983; Bolibrukh 1990]), a lattice $\Lambda \in \mathbf{A}$ is logarithmic if and only if

(i) there exists a basis (e) of Υ such that $(z^D e)$ is a basis of Λ , with $D = \mathbf{ED}_{\Lambda}(\Lambda)$,

(ii)
$$z^{-D}Lz^{D} \in M_{n}(\mathbb{C})$$
, where $L = Mat(\nabla_{z\frac{d}{dz}}, (e))$.

It results from (ii) that in this case, the matrix L is D-parabolic. Since the flag $F^{\Delta}(\Lambda)$ induced by Λ on $\mathbf{D} = \Delta/\mathfrak{m}\Delta$ is spanned by the images of the basis (e) in \mathbf{D} , it is stable under δ_{Δ} . Conversely, it is simply a matter of computation to show that any lattice in the Υ -fiber of a δ_{Δ} -stable flag of \mathbf{D} is logarithmic. \Box

A difference between our result and Sabbah's is that he only states this result as an equivalence of categories between the set of stable filtrations of \mathbf{D} and the logarithmic lattices, whereas we give the explicit correspondence based on the lifting of \mathbf{D} to a Birkhoff form. Although it would seem that the previous result has little value to effectively determine all logarithmic lattices, it is always possible to determine them in finite terms.

Lemma 53. Let $M \in \Lambda_{\log}$ and let $(\mathbf{F}, D) = \prod_{\Delta}(M)$. Let Y be a form of Δ , and let (e) be a basis of Y respecting the flag \mathbf{F} . Fix a coordinate z, and let $P = I + P_1 z + \cdots$ be the gauge from (e) to its Υ_z -basis (e_z). Then the Laurent polynomial gauge $Q \in \mathfrak{gl}(\mathbb{C}[z, z^{-1}])$ defined by

$$Q = (I + \dots + P_{d-1}z^{d-1})z^D$$
, where $d = d(\Delta, M)$,

sends the basis (e) of Δ to a basis of M.

Proof. This is an almost direct consequence of Lemma 15.

Note that the polynomial gauge Q can be explicitly computed from formula (18). On the other hand, one can also explicitly describe the set $\Xi_0(\Upsilon)$. For a linear map $f \in \text{End}(\mathbb{C}^n)$, say that an apartment \mathscr{A}_{Φ} is a *diagonalizing apartment* of f if the frame Φ is composed of eigenlines of f.

Lemma 54. Let $\delta_{\Delta} = \mathfrak{d} + \mathfrak{n}$ be the additive Jordan decomposition of the residue map $\delta_{\Delta} = \operatorname{Res}_{\Delta} \nabla$. The pair $(\mathbf{F}, D) \in \Xi(\Upsilon)$ is an element of $\Xi_0(\Upsilon)$ if and only if \mathbf{F} admits a complete flag refinement $\hat{\mathbf{F}}$ such $\hat{\mathbf{F}} \in \mathfrak{Fl}_{\mathfrak{n}}$ and there is a diagonalizing apartment \mathcal{A} for \mathfrak{d} that respects the flag $\hat{\mathbf{F}}$.

Proof. F is δ -stable if and only if it is stable under both ϑ and \mathfrak{n} . It is known that F is stable under ϑ if and only if every component F_i of F is a direct sum of ϑ -stable lines, and under \mathfrak{n} if and only if it admits a complete flag refinement $\hat{F} \in \mathfrak{Fl}_n$. \Box

5. The Riemann-Hilbert problem

This problem is by now very well-known, so we will just state the necessary notations and definitions, and refer to the classical paper of Bolibrukh [1990] and to the account he gives of the construction of the Deligne bundle (see [Sabbah 2002; Ilyashenko and Yakovenko 2007] and also [André and Baldassarri 2001] for a purely algebraic construction).

Let $\mathcal{G} = \{s_1, \dots, s_p\}$ be a prescribed set of singular points, $z_0 \notin \mathcal{G}$ be an arbitrary base point, and let χ denote a representation

(19)
$$\chi : \pi_1(X \setminus \mathcal{G}, z_0) \to \operatorname{GL}_n(\mathbb{C}).$$

The Riemann–Hilbert problem asks informally for a linear differential system having χ as monodromy representation. In the terms used in this paper, it asks for a regular meromorphic connection ∇ with singular set \mathscr{G} and monodromy χ on a holomorphic vector bundle \mathscr{C} . If the bundle is required to be logarithmic with respect to ∇ one speaks of a *weak solution* to RH. In its strongest form, the Riemann–Hilbert problem asks for a differential system Y' = A(z)Y having simple poles on \mathscr{G} as only singularities, and whose monodromy representation is globally conjugate to χ . This amounts to asking for a weak solution (\mathscr{C}, ∇) which is moreover *trivial*.

5.1. *The Röhrl–Deligne construction.* We briefly recall H. Röhrl's construction (as presented, for instance, in [Bolibrukh 1990; Bolibrukh et al. 2006]). Let $\mathfrak{U} = (U_i)_{i \in I}$ be a finite open cover of $X^* = X \setminus \mathscr{P}$ by connected and simply connected open subsets $U_i \subset X^*$ such that their intersection has the same property, and all triple intersections are empty. Consider arbitrary points $z_i \in U_i$ and $z_{ij} \in U_i \cap U_j$, and paths $\gamma_i : z_0 \to z_i$ and $\gamma_{ij} : z_i \to z_{ij}$, so that $\delta_{ij} = \gamma_i \gamma_{ij} \gamma_{ji}^{-1} \gamma_j^{-1}$ is a positively oriented loop around z_i having winding number 1. Then the cocycle $g = (g_{ij})$ defined over \mathfrak{U} by the constant functions $g_{ij} = \chi([\delta_{ij}])$ defines a flat vector bundle \mathscr{F} over X^* . Define the connection ∇ over U_i by the (0) matrix in the basis of sections corresponding to the cocycle g. The ∇ -horizontal sections of \mathscr{F} have by construction the prescribed monodromy behavior. This solves what we called the *topological* Riemann–Hilbert problem in our introduction.

Now add a small neighborhood D of each singular point $s \in \mathcal{G}$ to the cover \mathfrak{U} , in such a way that $D \setminus \{s\}$ is covered by k pairwise overlapping sectors $\Sigma_1 = D \cap U_{j_1}$, ..., $\Sigma_k = D \cap U_{j_k}$. On an arbitrarily chosen sector among the Σ_i , say, Σ_1 , let $\tilde{g}_{s_1} = z^L$ where z is a local coordinate at s and

$$L = \frac{1}{2i\pi} \log \chi(\delta),$$

normalized with eigenvalues having their real part in the interval [0, 1]. Since the open subset Σ_1 only intersects Σ_2 and Σ_k , the only necessary cocycle relations to satisfy are $\tilde{g}_{s2} = \tilde{g}_{s1}g_{12}$ and $\tilde{g}_{sk} = \tilde{g}_{s1}g_{1k}$, which we take as *definition* of the cocycle elements \tilde{g}_{s2} and \tilde{g}_{sk} . Define in this way the remaining elements of the cocycle \tilde{g} on $D \cap U_{j_i}$. By construction, the result defines a holomorphic vector bundle \mathfrak{D} on the whole of X, and the connection ∇ can be extended as $L\frac{dz}{z}$ in the basis of sections (σ) of \mathfrak{D} over D chosen to construct \tilde{g}_{s1} . The pair (\mathfrak{D}, ∇) is called the *Deligne bundle* of χ . This construction solves simultaneously the meromorphic and the weak Riemann–Hilbert problem.

Note 55. The basis (σ) is, in our terms, a basis of the Birkhoff form attached to the coordinate z at s.

5.2. Weak and strong solutions. The Riemann–Hilbert problem can be seen as involving three different levels. The topological level is only governed by the (analytic) monodromy around the prescribed singular set. The meromorphic level is essentially based on the solution of the *local* inverse problem. The third one, that we call *holomorphic* is global and asks for the existence of a trivial holomorphic vector bundle. In fact, separating these three aspects is not so easy to do, because the Röhrl–Deligne construction in fact yields a particular *holomorphic* vector bundle % with a connection ∇ that already respects the holomorphic prescribed behavior.

What makes the strong Riemann–Hilbert problem a difficult one is precisely this third level. The local meromorphic invariants added to the topological solution of the inverse monodromy specify up to *meromorphic equivalence class* the connection ∇ on X. In this respect, the natural category to state this construction is not the category of holomorphic vector bundles with meromorphic connections, but the meromorphic vector bundles, that is, pairs (\mathcal{V}, ∇) where \mathcal{V} is locally (but in fact globally) isomorphic to \mathcal{M}_X^n . This is why we call the second step *meromorphic*. The Riemann–Hilbert problem with the given data solved here corresponds to the *very weak* Riemann–Hilbert problem (as coined in [Sabbah 2002]): any subsheaf \mathcal{F} of locally free \mathbb{O}_X -modules contained in the (trivial) meromorphic bundle \mathcal{V} is endowed naturally with the connection ∇ , and therefore is a holomorphic vector bundle with a regular connection having the prescribed monodromy. As stated by the next result (and otherwise well known), all solutions to the weak problem are obtained as local modifications of the Deligne bundle. **Proposition 56.** Let $\tilde{\pi} : \tilde{E} \to X$ and $\tilde{\nabla} : \tilde{\mathscr{E}} \to \tilde{\mathscr{E}} \otimes_{\mathbb{C}} \Omega$ be a weak solution to the Riemann–Hilbert problem. Then there exist a finite set $S \subset X$, and local lattices M_x for $x \in S$ such that the pair $(\tilde{\mathscr{E}}, \tilde{\nabla})$ is holomorphically isomorphic to (\mathfrak{D}^M, ∇) .

The last step of the strong Riemann–Hilbert problem consists of searching the set of holomorphic vector bundles endowed with the connection ∇ for a bundle which at the same time has the required holomorphic invariants and is holomorphically trivial. A negative answer requires to know all the holomorphic vector bundles with this prescribed logarithmic property. Note that up to this point, the discussion presented in this section holds over an arbitrary compact Riemann surface.

5.2.1. *Plemelj's theorem.* In 1908, the Slovenian mathematician J. Plemelj (see [Plemelj 1964]) proved a first version of the strong Riemann–Hilbert problem, under the assumption that at least one monodromy is diagonalizable. Whereas his first proof used an analytic approach (Fredholm integrals) to construct the actual matrix of solutions, to thence deduce the differential system and prove that it has only simple poles, the general framework of vector bundles recalled so far allows to establish this fact in an amazingly concise way.

Theorem 57 (Plemelj). If one of the elementary monodromy maps from representation $\chi : \pi_1(X \setminus \mathcal{G}, z_0) \to \operatorname{GL}_n(\mathbb{C})$ is diagonalizable, then the Riemann–Hilbert problem has a strong solution.

Proof. Let (\mathfrak{D}, ∇) be the Röhrl–Deligne bundle attached to the representation χ . Let, say $G = \chi(\gamma)$ around $s \in \mathcal{G}$, be diagonalizable. Let Υ be a Birkhoff form at s, and let (e) be a basis of Υ where G is diagonal. According to condition (ii) in Section 4.2, the whole apartment \mathcal{A} spanned by (e) consists of logarithmic lattices, whereas Theorem 44 implies that \mathcal{A} contains a trivializing lattice M. The vector bundle \mathfrak{D}^M is therefore both logarithmic and trivial. \Box

Note 58. Here we have a solution by modifying the Deligne bundle only at one point. Note that the lattice M corresponds to a BG trivialization of \mathfrak{D} (see Theorem 59 below). Also note that this result also holds replacing \mathfrak{D} with any other weak solution to Riemann–Hilbert.

5.2.2. *Trivializations of weak solutions.* Let \mathscr{E} be a weak solution of the Riemann– Hilbert problem, and let \mathscr{F} be a trivialization of \mathscr{E} at $x \notin \mathscr{F}$. In a global basis of sections (e) of the bundle \mathscr{F} , the connection ∇ is expressed by the matrix of global meromorphic 1-forms Ω , which has a simple pole at every $s \in \mathscr{F}$, and an *a priori* uncontrolled pole at x. Assuming for simplicity that $x \notin \mathscr{F}$ is the point at infinity $\infty \in \mathbb{P}^1(\mathbb{C})$, there exist matrices $A_i \in M_n(\mathbb{C})$ for $1 \leq i \leq p$ and a matrix

$$B(z) = B_0 + \dots + B_t z^t$$

such that the connection has the matrix

$$\Omega = \left(\sum_{i=1}^{p} \frac{A_i}{z - s_i} + B(z)\right) dz.$$

The most surprising consequence of the permutation lemma, as we state it, concerns the analytic invariants of the weak solutions to the Riemann–Hilbert problem.

Theorem 59. Let \mathscr{C} be a weak solution to the Riemann–Hilbert problem for χ . Then, for any $x \notin \mathscr{G}$, there exists a BG trivialization \mathscr{F} of \mathscr{C} at x which is also logarithmic at x. Let $Y = \Gamma(X, \mathscr{F})$ and let $\psi_s = \operatorname{Res}_s^{\mathscr{F}} \nabla \in \operatorname{End}(Y)$.

- (1) The map $\Psi = \sum_{s \in \mathcal{S}} \psi_s = -\text{Res}_x \nabla$ is semisimple, and has integer eigenvalues, which are equal to the type of the bundle \mathscr{E} .
- (2) The image of the Harder–Narasimhan filtration of C in Y is equal to the flag induced by the eigenspaces of Ψ ordered by increasing values.

Proof. If $x \notin \mathcal{G}$, the monodromy at x is trivial, and the stalk \mathscr{C}_x of \mathscr{C} coincides with \mathfrak{D}_x . The Birkhoff form Υ of D (which is then unique) is equal to the space V^{∇} of horizontal sections at x. All flags in $\mathbf{D} = \mathfrak{D}_x/\mathfrak{m}_x\mathfrak{D}_x$ are stable under $\operatorname{Res}_x^{\mathfrak{D}} \nabla = 0$. According to Corollary 46, the Υ -lifting of the flag induced by any BG trivialization of \mathscr{C} at x is a *logarithmic BG trivialization* of \mathscr{C} at x. In a global basis of sections (e) of \mathscr{F} , the connection has the matrix

$$A = \sum_{s \in \mathcal{G} \setminus \{\infty\}} \frac{A_s}{z - s} + \frac{B}{z - x} \quad \left(\text{where } B = -\sum_{s \in \mathcal{G}} A_s \text{ if } x \neq \infty \right)$$
$$= \sum_{s \in \mathcal{G}} \frac{A_s}{z - s} \quad \text{if } x = \infty \notin \mathcal{G},$$

since ∇ has no other singularities outside $\mathcal{G} \cup \{x\}$. The eigenvalues of $-B = \sum_{s \in \mathcal{G}} A_s$ are therefore equal to the type of \mathcal{C} , and the Harder–Narasimhan filtration is defined by the blocks of equal eigenvalues ordered by increasing values. \Box

As a consequence, we deduce the following new sufficient condition for the solubility of the strong Riemann–Hilbert problem.

Corollary 60. Let $\mathscr{C} \in \mathscr{H}$ and let \mathfrak{D} be the Deligne lattice of (\mathscr{V}, ∇) . Let $x \in X$, such that $\mathscr{C}_x = \mathfrak{D}_x = \Delta$. Let $\mathbf{D} = \Delta/\mathfrak{m}_x \Delta$. Let $\mathscr{F} \in BG_x(\mathscr{C})$, and $M = \mathscr{F}_x$. If the flag $\mathbf{F}^{\Delta}(M)$ induced in D is stable under $\operatorname{Res}_x^{\mathfrak{D}} \nabla \in \operatorname{End}(\mathbf{D})$, then there exists $\widetilde{\mathscr{F}} \in BG_x(\mathscr{C})$ which is moreover logarithmic at x.

Proof. Let \tilde{M} be the Υ -lifting of the flag $F^{\Delta}(M)$, where Υ is a Birkhoff form of Δ . According to Proposition 52, the lattice \tilde{M} is logarithmic, and by the permutation lemma, it is a BG trivializing lattice. Therefore, the bundle $\mathscr{C}^{\tilde{M}}$ satisfies the conclusions of the corollary. At this point, we would like to sum up our findings about trivial bundles in the following proposition.

Proposition 61. Let $\mathcal{F} \in \mathcal{H}_0$ be a trivial bundle in \mathcal{V} , and let $Y = \mathcal{F}(X)$ be the \mathbb{C} -vector space of global sections. Let $x \in X$, and $\mathcal{E} \in H$ such that $\mathcal{F} \sim_x \mathcal{E}$. Let moreover $\Lambda = \mathcal{E}_x$ and $M = \mathcal{F}_x$.

- (i) Y has a well-defined flag HN_ℓ induced by the Harder–Narasimhan filtration of ℓ.
- (ii) If $\mathcal{F} \in BG_x(\mathcal{C})$, then $T(\mathcal{C}) \cong ED_{\Lambda}(M)$ and any Smith basis of Y for Λ ordered according to K^{\searrow} is strictly adapted to $HN_{\mathcal{C}}$.
- (iii) If \mathscr{F} is additionally logarithmic at x, and the stalk \mathscr{C}_x coincides with the Deligne lattice \mathfrak{D}_x , then the type $T(\mathscr{C})$ is given by the integer parts of the eigenvalues of the residue $\operatorname{Res}_x^{\mathscr{F}} \nabla \in \operatorname{End}(Y)$, that is, of the exponents of ∇ on \mathscr{F} at x.
- (iv) Finally, if $\mathscr{C} \in RH_{\chi}$ is moreover a weak solution to Riemann–Hilbert, then

$$\sum_{x \in X} \operatorname{Res}_{x}^{\mathcal{F}} \nabla = 0.$$

When $(\mathcal{E}, \mathcal{F})$ satisfy (i) to (iv), we say that \mathcal{F} is a good RH trivialization of \mathcal{E} at x.

Let \mathcal{F} be a good RH trivialization of \mathcal{C} at $x \notin \mathcal{F}$. Let (σ) be any basis of $Y = \Gamma(X, \mathcal{F})$. In (σ) , the connection has a matrix of the form (21). The identification of Y to \mathbb{C}^n by means of (σ) endows \mathbb{C}^n with p + 1 linear maps ψ_s for $s \in \mathcal{F}^* = \mathcal{F} \cup \{x\}$, that we can identify with the matrices \tilde{L}_s for $s \in \mathcal{F}$ and $-\sum_{s \in \mathcal{F}} \tilde{L}_s$ for s = x. With these notations, we set the following definition.

Definition 62. The space \mathbb{C}^n , endowed with the maps ψ_s for $s \in \mathcal{G}^*$, is called a *linear Fuchsian model* of \mathscr{C} .

With this notion, we can reduce some questions about vector bundles to linear algebra statements. For instance we can give the following computable version of a criterion due to Gabber for the reducibility of the triviality index originally appearing in [Sabbah 2002, Corollary I.4.14], that we state here only for the case of a logarithmic modification.

Corollary 63. Let $\mathscr{C} \in \operatorname{RH}_{\chi}$ be a weak solution, and consider a linear Fuchsian model at $x \notin \mathscr{G}$, given by p matrices A_s for $s \in \mathscr{G}$ such that

$$\sum_{s\in\mathscr{S}} A_s = \operatorname{diag}(t_1 I_{n_1}, \ldots, t_s I_{n_s})$$

where the integers \mathbf{t}_i satisfy $\mathbf{t}_i > \mathbf{t}_{i+1}$, in such a way that the flag HN is the flag $0 = F_0 \subset F_1 \subset \cdots \subset F_s = \mathbb{C}^n$ having signature (n_1, \ldots, n_s) in the canonical basis of \mathbb{C}^n . There exists a weak solution \mathcal{E}' adjacent to \mathcal{E} at $s \in X$ and such that

 $i(\mathscr{C}') < i(\mathscr{C})$ if and only if there exists an A_s -stable subspace $W \subset \mathbb{C}^n$ such that $W \cap F_1 = (0)$.

Proof. Let $T = \text{diag}(t_1 I_{n_1}, \dots, t_s I_{n_s}) = \text{diag}(t_1, \dots, t_n)$ be the type of \mathscr{E} . We have $i(\mathscr{E}) = \sum_{i=1}^{s} n_i(t_1 - t_i)$. According to Proposition 28, any adjacent weak solution \mathscr{E}' is given by an A_s -stable subspace $W \subset \mathbb{C}^n$. For any basis (e) of \mathbb{C}^n respecting the flag HN, the bundle \mathscr{E}' has type T' = T - K, where $k_i = 0$ when $e_i \in W$ and $k_i = 1$ otherwise; therefore $i(\mathscr{E}') = \sum_{i=1}^{n} (\max(t_i - k_i) - t_i + k_i)$ where t_i represent the elements of T without multiplicities. Accordingly, we have

$$i(\mathscr{C}) - i(\mathscr{C}') = \sum_{i=1}^{n} (t_1 - k_i - \max(t_i - k_i)).$$

Now, if there exists *i* such that $t_i = t_1$ and $k_i = 0$, then $\max(t_i - k_i) = t_1$; thus $i(\mathscr{E}) - i(\mathscr{E}') = \sum_{i=1}^n -k_i < 0$ (because we exclude the trivial case $W = \mathbb{C}^n$). Otherwise we have $\max(t_i - k_i) = t_1 - 1$, and then $i(\mathscr{E}) - i(\mathscr{E}') = \sum_{i=1}^n (1 - k_i) > 0$. Therefore \mathscr{E}' exists if and only if there exists *W* stable under some A_s such that $W \cap F_1 = 0$.

Proposition 64. Let \mathcal{F} be a BG trivialization of \mathfrak{D} at $x \notin \mathcal{F}$. If there exists a flag F in $Y = \mathcal{F}(X)$ which is transversal to $HN_{\mathfrak{D}}$, and is moreover stable under the action of one of the maps ψ_s for $s \in \mathcal{F}$, then the strong Riemann–Hilbert problem has a solution, which moreover coincides with \mathfrak{D} outside s.

Proof. Let F be a flag of Y, which is stable under ψ_s . Taking stalks at x of a \mathbb{C} -basis of F, we can see the flag F in $\mathbf{D} = \mathfrak{D}_s/\mathfrak{m}_s\mathfrak{D}_s$. According to Lemma 27(ii), there exists a BG trivialization \mathscr{C} of \mathfrak{D} at x, whose image in $D = \mathfrak{D}_s/\mathfrak{m}_s\mathfrak{D}_s$ is F. Let (e) be a BG basis of \mathfrak{D}_s with respect to \mathscr{C}_s . Consequently, its image in \mathbf{D} respects the flag F. Let Υ be a Birkhoff form of \mathfrak{D}_s , and let (e_{Υ}) be the Υ -basis of (e). Since the gauge from (e) to (e_{Υ}) is tangent to I, the lattice M induced from (e_{Υ}) by the elementary divisors K of \mathscr{C}_s in Λ is also a trivializing BG lattice for \mathfrak{D} at s. However, the lattice M is also logarithmic, since by construction it induces in \mathbf{D} the ψ_s -stable flag F, and moreover sits inside an apartment that contains the Birkhoff form Υ . Hence, the bundle \mathfrak{D}^M is both trivial and logarithmic. \square

We have represented the weak solutions to the Riemann–Hilbert problem as points in a product of subvarieties of stable flags.

Theorem 65. Let \mathfrak{D} be the Deligne bundle, and \mathfrak{F} a BG trivialization at an apparent singularity $x \notin \mathfrak{S}$. The set of weak solutions to the Riemann–Hilbert problem for χ is parametrized by the set

$$\operatorname{RH}_{\chi} = \{ (\boldsymbol{F}^{s}, D_{s})_{s \in \mathcal{G}} \mid \boldsymbol{F}^{s} \in \mathfrak{Fl}_{\psi_{s}}(Y), D_{s} \in \mathbb{Z}^{n}(\boldsymbol{F}^{s}) \},\$$

where $Y = \mathcal{F}(X)$ and $\psi_s = \operatorname{Res}_s^{\mathcal{F}} \nabla \in \operatorname{End}_{\mathbb{C}}(Y)$ for $s \in \mathcal{G}$.

5.3. *The type of the Deligne bundle.* The strong version of the Riemann–Hilbert problem would directly have a solution if the Deligne bundle were trivial. However, this is not the case, unless all singular points are apparent, since the exponents of ∇ are normalized in such a way that their sum is nonnegative. This means that the type of the Deligne bundle as a rule is not trivial. We have seen several ways to characterize this nontriviality. The *type* characterizes the isomorphism classes of holomorphic vector bundles, so it would seem possible to work with this sole information. However, we are not in the right category to do so, since we consider holomorphic bundles with an embedding in a meromorphic one, denoted by \mathcal{V} . This is the reason for which there are *several* trivial bundles in \mathcal{V} . From another point of view, it is not possible to determine on the sole basis of the sequence $T = (a_1, \ldots, a_n)$, what the effect of changing the stalk of \mathfrak{D} at *x* will be. Obviously the geometry of the Harder–Narasimhan filtration will play a decisive role.

5.3.1. *Trivializations of the Deligne bundle.* Let us examine in further detail the case of the Deligne bundle \mathfrak{D} . Let us say that δ_i is an elementary generator of the homotopy group $G = \pi_1(X \setminus \mathcal{G}, z_0)$, if δ_i is a closed path based at z_0 , having winding number +1 around the singularity s_i and 0 around the others. Let $G_i = \chi(\delta_i)$ and $L_i = \frac{1}{2i\pi} \log G_i$, normalized as for the Deligne lattice. Let (σ_i) be a basis of the Birkhoff form Υ_i at s_i described in Note 55, such that the connection has locally as matrix $\Omega_i = L_i \frac{dz}{z}$, on a neighborhood, say D_i of s_i . On the other hand, let D_0 be a neighborhood of z_0 , and consider a basis (σ_0) of the local Birkhoff form. According to what precedes, (σ_0) is a basis of local ∇ -horizontal sections of \mathfrak{D} over D_0 . One can moreover choose this basis in such a way that the monodromy of (σ_0) around s_i is exactly given by the matrix G_i .

Assume now for simplicity that $x \notin \mathcal{G}$ is the point at infinity $\infty \in \mathbb{P}^1(\mathbb{C})$, and let \mathcal{F} be a trivialization of \mathfrak{D} at x. In a global basis of sections (*e*) of the bundle \mathcal{F} , there exist matrices $B_i \in M_n(\mathbb{C})$ and a matrix

$$B(z) = B_0 + \dots + B_t z^t$$
 and $C_i \in \operatorname{GL}_n(\mathbb{C})$ for $1 \le i \le p$

such that the connection has the matrix

$$\Omega = \left(\sum_{i=1}^{p} \frac{C_i^{-1} L_i C_i}{z - s_i} + B(z)\right) dz.$$

Note 66. If the bundle \mathcal{F} is moreover logarithmic at ∞ — which can be achieved, for example, by Plemelj's theorem — then B = 0 and the residue at infinity, $L_{\infty} = -\sum_{i=1}^{p} C_i^{-1} L_i C_i$, is semisimple with integer eigenvalues (ssie). At the cost of a (harmless) global conjugation, we can already assume that

$$L_{\infty} = \operatorname{diag}(b_1 I_{n_1}, \dots, b_s I_{n_s}) \quad \text{with } b_1 < \dots < b_s.$$

Note that the sequence $\mathfrak{B} = (b_1 I_{n_1}, \dots, b_s I_{n_s})$ coincides with the elementary divisors of the stalk \mathcal{F}_{∞} in \mathfrak{D}_{∞} .

Definition 67. We say that $(C_1, \ldots, C_p) \in GL_n(\mathbb{C})^p$ is a *normalizing p-tuple* for χ if $\sum_{i=1}^p C_i^{-1}L_iC_i$ is ssie for some (and therefore any) normalized logarithms L_i of the generators $\chi(\gamma_i)$ of the monodromy group.

Normalizing *p*-tuples always exist. Putting *t* as the coordinate 1/z at infinity, the Taylor expansion of ∇ at $x = \infty$ has then the nice expression

(20)
$$\Omega = -\sum_{k\geq 0} \sum_{i=1}^{p} s_i^k \widetilde{L}_i t^k \frac{dt}{t} \quad \text{with } \widetilde{L}_i = C_i^{-1} L_i C_i.$$

We have thus reduced the computation of the type of the Deligne bundle to the computation of the matrices C_i (the so-called *connection matrices*, because they connect the different local expressions of ∇ on the local Birkhoff forms). It is however well known that the computation of the connection matrices is difficult. Any other trivialization of \mathfrak{D} at infinity is given by a *monopole gauge* (as coined in [Ilyashenko and Yakovenko 2007]), namely a unimodular polynomial matrix $\Pi \in GL_n(\mathbb{C}[z])$, that is, a matrix satisfying

$$\Pi = P_0 + P_1 z + \dots + P_k z^k \quad \text{such that } \det \Pi(z) = \operatorname{cst} \in \mathbb{C}^*.$$

Proposition 68. Given a family of points $s_1, \ldots, s_p \in \mathbb{C}$ and invertible matrices $C_1, \ldots, C_p \in \operatorname{GL}_n(\mathbb{C})$, there exists a monopole gauge $\Pi \in \operatorname{GL}_n(\mathbb{C}[z])$ such that $\Pi(s_i) = C_i$ for $1 \leq i \leq p$.

Proof. The group $GL_n(R)$ on a ring is generated by transformations $T_{ij}(\lambda) = I + \lambda E_{ij}$ where $\lambda \in R$ and E_{ij} is the (i, j) element of the canonical basis of the vector space \mathfrak{gl}_n . At the cost of introducing the trivial transformations $T_{ij}(0) = I$, one can assume that all the matrices C_i can be expressed as a product of the same transformations with different parameters:

$$C_i = T_1(\mu_1^i) \cdots T_s(\mu_s^i)$$
 with $\mu_t^i \in \mathbb{C}$.

Define then $\lambda_k \in \mathbb{C}[z]$ such that $\lambda_k(s_i) = \mu_k^i$ for $1 \le i \le p$. By construction, the product $\widetilde{\Pi} = T_1(\lambda_1) \cdots T_s(\lambda_s) \in SL_n(\mathbb{C}[z])$ indeed interpolates the matrices C_i at the points s_i .

As a consequence of this result, one can find a trivialization \mathscr{E} at infinity of the Deligne bundle such that the residues of the connection ∇ are expressed in a basis of $Y = \Gamma(X, \mathscr{E})$ as the actual matrices L_i (and not *conjugated* to them). Although the point at infinity of \mathscr{E} is still an apparent singularity, we have no control on the Poincaré rank of ∇ at ∞ .

The results of this section also hold (with the adequate modifications) if the apparent singularity is assumed to be located at $z_0 \notin \mathcal{G} \cup \{\infty\}$. We will refer to the trivialization \mathscr{E} as an *adapted trivialization of* \mathfrak{D} *at* z_0 .

5.3.2. A Deligne–Simpson-type problem. We know that there exists a family of invertible matrices (C_i) such that $\sum_{i=1}^{p} C_i^{-1} L_i C_i$ is semisimple with integer eigenvalues and that these eigenvalues are equal to the type of the Deligne bundle. This raises two questions:

- (1) Does there exist a logarithmic trivialization of \mathfrak{D} for any such family (C_i) ?
- (2) If there exist several families with this property, how do we recognize those that indeed give the type of the Deligne bundle?

This also raises an interesting computational problem akin to the well-known Deligne–Simpson problem (see, e.g., [Crawley-Boevey 2003]). Let \mathscr{C}_i be the conjugacy class of $L_i = \frac{1}{2i\pi} \log G_i$ under $\operatorname{GL}_n(\mathbb{C})$.

(DS) Determine all conjugates $\tilde{L}_i \in \mathscr{C}_i$ such that $\sum_{i=1}^p \tilde{L}_i = \text{diag}(b_1, \ldots, b_n) \in \mathbb{Z}^n$.

5.4. *The Bolibrukh–Kostov theorem.* The most celebrated recent result on the Riemann–Hilbert problem is the following fact, proved first independently by A. Bolibrukh and V. Kostov.

Theorem 69 (Bolibrukh–Kostov). *The strong Riemann–Hilbert problem is solvable for any* irreducible *monodromy representation* χ .

We give first an algebraic proof of a classical result of Bolibrukh [Anosov and Bolibrukh 1994, Proposition 4.2.1].

Proposition 70. If the representation χ is irreducible, then for any weak solution $\mathscr{C} \in \operatorname{RH}_{\chi}$, the type $T(\mathscr{C}) = (t_1, \ldots, t_n)$ of \mathscr{C} satisfies $|t_i - t_j| \le |i - j|(p - 2)$.

Proof. Assume here for simplicity that $x = \infty \notin \mathcal{G}$, and consider again the setting of Section 5.3.1. Let \mathscr{C} be any weak solution to Riemann–Hilbert, and \mathscr{F} be a logarithmic BG trivialization of \mathscr{C} at x. Let $T = (t_1, \ldots, t_n)$ be the type of \mathscr{C} . In a basis (e) of global sections of \mathscr{F} , there exist constant matrices \widetilde{L}_a for $a \in \mathscr{G}$ such that the connection ∇ has in (e) the matrix

(21)
$$\Omega = \sum_{a \in \mathcal{G}} \frac{\widetilde{L}_a}{z - a} dz = -\frac{d\widetilde{z}}{\widetilde{z}} \sum_{k \ge 0} \Omega_k \widetilde{z}^k \quad \text{with } \Omega_k = \sum_{a \in \mathcal{G}} a^k \widetilde{L}_a \text{ and } \widetilde{z} = \frac{1}{z}.$$

By Proposition 49, the shearing \tilde{z}^{-T} suppresses the singularity at *x*, since the basis $\tilde{z}^{-T}(e)$ spans the Deligne lattice. As a consequence, $\tilde{\Omega} = \tilde{z}^T \Omega(\tilde{z}) \tilde{z}^{-T} + T \frac{d\tilde{z}}{\tilde{z}}$ must satisfy $v(\tilde{\Omega}) \ge 0$. Therefore, the residue matrix $B = -\sum_{a \in \mathcal{G}} \tilde{L}_a$ of Ω at *x* is diagonal and equal to -T. We can assume further that

$$B = \operatorname{diag}(b_1 I_{n_1}, \dots, b_s I_{n_s}) \quad \text{with } b_1 = -t_1 < \dots < b_s = -t_s$$

where $(t_1 I_{n_1}, \ldots, t_s I_{n_s})$ represents the type of \mathscr{C} with multiplicities. Partition any matrix M according to the eigenvalue multiplicities of B, as $(M_{\ell,m})$ for $1 \le \ell, m \le s$. Then the matrix of the connection can be rewritten by blocks as

$$\widetilde{\Omega}_{\ell,m} = \Omega_{\ell,m} t^{t_{\ell}-t_m} + T \frac{d\widetilde{z}}{\widetilde{z}} = \left(-\sum_{j\geq 0} \Omega_{\ell,m}^{(j)} \widetilde{z}^{j+t_{\ell}-t_m} + \delta_{\ell,m} t_{\ell} I_{n_{\ell}} \right) \frac{d\widetilde{z}}{\widetilde{z}}.$$

For each (ℓ, m) block, this series must have strictly positive valuation. The sum $\sum_{a \in \mathcal{G}} \tilde{L}_a = T$ imposes conditions on all blocks of the residues \tilde{L}_a , while when $\ell > m$ we get the following equations:

(22)
$$\Omega_{\ell,m}^{(j)} = \sum_{a \in \mathcal{G}} a^j (\tilde{L}_a)_{\ell,m} = 0 \text{ for } 0 \le j \le t_m - t_\ell \quad \text{when } \ell > m.$$

For a fixed pair (ℓ, m) , let $k = \max(0, t_m - t_\ell)$, and let $X_i \in \mathbb{C}^{n_\ell \times n_m}$ be the (ℓ, m) -block of the matrix \tilde{L}_{s_i} , for $1 \le i \le p$. For $1 \le \alpha \le n_\ell$ and $1 \le \beta \le n_m$, let $v_{\alpha,\beta} \in \mathbb{C}^p$ be the vector constructed by taking the coefficient of index (α, β) of X_i , for $1 \le i \le p$. Then, the equations (22) can be reformulated as

$$v_{\alpha,\beta} \in \ker M_k(\underline{s})$$
 where $M_k(\underline{s}) = \begin{pmatrix} 1 & \cdots & 1 \\ s_1 & \cdots & s_p \\ \vdots & & \vdots \\ s_1^k & \cdots & s_p^k \end{pmatrix}$.

The matrix $M_k(\underline{s})$ is an upper-left submatrix of a Vandermonde matrix with coefficients

$$\underline{s} = (s_1, \ldots, s_p) \in \mathbb{C}^p \setminus \bigcup_{i \neq j} \{x_i \neq x_j\}.$$

Since all the s_i are distinct, this matrix has always full rank. In particular, as soon as $t_m - t_{\ell} \ge p - 1$, it has a null kernel, and so all the blocks X_i are zero. Due to the ordering of the t_i , we also have $t_{m'} - t_{\ell'} \ge p - 1$ for $m' \le m$ and $\ell' \ge \ell$; thus all matrices \tilde{L}_a have a lower-left common zero block. If $m = \ell + 1$, this means that the representation χ is reducible.

Proof of Theorem 69. Consider the Deligne bundle \mathfrak{D} and an arbitrary singularity $s \in \mathcal{G}$. Put $\Delta = \mathfrak{D}_s$ and let \mathbf{F} be a complete flag in $\Delta/\mathfrak{m}\Delta$ which is stable under $\delta = \operatorname{Res}_s^{\mathfrak{D}} \nabla$. Let D = (0, d, 2d, ..., (n-1)d). By Proposition 52, the lattice Λ obtained by lifting the pair (\mathbf{F}, D) in the *z*-Birkhoff form $\Upsilon \subset \Delta$ is logarithmic. Let (e) be a basis of Υ respecting the flag \mathbf{F} . Since the residue δ is upper-triangular in (e), and the elements of D are distinct, the elements of the matrix $\Omega = \operatorname{Mat}(\nabla, (z^D e))$ outside the diagonal have valuation at least d. Hence $\Omega = (A + z^d U) \frac{dz}{z}$ where A is diagonal and U holomorphic. By Theorem 44, the

apartment \mathcal{A} spanned by (e) contains a Birkhoff–Grothendieck trivialization M of Λ , as in the following figure:

$$\Delta \xrightarrow{z^D} \Lambda \in \Lambda_{\log} \xrightarrow{z^{\widetilde{D}}} M \in \mathrm{BG}(\Lambda).$$

According to Proposition 70, the sequence \tilde{D} satisfies $\Delta \tilde{D} \leq (n-1)(p-2)$, while the Poincaré rank \mathfrak{p} of ∇ on M satisfies $\mathfrak{p} = \min(0, -v(z^{-\tilde{D}}\Omega z^{\tilde{D}}))$. Since $v(z^{-\tilde{D}}\Omega z^{\tilde{D}}) \geq d - \Delta \tilde{D}$ holds, it is sufficient to impose $d \geq (n-1)(p-2)$ to ensure that $\mathfrak{p} = 0$, that is, that the lattice M is logarithmic. Then \mathfrak{D}^M is a strong solution.

The previous proof holds in fact for any representation χ for which there exists a constant *R* such that for any weak solution \mathscr{C} of RH for χ , the type of \mathscr{C} satisfies $\Delta T(\mathscr{C}) \leq R$. Finally, as a byproduct of the proof of Proposition 70, we have the following result.

Corollary 71. Let $\tilde{L}_i \in \mathcal{C}_i$ such that $\sum_{i=1}^{p} \tilde{L}_i = \text{diag}(t_1, \ldots, t_n) \in \mathbb{Z}^n$. Then the sequence $T = (t_1, \ldots, t_n)$ represents the type of the Deligne bundle of the monodromy representation χ of the Fuchsian system (with singular locus $\mathcal{G} = \{s_1, \ldots, s_p\} \subset \mathbb{C}$)

$$Y' = \sum_{i=1}^{p} \frac{\tilde{L}_i}{z - s_i} Y$$

if and only if $\sum_{i=1}^{p} s_i^j (\tilde{L}_i)_{\ell,m} = 0$ holds for $0 \le j \le k_m - k_\ell$ and $\ell > m$.

This gives a partial answer to our question (2) from Section 5.3.2. For any explicit solution $L = (L_1, \ldots, L_p)$ to the generalized Deligne–Simpson problem, we can generate nontrivial explicit examples of Deligne bundles (and their Harder–Narasimhan filtrations) corresponding to monodromy representations which are *locally conjugate* to the original one. In the following section, we give the algorithmic procedures that can then be used to determine effectively if the corresponding representation admits or not a strong solution to the Riemann–Hilbert problem. Note that the equations in Corollary 71 show that the set of singular loci \mathcal{G} for which L corresponds to a logarithmic BG trivialization of a Deligne bundle is an algebraic (projective) subvariety of $\mathbb{C}^p \setminus \bigcup_{i \neq j} \{x_i = x_j\}$.

5.5. *Testing the solubility of the Riemann–Hilbert problem.* In this section, we apply the results of this paper to the experimental investigation of the solubility of the Riemann–Hilbert problem. We present two ways to search the space of weak solutions, which are completely effective (up to the known problem of connection matrices): one that follows paths of adjacent logarithmic lattices, based on Lemma 51, the other that uses the characterization as stable flags given in

Proposition 52. Note that, if any (not necessarily logarithmic) trivial holomorphic bundle of the meromorphic solution to Riemann–Hilbert is explicitly given, the procedures that we present, coupled with classical Poincaré rank reduction methods, implemented on a computer algebra system, allow to make the actual computations. We however do not know if this bypasses the problem of the connection matrices.

Let \mathfrak{D} be the Deligne bundle of the representation χ . Let $x \notin \mathcal{G}$, and consider a logarithmic BG trivialization \mathcal{F} of \mathfrak{D} at x. Let $Y = \Gamma(X, \mathcal{F})$ and choose a basis (σ) of Y in which the residue matrix at x is equal to the diagonal that represents the type of \mathfrak{D} :

$$\operatorname{Mat}(\operatorname{Res}_{x}^{\mathcal{F}}\nabla, (\sigma)) = -D = \operatorname{diag}(-k_{1}I_{n_{1}}, \dots, -k_{s}I_{n_{s}}) \quad \text{where } k_{1} > \dots > k_{s}.$$

In the basis (σ), the connection has a matrix of the form (21), and the Harder– Narasimhan filtration is expressed as the flag HN_Y of signature (n_1, \ldots, n_s) of Y. Let $V = \Gamma(X, \mathcal{V})$ be the \mathcal{K} -vector space of meromorphic sections of \mathcal{V} , where $\mathcal{K} = \Gamma(X, \mathcal{M}_X)$ is the field of meromorphic functions on X.

For $s \in \mathcal{G}$, let t be a coordinate at x with divisor (t) = x - s, and $(\tilde{\sigma}) = t^{-D}(\sigma)$. Recall that t^{-1} is a coordinate at s. For clarity's sake, we will put $t_x = t$ and $t_s = t^{-1}$ when we are dealing with local sections. Let $\tilde{\mathcal{F}} = t_s(\mathcal{F})$ be the transport of \mathcal{F} at s and $\tilde{Y} = \Gamma(X, \tilde{\mathcal{F}})$. We regard Y and \tilde{Y} as sub- \mathbb{C} -vector spaces of V, spanned respectively by the \mathcal{R} -bases (σ) and $(\tilde{\sigma})$ of V. The relation $(\tilde{\sigma}) = t^{-D}(\sigma)$ induces a well-defined fixed isomorphism between Y and \tilde{Y} .

Claim 72. The trivial bundle $\tilde{\mathcal{F}}$ is a BG trivialization of \mathfrak{D} at s.

Claim 73. The flag $HN_{\tilde{Y}}$ is the flag of signature (n_1, \ldots, n_s) spanned by $(\tilde{\sigma})$.

Claim 74. The germ (σ_s) of the global basis of Y at s is a local basis of \mathfrak{D}_s .

Indeed, we have the two dual schematic representations, where (σ_x) : \mathscr{C}_x means that (σ) is a local basis of \mathscr{C} at x and (σ) : Y means that (σ) is a global basis of the form Y:

$$(\widetilde{\sigma}_x): \mathfrak{D}_x \xrightarrow{t_x^D} (\sigma): Y \text{ and } (\sigma_s): \mathfrak{D}_s \xrightarrow{t_s^D} (\widetilde{\sigma}): \widetilde{Y}.$$

5.5.1. *Bolibrukh's first counterexample.* Bolibrukh's first published counterexample is the 3×3 system

(23)
$$dX/dz = AX,$$
$$A = \begin{pmatrix} 0 & 1 & 0 \\ 0 & z & 0 \\ 0 & 0 & -z \end{pmatrix} \frac{1}{z^2} + \begin{pmatrix} 0 & 1 & 0 \\ 0 & -\frac{1}{6} & \frac{1}{6} \\ 0 & -\frac{1}{6} & \frac{1}{6} \end{pmatrix} \frac{1}{z+1} + \begin{pmatrix} 0 & 0 & 1 \\ 0 & -\frac{1}{2} & -\frac{1}{2} \\ 0 & \frac{1}{2} & \frac{1}{2} \end{pmatrix} \frac{1}{z-1} + \begin{pmatrix} 0 & -1 & -1 \\ 0 & -\frac{1}{3} & \frac{1}{3} \\ 0 & -\frac{1}{3} & \frac{1}{3} \end{pmatrix} \frac{1}{z-\frac{1}{2}}.$$

Let us show what the different notions introduced in the paper are in this case. We consider the system (23) to be the expression of the connection ∇ on a trivial bundle \mathscr{C} , in a basis (σ) of global sections. The singular divisor is $\mathfrak{D} = 2 \cdot 0 + 1 - 1 + \frac{1}{2}$ and the matrices at 1, -1 and $\frac{1}{2}$ are nilpotent. The point at infinity is not singular, since, putting t = 1/z, we have

(24)
$$-\frac{1}{t^2}A\left(\frac{1}{t}\right) = \frac{1}{2}\begin{pmatrix}-1 & 0 & -1\\ -1 & 0 & 1\\ 1 & 0 & 1\end{pmatrix} + o(1).$$

Therefore the stalks \mathscr{C}_x for $x \in \mathbb{P}^1(\mathbb{C}) \setminus \{0\}$ coincide with the Deligne bundle \mathfrak{D} of ∇ : in the terms of Section 1, the bundle \mathscr{C} is a trivialization of \mathfrak{D} at 0. However, the singular point 0 is not an apparent singularity. The gauge P = diag(1, z, 1/z) brings the system to the form

$$A_{[P]} = \frac{1}{z}\tilde{A}_0 + \frac{1}{z+1}\tilde{A}_{-1} + \frac{1}{z-1}\tilde{A}_1 + \frac{1}{z-\frac{1}{2}}\tilde{A}_{\frac{1}{2}} + B$$

where

$$\begin{split} \widetilde{A}_{0} &= \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{pmatrix} \widetilde{A}_{-1} = \frac{1}{6} \begin{pmatrix} 0 & -6 & 0 \\ 0 & -1 & -1 \\ 0 & -1 & 1 \end{pmatrix} \widetilde{A}_{1} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 2 \\ 0 & -1 & 1 \\ 0 & -1 & 1 \end{pmatrix} ,\\ \widetilde{A}_{\frac{1}{2}} &= \frac{1}{12} \begin{pmatrix} 0 & -6 & -24 \\ 0 & -4 & 16 \\ 0 & -1 & 4 \end{pmatrix} \quad \text{and} \quad B = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}, \end{split}$$

The bundle \mathcal{F} spanned by the global sections $(\sigma') = (\sigma_1, z\sigma_2, \sigma_3/z)$ is trivial by construction, and since all residues over $\mathcal{F} = \{0, 1, -1, 1/2\}$ are nilpotent, \mathcal{F} is a trivialization of \mathfrak{D} at ∞ . Moreover, the stalk \mathfrak{D}_{∞} is spanned by (σ_{∞}) , so the germ of basis (σ'_{∞}) is a Smith basis of \mathfrak{D}_{∞} ; hence \mathcal{F} is actually a BG trivialization of the Deligne bundle \mathfrak{D} at $x = \infty$. In general, such a gauge can be found explicitly by combining a Poincaré rank reducing method at all finite singularities (e.g., Gérard–Levelt saturation [1973]) and the BG trivialization algorithm from Section 3.2.1.

Accordingly, the type of the Deligne bundle is $T = T(\mathfrak{D}) = (1, 0, -1)$. However, an apparent singularity of Poincaré rank 2 appears at ∞ ; hence the BG trivialization \mathcal{F} is not logarithmic at ∞ . To get one, we apply the permutation lemma (Proposition 38). We reorder the basis at ∞ as $(\tilde{\sigma}) = (\sigma_3, \sigma_1, \sigma_2)$ according to the decreasing elements of the type. Putting $\Lambda = \mathfrak{D}_{\infty}$ and $M = \mathcal{F}_{\infty}$, for any lattice gauge $Q \in GL_3(\mathbb{C}[t])$, we can find $P', \tilde{Q} \in GL_3(\mathbb{C}[t])$ and a monopole $\Pi \in GL_3(\mathbb{C}[z])$ as in the following diagram:



We will get a BG logarithmic trivialization if (ε) is a basis of the Birkhoff form Υ of \mathfrak{D}_{∞} . Since ∞ is a regular point, the gauge Q is a holomorphic fundamental matrix of solutions of the system (24). Lemma 15 ensures that we can actually truncate Q at order $\Delta T - 1 = 1$; hence we can take

$$Q = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} -1 & 0 & -1 \\ -1 & 0 & 1 \\ 1 & 0 & 1 \end{pmatrix} t.$$

The gauge P' is obtained as in the proof of Proposition 38:

$$P' = \begin{pmatrix} 1 & 0 & -\frac{t}{2} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \text{ so that } \Pi = \begin{pmatrix} 1 & 0 & -\frac{1}{2t} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Finally, we get

$$A_{[P't^T]} = \frac{1}{z}A_0 + \frac{1}{z+1}A_{-1} + \frac{1}{z-1}A_1 + \frac{1}{z-\frac{1}{2}}A_{\frac{1}{2}},$$

with $A_0 = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 1 \\ -1 & 0 & 0 \end{pmatrix}, \quad A_{-1} = \frac{1}{24}\begin{pmatrix} 2 & 0 & -1 \\ 0 & 0 & -24 \\ 4 & 0 & -2 \end{pmatrix},$
 $A_1 = \frac{1}{8}\begin{pmatrix} 2 & 0 & 1 \\ 8 & 0 & -4 \\ -4 & 0 & -2 \end{pmatrix}$ and $A_{\frac{1}{2}} = \frac{1}{3}\begin{pmatrix} 2 & 0 & -1 \\ -6 & 0 & 0 \\ 4 & 0 & -2 \end{pmatrix}.$

This is a Fuchsian linear model of (23) at $x = \infty$, and we check that the eigenvalues of

$$A_0 + A_{-1} + A_1 + A_{\frac{1}{2}} = \begin{pmatrix} 1 & 0 & -\frac{1}{4} \\ 0 & 0 & -\frac{1}{2} \\ 0 & 0 & -1 \end{pmatrix}$$

give indeed the type of \mathfrak{D} . The residues A_s are all nilpotent with maximal rank; hence there is a unique complete flag $F^{(s)}$ which is stable by A_s . In the canonical basis (e) of \mathbb{C}^3 , the Harder–Narasimhan filtration corresponds to the coordinate flag

$$\boldsymbol{H}:(0)\subset \langle e_1\rangle\subset \langle e_1,e_2\rangle\subset \mathbb{C}^3$$

and we have

$$F^{(1)} = F^{\left(\frac{1}{2}\right)} : (0) \subset \langle e_2 \rangle \subset \langle e_2, e_1 + 2e_3 \rangle \subset \mathbb{C}^3,$$

$$F^{(-1)} : (0) \subset \langle e_2 \rangle \subset \langle e_2, e_1 - 2e_3 \rangle \subset \mathbb{C}^3,$$

$$F^{(0)} : (0) \subset \langle e_2 \rangle \subset \langle e_2, e_3 \rangle \subset \mathbb{C}^3.$$

We see that no stable flag under any A_s is transversal to H, which is a necessary condition to be a counterexample (by Proposition 64). However, the condition of Corollary 63 is satisfied (at each $s \in \mathcal{S}$), which means that there is an adjacent weak solution \mathcal{E} with strictly smaller default $i(\mathcal{E}) < i(\mathfrak{D}) = 3$.

5.5.2. Adjacent lattices. In this section, we consider a weak solution $\mathscr{C} \in \operatorname{RH}_{\chi}$. In the following proposition, we describe a procedure which allows to read off at an apparent singularity $x \notin \mathscr{G}$, fixed once and for all, the effect on the weak solution \mathscr{C} of a change of logarithmic adjacent lattice at any singularity $s \in \mathscr{G}$. More precisely, let (σ) be a global basis of a logarithmic BG trivialization of \mathscr{C} at x, and Ω the matrix in Fuchsian form (21) of the connection ∇ in (σ), whose residue at x gives precisely the type of \mathscr{C} . Let M be a logarithmic lattice at s that is adjacent to \mathscr{C}_s . We determine explicitly a gauge transform Π_M which is a monopole at x, such that $\Omega_{[\Pi_M]}$ has again Fuchsian form (21). From its semisimple residue at x we read directly the type of the modified bundle \mathscr{C}^M , equal to the eigenvalues, and the Harder–Narasimhan filtration of \mathscr{C}^M , spanned by the eigenspaces ordered by increasing values. This procedure is completely effective once the connection matrices C_s that relate the *local* residue matrices $L_s = \frac{1}{2i\pi} \log G_s$ in the Birkhoff form at s and the global residue matrices $\widetilde{L}_s = C_s^{-1} L_s C_s$ in the basis (σ) have been determined.

Let M be a lattice at s that is adjacent to \mathscr{C}_s . This lattice is uniquely characterized by its image $W = M/\mathfrak{m}_s \mathscr{C}_s$, that can be seen as a sub- \mathbb{C} -vector space $W \subset Y$. It is logarithmic if and only if W is stable under the map $\operatorname{Res}_s^{\mathscr{C}} \nabla$.

According to Proposition 28, a BG trivialization of \mathscr{C}^M is obtained from a basis of \mathscr{C}_s that simultaneously respects the space W and the flag HN. Moreover, we can choose (ε) in the $\operatorname{GL}_n(\mathbb{C})$ -orbit of (σ) .

Claim 75. There exists a basis (ε) of Y such that $t_s^D(\varepsilon)$ spans a BG trivialization of both \mathscr{E} and \mathscr{E}^M at s.

Claim 76. The matrix $P \in GL_n(\mathbb{C})$ of the basis change from (σ) to (ε) is (-D)-parabolic.

Claim 77. The gauge $t_s^{-D} P t_s^{D} = t_x^{D} P t_x^{-D}$ is a monopole at s and an element of $GL_n(\mathbb{O}_x)$.

$$\begin{array}{c} (\widetilde{\sigma}) : \widetilde{Y} \xrightarrow{t_s^{-D} P t_s^{D}} (\widetilde{\varepsilon}) : Y' \\ t_s^{D} \uparrow & t_s^{D} \uparrow \\ (\sigma) : \mathscr{C}_s \xrightarrow{P} (\varepsilon) : \mathscr{C}_s \xrightarrow{t_s^{T}} (\sigma') : M \\ \downarrow \pi & \downarrow \pi \\ E = \mathscr{C}_s / \mathfrak{m}_s \mathscr{C}_s \xrightarrow{P} W \end{array}$$

Claim 78. The basis (σ') generates M at s and \mathscr{C}_y at $y \neq x$.

Claim 79. The trivial bundle \mathcal{F}' spanned by (σ') is a BG trivialization of \mathcal{E}^M at x.

Claim 80. The gauge transform from (σ) to (σ') is $Pt_s^T = Pt_x^{-T}$.

Claim 81. The Harder–Narasimhan filtration of \mathscr{E}^M is given by the flag of Y' spanned by (σ') according to D - T.

Indeed, the last arrow on the right implies that at x, we have

$$(\tilde{\varepsilon}_x): \mathscr{C}_x = \mathscr{C}_x^M \xrightarrow{t_x^{D-T}} (\sigma'): Y' \text{ where } Y' \subset V \text{ is spanned over } \mathbb{C} \text{ by } (\sigma').$$

Therefore the type of \mathscr{C}^M is, as expected, equal to D - T.

Proposition 82. Assume that $\mathcal{G} \subset \mathbb{C}$ and $x = \infty$. Let $\mathcal{E} \in \operatorname{RH}_{\chi}$ be a weak solution to the Riemann–Hilbert problem. Let the connection ∇ have a matrix Ω of the form (21) in a basis (σ) of a logarithmic BG trivialization \mathcal{F} of \mathcal{E} at x. Then, for any \tilde{L}_s -stable subspace W_s of \mathbb{C}^n , there exists a computable monopole gauge $\Pi \in \operatorname{GL}_n(\mathbb{C}[z])$, a constant matrix $P_0 \in \operatorname{GL}_n(\mathbb{C})$ and a diagonal matrix T with only 0, 1 elements such that $\Omega_{[P_0(z-s)^T \Pi]}$ has again the form (21) corresponding to the modification \mathcal{E}^M , where M is the lattice of \mathcal{V}_s adjacent to \mathcal{E}_s canonically defined by W_s .

Proof. We identify $\Gamma(X, \mathcal{F})$ with \mathbb{C}^n by means of the basis (σ) . The residue of ∇ at *s* is then equal to the matrix $L = \tilde{L}_s$ of formula (21). A logarithmic adjacent lattice *M* is uniquely defined by an *L*-stable subspace $W \subset \mathbb{C}^n$. Let (ε) be a basis respecting both *W* and the Harder–Narasimhan flag *H*, and let $P \in GL_n(\mathbb{C})$ be the basis change from (σ) to (ε) . Assume for simplicity that we have ordered the vectors $\varepsilon_1, \ldots, \varepsilon_n$ in such a way that if $\varepsilon_i \in H_k \cap W$ and $\varepsilon_{i+1} \notin W$ then $\varepsilon_{i+1} \notin H_k$. Let $T = \text{diag}(t_1, \ldots, t_n)$ be the diagonal integer matrix defined by $t_i = 1$ if and only if $\varepsilon_i \notin W$. With the simplifying assumption, the type of \mathscr{C}^M

is equal to D - T, including the ordering condition, and the Harder–Narasimhan filtration is exactly obtained by putting together the groups of vectors corresponding to equal values of D-T. Therefore the basis $(\sigma') = (z-s)^{D-T}(\varepsilon)$ spans a BG trivialization \mathcal{F}' of \mathcal{C}^M at s, and it is simultaneously a global basis of V. The transport $t_x(\mathcal{F}')$ is again a BG trivialization of \mathcal{C}^M at x, but it needs not be logarithmic anymore. Since \mathscr{E} is a weak solution, we have $\mathscr{E}_x = \mathscr{E}_x^M = \mathfrak{D}_x$. Therefore, there exists a lattice gauge transformation $P = I + P_1 t_x + P_2 t_x^2 + \cdots$ which sends the basis (σ') into its Υ -basis (ε'), where Υ is the Birkhoff form at x. The lattice M' spanned by $t_x^{D-T}(\sigma')$ is then necessarily logarithmic, according to Lemma 15. We can effectively determine M' by truncating the gauge P at order $d(M', D) - 1 = k_n - k_1 - 2$, and then applying Gantmacher's classical recursive formulæ (18). Then, the permutation lemma yields a monopole gauge transform Π at x so that the resulting trivialization $\overline{\mathcal{F}}$ is both BG and logarithmic. In this last basis, the connection has again the form (21), where the spectrum of the residue at x gives the type of the modified logarithmic bundle \mathscr{C}^M .

Proposition 56 implies that iterated applications of this procedure describe the set of all weak solutions to the Riemann–Hilbert problem. The strong problem is solvable if and only if one of the bundles $\overline{\mathcal{F}}$ has a 0 residue at x in the orbit under these transformations.

5.5.3. The general case. For an arbitrary weak solution \mathcal{E} , we must start with the Deligne bundle \mathfrak{D} , for we only have the complete description of the local logarithmic lattices from the Deligne lattice. According to the description given in Proposition 52, any logarithmic lattice $N_s \in \Lambda_s$ is given by an admissible pair (F, T) where F is a Res^D_s ∇ -stable flag. If we put ourselves in the situation of Section 5.5.2, and consider a logarithmic BG trivialization \mathcal{F} of \mathfrak{D} at x, and identify $Y = \Gamma(X, \mathcal{F})$ to \mathbb{C}^n by means of the basis (σ), then the flag F can be viewed as a flag in \mathbb{C}^n stable under the matrix \widetilde{L}_s . As exposed in Theorem 65, the bundle \mathscr{C} is then described by an element $(F^s, D_s)_{s \in \mathscr{G}} \in \mathrm{RH}_{\chi}$ such that $F^s \in \mathfrak{Fl}_{\psi_s}(Y)$ and $D_s \in \mathbb{Z}^n(\mathbf{F}^s)$. In order to actually construct the lattice N_s , one should in principle reach first a Birkhoff form Υ_z in $\Lambda = \mathfrak{D}_s$. Since the Deligne lattice is nonresonant, it is possible to do so by a lattice gauge P_s tangent to the identity, as described in Section 4.1.1. We know from Definition 14 that if we put $d_s = \Delta D_s - 1$, the local gauge P_s can be truncated to t_s -degree d_s , as remarked in the proof of Proposition 82. Assume for simplicity that $x = \infty \notin \mathcal{G}$. Taking $Q_s \in GL_n(\mathbb{C})$ that brings (σ) to a basis respecting F^{s} , the local gauge can be written as

$$P_s = Q_s \left(I + \dots + P_{d_s}^{(s)} (z-s)^{d_s} \right).$$

Take a rational interpolation $Z \in GL_n(\mathbb{C}(z))$ of the local gauges P_s to the prescribed orders d_s , and having only a singularity outside \mathcal{G} at x (see, e.g., [van

Barel et al. 1994]). The global basis of sections ($\tilde{\sigma}$) whose matrix in (σ) is *P* spans by construction a trivialization $\tilde{\mathcal{F}}$ of \mathcal{E} at *x*. Everything can be seen at *x* as explained in the following scheme. Putting $\Lambda = \mathfrak{D}_x = \mathcal{E}_x$, we have



The effect of having changed the stalks over \mathscr{G} is translated in a purely local fashion by a change in the set Λ^0 . Indeed, when Λ represented the class of $[\mathfrak{D}]_X$, the germ (σ_X) was a diagonal shift of a global basis. In this scheme however, Λ represents the class $[\mathscr{C}]_X$, and we have now to apply the gauge Z, considered as an element of $\operatorname{GL}_n(\mathbb{C}((t_X)))$, to get a basis of global sections of $\widetilde{Y} = \Gamma(X, \widetilde{\mathscr{F}})$. The monopole Π_1 corresponds to the construction of a BG trivialization, as in Section 3.2.1. The second gauge Π_2 , which can be constructed by the enhanced permutation lemma (Proposition 38), brings the system to an optional logarithmic BG trivialization, where the system is again Fuchsian.

References

[Anosov and Bolibrukh 1994] D. V. Anosov and A. A. Bolibruch, *The Riemann–Hilbert problem*, Aspects of Mathematics **22**, Friedr. Vieweg & Sohn, Braunschweig, 1994. MR 95d:32024 Zbl 0801.34002

[Babbitt and Varadarajan 1983] D. G. Babbitt and V. S. Varadarajan, "Formal reduction theory of meromorphic differential equations: A group theoretic view", *Pacific J. Math.* **109**:1 (1983), 1–80. MR 86b:34010 Zbl 0533.34010

[van Barel et al. 1994] M. van Barel, B. Beckermann, A. Bultheel, and G. Labahn, "Matrix rational interpolation with poles as interpolation points", pp. 137–148 in *Nonlinear numerical methods and rational approximation, II: Proceedings of the 3rd International Conference* (University of Antwerp, Belgium, 5–11 September 1993), edited by A. Cuyt, Mathematics and its Applications **296**, Kluwer Academic Publishers, Dordrecht, 1994. Zbl 0803.41006

[Boalch 2011] P. P. Boalch, "Riemann–Hilbert for tame complex parahoric connections", *Transform. Groups* **16**:1 (2011), 27–50. MR 2012m:14020 Zbl 1232.34117

[[]Abramenko and Brown 2008] P. Abramenko and K. S. Brown, *Buildings: Theory and applications*, Graduate Texts in Mathematics **248**, Springer, New York, 2008. MR 2009g:20055 Zbl 1214.20033

[[]André and Baldassarri 2001] Y. André and F. Baldassarri, *De Rham cohomology of differential modules on algebraic varieties*, Progress in Mathematics **189**, Birkhäuser, Basel, 2001. MR 2002h:14031 Zbl 0995.14003

- [Bolibrukh 1990] A. A. Bolibrukh, "The Riemann–Hilbert problem", *Uspekhi Mat. Nauk* **45**:2(272) (1990), 3–47. In Russian; translated in *Russ. Math. Surv.* **45**:2 (1990), 1–58. MR 92j:14014 Zbl 0706.34005
- [Bolibrukh et al. 2006] A. A. Bolibruch, S. Malek, and C. Mitschi, "On the generalized Riemann– Hilbert problem with irregular singularities", *Expo. Math.* **24**:3 (2006), 235–272. MR 2007i:34148 Zbl 1106.34061
- [Corel 2004] E. Corel, "On Fuchs' relation for linear differential systems", *Compos. Math.* **140**:5 (2004), 1367–1398. MR 2005e:34270 Zbl 1070.34120
- [Crawley-Boevey 2003] W. Crawley-Boevey, "On matrices in prescribed conjugacy classes with no common invariant subspace and sum zero", *Duke Math. J.* 118:2 (2003), 339–352. MR 2004g:16014 Zbl 1046.15013
- [Deligne 1970] P. Deligne, *Équations différentielles à points singuliers réguliers*, Lecture Notes in Mathematics **163**, Springer, Berlin, 1970. MR 54 #5232 Zbl 0244.14004
- [Gantmacher 1959] F. R. Gantmacher, *The theory of matrices, II*, Chelsea, New York, 1959. MR 21 #6372c Zbl 0927.15002
- [Garrett 1997] P. Garrett, *Buildings and classical groups*, Chapman-Hall, London, 1997. MR 98k: 20081 Zbl 0933.20019
- [Gérard and Levelt 1973] R. Gérard and A. H. M. Levelt, "Invariants mesurant l'irrégularité en un point singulier des systèmes d'équations différentielles linéaires", *Ann. Inst. Fourier (Grenoble)* **23**:1 (1973), 157–195. MR 49 #10947 Zbl 0243.35016
- [Ilyashenko and Yakovenko 2007] Y. Ilyashenko and S. Yakovenko, *Lectures on analytic differential equations*, Graduate Studies in Mathematics 86, American Mathematical Society, Providence, RI, 2007. MR 2009b:34001 Zbl 1186.34001
- [Kostov 1992] V. P. Kostov, "Fuchsian linear systems on CP¹ and the Riemann–Hilbert problem", *C. R. Acad. Sci. Paris Sér. I Math.* **315**:2 (1992), 143–148. MR 94a:34007 Zbl 0772.34007
- [Krause 1996] M. Krause, "A simple proof of the Gale–Ryser theorem", *Amer. Math. Monthly* **103**:4 (1996), 335–337. MR 1383671 Zbl 0855.05032
- [Levelt 1961] A. H. M. Levelt, "Hypergeometric functions, II", *Nederl. Akad. Wetensch. Proc. Ser. A* 64 = Indag. Math. 23 (1961), 373–385. MR 25 #1302 Zbl 0124.03602
- [Malgrange 1996] B. Malgrange, "Connexions méromorphes, II: Le réseau canonique", *Invent. Math.* **124**:1-3 (1996), 367–387. MR 97h:32060 Zbl 0849.32003
- [Plemelj 1964] J. Plemelj, *Problems in the sense of Riemann and Klein*, Interscience Tracts in Pure and Applied Mathematics **16**, Wiley, New York, 1964. MR 30 #5008 Zbl 0124.28203
- [Ronan 1989] M. Ronan, *Lectures on buildings*, Perspectives in Mathematics 7, Academic Press, Boston, MA, 1989. MR 90j:20001 Zbl 0694.51001
- [Sabbah 2002] C. Sabbah, *Déformations isomonodromiques et variétés de Frobenius*, EDP Sciences, Paris, 2002. MR 2003m:32013 Zbl 1101.14001
- [Treibich Kohn 1983] A. Treibich Kohn, "Un résultat de Plemelj", pp. 307–312 in *Mathematics and physics* (Paris, 1979/1982), edited by L. Boutet de Monvel et al., Progr. Math. **37**, Birkhäuser, Boston, 1983. MR 87c:32016 Zbl 0527.34007

Received December 3, 2010. Revised March 14, 2013.

ELIE COMPOINT

ELIE COMPOINT AND EDUARDO COREL

UNIVERSITÉ DES SCIENCES ET TECHNOLOGIES - LILLE 1 CITÉ SCIENTIFIQUE 59655 VILLENEUVE D'ASCQ FRANCE elie.compoint@math.univ-lille1.fr

EDUARDO COREL UNIVERSITÉ D'EVRY-VAL-D'ESSONNE IBGBI, 23 BOULEVARD DE FRANCE 91037 EVRY FRANCE

eduardo.corel@genopole.cnrs.fr

ELLIPTIC ALIQUOT CYCLES OF FIXED LENGTH

NATHAN JONES

Silverman and Stange define the notion of an aliquot cycle of length L for a fixed elliptic curve E over \mathbb{Q} , and conjecture an order of magnitude for the function which counts such aliquot cycles. In the present note, we combine heuristics of Lang-Trotter with those of Koblitz to refine their conjecture to a precise asymptotic formula by specifying the appropriate constant. We give a criterion for positivity of the conjectural constant, as well as some numerical evidence for our conjecture.

1. Introduction

Let *E* be an elliptic curve over \mathbb{Q} and fix a positive integer $L \ge 2$. In analogy with the classical notion of an aliquot cycle, Silverman and Stange [2011] define an *L*-tuple (p_1, p_2, \ldots, p_L) of distinct positive integers to be an *aliquot cycle of length L for E* if each p_i is a prime number of good reduction for *E*,

 $p_1 = |E(\mathbb{F}_{p_L})|$ and $p_{i+1} = |E(\mathbb{F}_{p_i})|$ for all $i \in \{1, 2, \dots, L-1\}$,

which may be more succinctly written as

(1)
$$p_{i+1} = |E(\mathbb{F}_{p_i})|$$
 for all $i \in \mathbb{Z}/L\mathbb{Z}$.

When L = 2, an aliquot cycle is also referred to as an *amicable pair for E*. As observed in [Silverman and Stange 2011, Remark 1.5], there is an intimate connection between aliquot cycles for *E* and elliptic divisibility sequences, which relate to generalizations of classical index divisibility questions about Lucas sequences (see also [Gottschlich 2012], which studies some distributional aspects of elliptic divisibility sequences).

It is of interest to know how common such aliquot cycles are, so we presently consider the function which counts aliquot cycles of fixed length for a fixed elliptic curve *E* over \mathbb{Q} . More precisely, define an aliquot cycle (p_1, p_2, \ldots, p_L) to be *normalized* if $p_1 = \min\{p_i : 1 \le i \le L\}$, and then write

 $\pi_{E,L}(x) := \left| \left\{ p_1 \le x : \exists \text{ a normalized aliquot cycle } (p_1, p_2, \dots, p_L) \text{ for } E \right\} \right|.$

Work partially supported by the National Security Agency under grant H98230-12-1-0210. *MSC2010:* 11G05.

Keywords: elliptic curve, aliquot cycle, amicable pair.

NATHAN JONES

The behavior of $\pi_{E,L}(x)$ for large x depends heavily on whether or not E has complex multiplication (CM), as the following conjecture indicates.

Conjecture 1.1 (Silverman–Stange). Let *E* be an elliptic curve over \mathbb{Q} and $L \ge 2$ a fixed integer, and assume that there are infinitely many primes *p* such that $|E(\mathbb{F}_p)|$ is prime. Then, as $x \to \infty$, one has

$$\pi_{E,L}(x) \quad \begin{cases} \asymp \frac{\sqrt{x}}{(\log x)^L} & \text{if } E \text{ has no CM,} \\ \sim A_E \frac{x}{(\log x)^2} & \text{if } E \text{ has CM and } L = 2, \end{cases}$$

where the implied constants in \asymp are both positive and depend only on *E* and *L*, and *A_E* is a positive constant.

Remark 1.2. We may interpret the case L = 1 of (1) as describing primes p_1 for which $p_1 = |E(\mathbb{F}_{p_1})|$. Such primes are called *anomalous* primes and have been considered in [Mazur 1972]. The asymptotic count for anomalous primes up to x is a special case of a conjecture of Lang and Trotter [1976].

Silverman and Stange [2011] focus on the intricacies of the CM case, proving that if *E* has CM, $j_E \neq 0$ and $L \ge 3$, then any normalized aliquot cycle $(p_1, p_2, ..., p_L)$ for *E* must have $p_1 < 5$ (so, in particular, $\pi_{E,L}(x) = O(1)$). The case $j_E = 0$ is apparently more complicated, and no proof is given that $\pi_{E,L}(x) = O(1)$ when $j_E = 0$ and L > 3.

In this note, we refine Conjecture 1.1 to an asymptotic formula in the non-CM case. Heuristics will be developed which lead to the following conjecture.

Conjecture 1.3. Let *E* be an elliptic curve over \mathbb{Q} without complex multiplication and $L \ge 2$ a fixed integer. Then there is a nonnegative real constant $C_{E,L} \ge 0$ (see (5) below) so that, as $x \longrightarrow \infty$,

$$\pi_{E,L}(x) \sim C_{E,L} \int_2^x \frac{1}{2\sqrt{t}(\log t)^L} dt.$$

Remark 1.4. It is possible for the constant $C_{E,L}$ to be zero, in which case the limit $\lim_{x\to\infty} \pi_{E,L}(x)$ is provably finite. Thus, in case $C_{E,L} = 0$, let us interpret the above asymptotic to mean that $\lim_{x\to\infty} \pi_{E,L}(x) < \infty$.

Remark 1.5. By integration by parts, one has

$$\int_{2}^{x} \frac{1}{2\sqrt{t}(\log t)^{L}} dt = \frac{\sqrt{x}}{(\log x)^{L}} + O\left(\frac{\sqrt{x}}{(\log x)^{L+1}}\right).$$

Thus, Conjecture 1.3 is consistent with Conjecture 1.1. In practice, the error term

$$\left| \pi_{E,L}(x) - C_{E,L} \int_{2}^{x} \frac{1}{2\sqrt{t}(\log t)^{L}} \, dt \right|$$

E	$x = 10^{6}$	$x = 10^8$	$x = 10^{10}$	$x = 10^{12}$	$x = 10^{13}$
$E_1: y^2 + y = x^3 - x$	0	1	16	115	332
$E_2: y^2 = x^3 + 6x - 2$	0	5	32	208	564
$E_3: y^2 = x^3 - 3x + 4$	0	0	0	0	0

Table 1.	Values	of $\pi_{E,2}$	(x)
----------	--------	----------------	-----

should be smaller than $\left|\pi_{E,L}(x) - C_{E,L}\frac{\sqrt{x}}{(\log x)^L}\right|$, just as in the case of the prime number theorem.

Consider Table 1, which lists the values of $\pi_{E,2}(x)$ for a few non-CM curves *E* and various magnitudes *x*. Note that $\pi_{E_2,2}(x)$ is larger than $\pi_{E_1,2}(x)$. This difference is explained by the associated constants appearing in Conjecture 1.3. Indeed, a computation shows that

$$\frac{C_{E_2,2}}{C_{E_1,2}} \approx 1.714$$

Also note that $\pi_{E_{3,2}}(10^{13}) = 0$. The additional fact that

 $|\{p \le 10^{12} : p \text{ is of good reduction for } E_3 \text{ and } |E_3(F_p)| \text{ is prime}\}|=715, 698, 540$

indicates that there probably are infinitely many primes p for which $|E_3(\mathbb{F}_p)|$ is prime, in which case the above data suggests that E_3 might be a counterexample to Conjecture 1.1. We will later see that $C_{E_3,2} = 0$, and that E_3 is indeed a counterexample, assuming a conjecture of Koblitz on the primality of $|E(\mathbb{F}_p)|$.

Remark 1.6. The heuristics which lead to Conjecture 1.3 are in the style of Koblitz and Lang–Trotter, whose conjectures have been proven "on average over elliptic curves E" (see [Balog et al. 2011; David and Pappalardi 1999]). It might be interesting to see if one could also prove an average version of Conjecture 1.3.

1.1. *Positivity of* $C_{E,L}$ *and a directed graph* \mathcal{G}_E . In the interest of characterizing the non-CM elliptic curves which have infinitely many aliquot cycles of length L, we will state a graph-theoretic criterion for positivity of $C_{E,L}$. Recall that a *directed graph* \mathcal{G} is a pair $(\mathcal{V}, \mathcal{C})$, where $\mathcal{V} = \mathcal{V}(\mathcal{G})$ is an arbitrary set of *vertices* and $\mathcal{C} = \mathcal{E}(\mathcal{G}) \subseteq \mathcal{V} \times \mathcal{V}$ is a subset of *directed edges*. The sequence of vertices $(v_1, v_2, v_3, \ldots, v_n)$ is a *closed walk of length* n if and only if $(v_i, v_{i+1}) \in \mathcal{C}$ for each $i \in \mathbb{Z}/n\mathbb{Z} = \{1, 2, 3, \ldots, n\}$. Note that closed walks may have repeated vertices. For instance, if $(v, v) \in \mathcal{C}$ for some vertex v (i.e., if \mathcal{G} has a *loop* at a vertex v), then \mathcal{G} has closed walks of any length.

We will associate to an elliptic curve E a directed graph \mathscr{G}_E . First, consider the *n*-th division field $\mathbb{Q}(E[n])$ of E, obtained by adjoining to \mathbb{Q} the *x* and *y*-coordinates

NATHAN JONES

of the *n*-torsion E[n] of a given Weierstrass model of E. The extension $\mathbb{Q}(E[n])$ is Galois over \mathbb{Q} , and once we fix a basis over $\mathbb{Z}/n\mathbb{Z}$ of E[n], we may view

(2)
$$\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q}) \subseteq \operatorname{GL}_2(\mathbb{Z}/n\mathbb{Z}).$$

We will now attach to Gal($\mathbb{Q}(E[n])/\mathbb{Q}$) a directed graph $\mathscr{G}_E(n)$. Viewing Galois automorphisms as 2×2 matrices via (2), the vertex set $\mathscr{V}(n)$ of our graph $\mathscr{G}_E(n)$ is

$$\mathscr{V}(n) := \{(t, d) \in \mathbb{Z}/n\mathbb{Z} \times (\mathbb{Z}/n\mathbb{Z})^{\times} : \exists g \in \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q}) \text{ with } \operatorname{tr} g = t, \det g = d\}.$$

We define the set $\mathscr{E}(n) \subseteq \mathscr{V}(n) \times \mathscr{V}(n)$ of directed edges by declaring that $(v_1, v_2) \in \mathscr{E}(n)$ if and only if $d_1 + 1 - t_1 = d_2$, where $v_i = (t_i, d_i) \in \mathscr{V}(n)$.

Let m_E denote the *torsion conductor* of E, which is defined as the smallest positive integer m for which

$$\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q}) = \pi^{-1} \left(\operatorname{Gal}(\mathbb{Q}(E[\operatorname{gcd}(m, n)])/\mathbb{Q}) \right) \text{ for all } n \in \mathbb{Z}_{>0},$$

where $\pi : \operatorname{GL}_2(\mathbb{Z}/n\mathbb{Z}) \to \operatorname{GL}_2(\mathbb{Z}/\operatorname{gcd}(m, n)\mathbb{Z})$ is the canonical projection. (The existence of a torsion conductor m_E for a non-CM elliptic curve *E* is a celebrated theorem of Serre [1972].) Finally, we define the directed graph \mathscr{G}_E to be the above graph at level m_E :

$$\mathscr{G}_E := \mathscr{G}_E(m_E).$$

The following version of Conjecture 1.3 states a criterion for positivity of $C_{E,L}$ in terms of the directed graph \mathcal{G}_E .

Conjecture 1.7. Let *E* be an elliptic curve over \mathbb{Q} without complex multiplication and $L \ge 2$ a fixed integer. Suppose that the directed graph \mathscr{G}_E has a closed walk of length *L*. Then there are infinitely many aliquot cycles of length *L* for *E*. More precisely, there is a positive constant $C_{E,L} > 0$ so that, as $x \longrightarrow \infty$,

$$\pi_{E,L}(x) \sim C_{E,L} \int_2^x \frac{1}{2\sqrt{t}(\log t)^L} dt.$$

Remark 1.8. If \mathscr{G}_E does not have a closed walk of length *L*, then $C_{E,L} = 0$ and there are at most finitely many aliquot cycles of length *L* for *E* (see Proposition 2.6).

In Section 2, we will write down the constant $C_{E,L}$ explicitly as an "almost Euler product" and discuss its positivity in terms of the graph \mathscr{G}_E . In Section 3, we will develop the heuristics which lead to Conjecture 1.3. In Section 4, we will provide some numerical evidence for Conjecture 1.3 by examining the order of magnitude of $\pi_{E,L}(x) - C_{E,L} \int_2^x \frac{1}{2\sqrt{t}(\log t)^L} dt$ for various elliptic curves *E* and $L \in \{2, 3\}$.

2. The constant

We now describe in detail the constant $C_{E,L}$. The next lemma allows us to interpret (1) in terms of the Frobenius automorphisms¹ $\operatorname{Frob}_{\mathbb{Q}(E[n])}(p_i) \in \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})$ attached to the various primes p_i . Recall the trace of Frobenius $a_p(E) \in \mathbb{Z}$, which satisfies the equation

$$|E(\mathbb{F}_p)| = p + 1 - a_p(E)$$

as well as the Hasse bound

 $|a_p(E)| \le 2\sqrt{p}.$

Lemma 2.1 [Serre 1968, IV-4–IV-5]. For any positive integer n and any prime p of good reduction for E which does not divide n, p is unramified in $\mathbb{Q}(E[n])$ and, for any Frobenius automorphism $\operatorname{Frob}_{\mathbb{Q}(E[n])}(p) \in \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})$, we have

 $\operatorname{tr}(\operatorname{Frob}_{\mathbb{Q}(E[n])}(p)) \equiv a_p(E) \mod n \quad and \quad \det(\operatorname{Frob}_{\mathbb{Q}(E[n])}(p)) \equiv p \mod n.$

For any subset $G \subseteq GL_2(\mathbb{Z}/n\mathbb{Z})$, define

$$G_{\text{ali-cycle}}^L := \left\{ (g_1, g_2, \dots, g_L) \in G^L : \forall i \in \mathbb{Z}/L\mathbb{Z}, \det(g_{i+1}) = \det(g_i) + 1 - \operatorname{tr}(g_i) \right\}.$$

Note that, by Lemma 2.1, if $(p_1, p_2, ..., p_L)$ is an aliquot cycle of length L for E, then

(4)
$$(\operatorname{Frob}_{\mathbb{Q}(E[n])}(p_1), \operatorname{Frob}_{\mathbb{Q}(E[n])}(p_2), \dots, \operatorname{Frob}_{\mathbb{Q}(E[n])}(p_L))$$

 $\in \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^L_{\operatorname{ali-cycle}}.$

Next, let $\phi(x) := \frac{2}{\pi}\sqrt{1-x^2}$ be the distribution function of Sato–Tate, which (assuming *E* has no CM) conjecturally² satisfies

$$\lim_{x \to \infty} \frac{\left| \left\{ p \le x : \frac{a_p(E)}{2\sqrt{p}} \in I \subseteq [-1, 1] \right\} \right|}{|\{p \le x\}|} = \int_I \phi(x) \, dx.$$

In other words, ϕ is the density function of $a_p(E)/2\sqrt{p}$, viewed as a random variable. Denote by $\phi_L := \phi * \phi * \cdots * \phi$ the *L*-fold convolution of ϕ with itself,

¹The Frobenius automorphism in

$$\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})$$

attached to an unramified rational prime p is only defined up to conjugation in $Gal(\mathbb{Q}(E[n])/\mathbb{Q})$. Here and throughout the paper, we understand $\operatorname{Frob}_{\mathbb{Q}(E[n])}(p)$ to be any choice of such a Frobenius automorphism.

²Assuming *E* has nonintegral *j*-invariant, the Sato–Tate conjecture is now a theorem of L. Clozel, M. Harris, N. Shepherd-Barron, and R. Taylor (see [Taylor 2008] and the references therein).

NATHAN JONES

which (again assuming the Sato–Tate conjecture) is the density function of the random variable

$$\sum_{i=1}^{L} \frac{a_{p_i}(E)}{2\sqrt{p_i}},$$

provided the various terms $a_{p_i}(E)/2\sqrt{p_i}$ are "statistically independent." Since the primes p_1, p_2, \ldots, p_L belonging to an aliquot cycle must be close to one another (i.e., within $\approx L\sqrt{t}$ of one another where $p_1 \approx t$, by the Hasse bound (3)), we are really assuming statistical independence *in short intervals* of the various terms $a_{p_i}(E)/2\sqrt{p_i}$. Finally, for a positive integer k, put

$$n_k := \prod_{p \le k} p^k$$

In Section 3, we will develop heuristics which predict Conjecture 1.3, with

(5)
$$C_{E,L} := \frac{\phi_L(0)}{L} \cdot \lim_{k \to \infty} \frac{n_k^L |\operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})_{\operatorname{ali-cycle}}^L|}{|\operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})^L|}$$

2.1. *The constant as a product.* We will presently prove the following proposition, which gives a more explicit expression of $C_{E,L}$ as a convergent Euler product. Recall that m_E denotes the torsion conductor of E, i.e., the smallest positive integer m for which

$$\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q}) = \pi^{-1} \big(\operatorname{Gal}(\mathbb{Q}(E[\operatorname{gcd}(m, n)])/\mathbb{Q}) \big) \quad \text{for all } n \in \mathbb{Z}_{>0},$$

where $\pi : \operatorname{GL}_2(\mathbb{Z}/n\mathbb{Z}) \to \operatorname{GL}_2(\mathbb{Z}/\operatorname{gcd}(m, n)\mathbb{Z})$ is the canonical projection.

Proposition 2.2. For a positive integer k, let $n_k := \prod_{p \le k} p^k$. Then one has

$$\lim_{k \to \infty} \frac{n_k^L |\operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})_{\operatorname{ali-cycle}}^L|}{|\operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})^L|} = \frac{m_E^L |\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})_{\operatorname{ali-cycle}}^L|}{|\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L|} \cdot \prod_{l \nmid m_E} \frac{l^L |\operatorname{GL}_2(\mathbb{F}_l)_{\operatorname{ali-cycle}}^L|}{|\operatorname{GL}_2(\mathbb{F}_l)^L|}.$$

Furthermore,

(6)
$$0 < \frac{l^L \left| \operatorname{GL}_2(\mathbb{F}_l)_{\operatorname{ali-cycle}}^L \right|}{\left| \operatorname{GL}_2(\mathbb{F}_l)^L \right|} = 1 + O_L \left(\frac{1}{l^2} \right),$$

so the infinite product
$$\prod_{l \nmid m_E} \frac{l^L |\operatorname{GL}_2(\mathbb{F}_l)_{\operatorname{ali-cycle}}^L|}{|\operatorname{GL}_2(\mathbb{F}_l)^L|} \text{ converges absolutely}$$

The proof of Proposition 2.2 involves the following two lemmas.

Lemma 2.3. Let n_1 and n_2 be relatively prime positive integers, and pick any subgroups $G_1 \subseteq GL_2(\mathbb{Z}/n_1\mathbb{Z})$ and $G_2 \subseteq GL_2(\mathbb{Z}/n_2\mathbb{Z})$. Then, viewing $G_1 \times G_2 \subseteq GL_2(\mathbb{Z}/n_1n_2\mathbb{Z})$, one has

$$(G_1 \times G_2)^L_{\text{ali-cycle}} = (G_1)^L_{\text{ali-cycle}} \times (G_2)^L_{\text{ali-cycle}}$$

Proof. Let ι : GL₂($\mathbb{Z}/n_1\mathbb{Z}$) × GL₂($\mathbb{Z}/n_2\mathbb{Z}$) → GL₂($\mathbb{Z}/n_1n_2\mathbb{Z}$) be the isomorphism of the Chinese remainder theorem, and set $G := \iota(G_1 \times G_2)$. For each *L*-tuple $(g_i)_i \in G^L$, we have

$$\det g_{i+1} \equiv \det g_i + 1 - \operatorname{tr} g_i \pmod{n_1 n_2} \quad \text{for all } i \in \mathbb{Z}/L\mathbb{Z}$$
$$\iff \begin{cases} \det g_{i+1} \equiv \det g_i + 1 - \operatorname{tr} g_i \pmod{n_1} \\ \det g_{i+1} \equiv \det g_i + 1 - \operatorname{tr} g_i \pmod{n_2} \end{cases} \quad \text{for all } i \in \mathbb{Z}/L\mathbb{Z}.$$

This implies the conclusion of Lemma 2.3.

Lemma 2.4. Let *n* be a positive integer and *n'* any multiple of *n* such that, for every prime number $l, l \mid n' \Rightarrow l \mid n$. Let $\pi : \operatorname{GL}_2(\mathbb{Z}/n'\mathbb{Z}) \to \operatorname{GL}_2(\mathbb{Z}/n\mathbb{Z})$ denote the canonical projection and let $G \subseteq \operatorname{GL}_2(\mathbb{Z}/n\mathbb{Z})$ be any subgroup. Then one has

(7)
$$\frac{(n')^{L} |(\pi^{-1}(G))_{\text{ali-cycle}}^{L}|}{|\pi^{-1}(G)^{L}|} = \frac{n^{L} |G_{\text{ali-cycle}}^{L}|}{|G^{L}|}$$

Proof. By induction, it suffices to check the case n' = ln, where *l* is some prime dividing *n*. In this case, since $|\pi^{-1}(G)| = l^4 |G|$, (7) is equivalent to

(8)
$$\left| (\pi^{-1}(G))_{\text{ali-cycle}}^L \right| = l^{3L} \left| G_{\text{ali-cycle}}^L \right|,$$

which we now show. Fix an element $g = (g_1, g_2, ..., g_L) \in G_{ali-cycle}^L$, and note that any element $g' \in \pi^{-1}(g)$ has the form

$$g' = (g'_1, g'_2, \dots, g'_L) = \left(\tilde{g}_1(I + nA_1), \tilde{g}_2(I + nA_2), \dots, \tilde{g}_L(I + nA_L)\right) \in \pi^{-1}(g),$$

where for each *i*, \tilde{g}_i is any fixed lift to $\operatorname{GL}_2(\mathbb{Z}/\ln\mathbb{Z})$ of g_i , and $A_i \in M_{2\times 2}(\mathbb{F}_l)$ is arbitrary. We will presently determine the exact conditions on the A_i which force $(g'_1, g'_2, \ldots, g'_L) \in (\pi^{-1}(G))^L_{\operatorname{ali-cycle}}$. First, since $(g_1, g_2, \ldots, g_L) \in G^L_{\operatorname{ali-cycle}}$, we must have

(9)
$$g_i \pmod{l} \notin \{0, I\}$$
 for each $i \in \mathbb{Z}/L\mathbb{Z}$

and furthermore, the quantity

$$\gamma_i := \frac{\det \tilde{g}_{i+1} - \det \tilde{g}_i - 1 + \operatorname{tr} \tilde{g}_i}{n} \in \mathbb{F}_i$$

is well-defined. One checks that

(10)
$$\det g'_{i+1} \equiv \det g'_i + 1 - \operatorname{tr} g'_i \mod ln$$
$$\iff \gamma_i \equiv -\det g_{i+1} \cdot \operatorname{tr} A_{i+1} + \det g_i \cdot \operatorname{tr} A_i - \operatorname{tr}(g_i A_i) \mod l.$$

The condition on the right-hand side is (affine) linear in the coefficients of A_{i+1} and A_i . We consider the linear transformation

$$T: \mathbb{F}_l^{4L} \simeq M_{2 \times 2}(\mathbb{F}_l)^L \to \mathbb{F}_l^L,$$

given by

$$(A_i)_{i=1}^L \mapsto \left(-\det g_{i+1} \cdot \operatorname{tr} A_{i+1} + \det g_i \cdot \operatorname{tr} A_i - \operatorname{tr}(g_i A_i) \right)_{i=1}^L$$

In light of (10), the condition (8) will follow from the surjectivity of the above linear transformation, which we now verify. Writing coordinates as

$$g_i =: \begin{pmatrix} x_i & y_i \\ z_i & w_i \end{pmatrix}$$
 and $A_i =: \begin{pmatrix} a_i & b_i \\ c_i & d_i \end{pmatrix}$,

we have

$$T((A_i)) = ((\det g_i - x_i)a_i + (\det g_i - w_i)d_i - y_ic_i - z_ib_i - \det g_{i+1}a_{i+1} - \det g_{i+1}d_{i+1}).$$

By (9), at least one of det $g_i - x_i$, det $g_i - w_i$, y_i and z_i must be nonzero modulo l, and so

$$T(\{0\}\times\cdots\times\{0\}\times M_{2\times 2}(\mathbb{F}_l)\times\{0\}\times\cdots\times\{0\})=\{0\}\times\cdots\times\{0\}\times\mathbb{F}_l\times\{0\}\times\cdots\times\{0\},$$

where the nonzero entries correspond to the same index i. In particular, the linear transformation in question is surjective and we have verified (8), finishing the proof of Lemma 2.4.

Proof of Proposition 2.2. Choose k large enough so that $m_E \mid n_k$, and write $n_k = n_k^{(1)} \cdot n_k^{(2)}$, where $n_k^{(1)}$ is divisible by primes dividing m_E and $gcd(m_E, n_k^{(2)}) = 1$. By definition of m_E , we then have

$$\operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q}) \simeq \pi^{-1} \left(\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q}) \right) \times \prod_{\substack{l^k \parallel n_k \\ l \nmid m_F}} \operatorname{GL}_2(\mathbb{Z}/l^k \mathbb{Z}),$$

where $\pi : \operatorname{GL}_2(\mathbb{Z}/n_k^{(1)}\mathbb{Z}) \to \operatorname{GL}_2(\mathbb{Z}/m_E\mathbb{Z})$ is the canonical projection. By Lemmas 2.3 and 2.4, we have

$$\frac{n_k^L |\operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})_{\operatorname{ali-cycle}}^L|}{|\operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})^L|} = \frac{m_E^L |\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})_{\operatorname{ali-cycle}}^L|}{|\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L|} \cdot \prod_{\substack{l \mid n_k \\ l \nmid m_E}} \frac{l^L |\operatorname{Gal}(\mathbb{Q}(E[l])/\mathbb{Q})_{\operatorname{ali-cycle}}^L|}{|\operatorname{Gal}(\mathbb{Q}(E[l])/\mathbb{Q})^L|}.$$

Taking the limit as $k \to \infty$, we arrive at the product representation of $C_{E,L}$ stated in Proposition 2.2. We leave the verification of (6) as an exercise.

2.2. *Positivity of the constant.* We will now discuss the positivity of $C_{E,L}$. The following corollary of Proposition 2.2 is immediate.

Corollary 2.5. One has

 $C_{E,L} > 0 \iff \operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L_{\operatorname{ali-cycle}} \neq \emptyset.$

We will now prove the following proposition, which allows one to deduce Conjecture 1.7 from Conjecture 1.3.

Proposition 2.6. For any non-CM elliptic curve E over \mathbb{Q} , one has

(11) $C_{E,L} > 0 \iff \mathscr{G}_E$ has a closed walk of length L.

Furthermore, if \mathscr{G}_E has no closed walks of length L, then there are only finitely many aliquot cycles (p_1, p_2, \ldots, p_L) of length L for E.

Proof. First we prove (11). By Corollary 2.5, we are reduced to showing that

(12) $\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})_{\operatorname{ali-cycle}}^L \neq \emptyset \iff \mathscr{G}_E$ has a closed walk of length L.

The mapping

$$Gal(\mathbb{Q}(E[m_E])/\mathbb{Q}) \to \mathcal{V}(\mathcal{G}_E),$$
$$g \mapsto (\operatorname{tr} g, \det g)$$

induces a mapping $\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L_{\operatorname{ali-cycle}} \longrightarrow \{\text{closed walks of length } L \text{ in } \mathcal{G}_E \}$. Thus, if $\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L_{\operatorname{ali-cycle}} \neq \emptyset$ then \mathcal{G}_E has a closed walk of length L. Conversely, suppose \mathcal{G}_E has a closed walk $(v_1, v_2, v_3, \ldots, v_L)$ of length L. Recall that $\mathcal{V} = \mathbb{Z}/m_E\mathbb{Z} \times (\mathbb{Z}/m_E\mathbb{Z})^{\times}$ and write $v_i = (t_i, d_i)$. Choosing any element $g_i \in$ $\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})$ with tr $g_i = t_i$ and det $g_i = d_i$, we have then constructed an element $(g_1, g_2, \ldots, g_L) \in \operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L_{\operatorname{ali-cycle}}$, so $\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L_{\operatorname{ali-cycle}} \neq \emptyset$. By Corollary 2.5, we conclude the proof of (11).

To see why the nonexistence of closed walks of length L in \mathscr{G}_E implies that $\lim_{x\to\infty} \pi_{E,L}(x) < \infty$, note that, by (12), one has $\operatorname{Gal}(\mathbb{Q}(E[m_E])/\mathbb{Q})^L_{\operatorname{ali-cycle}} = \emptyset$. But then (4) implies that $\lim_{x\to\infty} \pi_{E,L}(x) < \infty$, and the proof of Proposition 2.6 is complete.

3. Heuristics

We will construct a probabilistic model in the style of [Koblitz 1988] and [Lang and Trotter 1976]. We shall call the *L*-tuple $(p_1, p_2, ..., p_L)$ of distinct prime numbers an *aliquot sequence of length L for E* if it satisfies

$$p_{i+1} = |E(\mathbb{F}_{p_i})|$$
 for all $i \in \{1, 2, \dots L-1\}$.

Thus, an aliquot cycle of length L is an aliquot sequence of length L which additionally satisfies $p_1 = |E(\mathbb{F}_{p_L})|$. Suppose that (p_1, p_2, \dots, p_L) is an aliquot

NATHAN JONES

sequence of length *L* for *E*. By substituting $p_2 = p_1 + 1 - a_{p_1}(E)$ into the equation $p_3 = p_2 + 1 - a_{p_2}(E)$, one finds that $p_3 = p_1 + 2 - (a_{p_1}(E) + a_{p_2}(E))$, and continuing in this manner one obtains

(13)
$$p_1 = |E(\mathbb{F}_{p_L})| \iff \sum_{j=1}^L a_{p_j}(E) = L.$$

Thus, a given L-tuple $(p_1, p_2, ..., p_L)$ of positive integers is an aliquot cycle of length L for E if and only if the following conditions hold:

 (1_L) the *L*-tuple (p_1, p_2, \ldots, p_L) is an aliquot sequence of length *L* for *E*;

(2_L) one has
$$\sum_{j=1}^{L} a_{p_j}(E) = L$$
.

Consider the following condition, which generalizes condition (2_L) above by replacing *L* with an arbitrary fixed integer *r*:

 $(2'_{L})$ one has $\sum_{j=1}^{L} a_{p_j}(E) = r$.

We now develop the heuristic "probability" that a given *L*-tuple $(p_1, p_2, ..., p_L)$ of positive integers satisfies (1_L) and $(2'_L)$. First, we must gather some notation. Fix a positive integer *n* and elements $a, b \in \mathbb{Z}/n\mathbb{Z}$. For any subset $S \subseteq GL_2(\mathbb{Z}/n\mathbb{Z})$, let

$$S_{\mathcal{N}=a} := \{g \in S : \det(g) + 1 - \operatorname{tr}(g) = a\} = \{g \in S : \det(I - g) = a\},\$$

$$S^{\det=b} := \{g \in S : \det(g) = b\}, \quad S^{\det=b}_{\mathcal{N}=a} := S_{\mathcal{N}=a} \cap S^{\det=b}.$$

Finally, for $L \ge 1$ and $G \subseteq GL_2(\mathbb{Z}/n\mathbb{Z})$, put

$$G_{\text{ali-sequence}}^{L} := \{ (g_1, g_2, \dots, g_L) \in G^L : \\ \text{for all } i \in \{1, 2, \dots, L-1\}, \det(g_{i+1}) = \det(g_i) + 1 - \operatorname{tr}(g_i) \}.$$

Note that when L = 1, the defining conditions become empty and we have $G_{\text{ali-sequence}}^{L=1} = G$. For a general $L \ge 1$, note that any aliquot sequence (p_1, p_2, \dots, p_L) for E will satisfy

$$(\operatorname{Frob}_{\mathbb{Q}(E[n])}(p_1), \operatorname{Frob}_{\mathbb{Q}(E[n])}(p_2), \dots, \operatorname{Frob}_{\mathbb{Q}(E[n])}(p_L))$$

 $\in \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^L_{\operatorname{ali-sequence}}.$

Finally, for a fixed integer r, define

$$G_{\text{ali-sequence}}^{L,\sum \text{tr}=r} := \left\{ (g_1, g_2, \dots, g_L) \in G_{\text{ali-sequence}}^L : \sum_{i=1}^L \text{tr}(g_i) \equiv r \mod n \right\}.$$

We will presently derive an expression for the probability

$$\mathcal{P}_{(1_L),(2'_L)}(t) := \operatorname{Prob}((p_1, p_2, \dots, p_L) \text{ satisfies } (1_L) \text{ and } (2'_L), \text{ given that } p_1 \approx t).$$
Putting $\mathcal{P}_{(1_L)}(t)$ for the probability that (p_1, p_2, \ldots, p_L) satisfies (1_L) above, and $\mathcal{P}_{(2'_L)}^{given}(t)$ for the conditional probability that (p_1, p_2, \ldots, p_L) satisfies $(2'_L)$, given that it satisfies (1_L) , we have

(14)
$$\mathscr{P}_{(1_L),(2'_L)}(t) = \mathscr{P}_{(1_L)}(t) \cdot \mathscr{P}_{(2'_L)}^{\text{given } (1_L)}(t).$$

In Section 3.1 below, we will derive the probability formula

(15)
$$\mathcal{P}_{(1_L)}(t) \approx \frac{n^{L-1} \cdot \left| \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\operatorname{ali-sequence}}^L \right|}{\left| \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^L \right|} \cdot \frac{1}{(\log t)^L}$$

Following this, in Section 3.2, we will derive

(16)
$$\mathcal{P}_{(2'_{L})}^{\text{given }(1_{L})}(t) \approx \phi_{L}\left(\frac{r}{2\sqrt{t}}\right) \frac{n \cdot \left|\text{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\text{ali-sequence}}^{L, \sum \text{tr}=r}\right|}{\left|\text{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\text{ali-sequence}}^{L}\right|} \cdot \frac{1}{2\sqrt{t}}.$$

Before deriving (15) and (16), we will now observe that, taken together, they lead to Conjecture 1.3. Indeed, using (14), (15) and (16), one concludes

$$\mathcal{P}_{(1_L),(2'_L)}(t) \approx \phi_L\left(\frac{r}{2\sqrt{t}}\right) \cdot \frac{n^L \left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\text{ali-sequence}}^{L,\sum \operatorname{tr}=r}\right|}{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^L\right|} \cdot \frac{1}{2\sqrt{t}(\log t)^L}.$$

Just as with (13), one verifies that, for each $(g_1, g_2, ..., g_L) \in GL_2(\mathbb{Z}/n\mathbb{Z})^L_{ali-sequence}$, one has

$$\det(g_L) + 1 - \operatorname{tr}(g_L) = \det g_1 \iff \sum_{i=1}^L \operatorname{tr}(g_i) \equiv L \mod n$$

It follows that $\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\operatorname{ali-cycle}}^{L} = \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\operatorname{ali-sequence}}^{L,\sum \operatorname{tr}=L}$. Thus, putting $r = L, n = n_k$ and taking the limit as $k \to \infty$, one arrives at

$$\mathcal{P}_{(1_L),(2_L)}(t) \approx \phi_L\left(\frac{L}{2\sqrt{t}}\right) \cdot \lim_{k \to \infty} \frac{n_k^L \left| \operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})_{\operatorname{ali-cycle}}^L \right|}{\left| \operatorname{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})^L \right|} \cdot \frac{1}{2\sqrt{t}(\log t)^L}.$$

Thus, using

$$\pi_{E,L}(x) \approx \frac{1}{L} \int_2^x \mathcal{P}_{(1_L),(2_L)}(t) dt$$

one arrives at Conjecture 1.3. The reason for the extra factor of *L* in the denominator above is that $\pi_{E,L}(x)$ counts *normalized* aliquot cycles, whereas the heuristic probabilities above do not take normalization into account. Also, since *L* is fixed, one verifies that the estimation $\phi_L(L/(2\sqrt{t})) \approx \phi_L(0)$ does not affect the asymptotic.

NATHAN JONES

3.1. The probability that $(p_1, p_2, ..., p_L)$ satisfies $(\mathbf{1}_L)$. We will now derive a refined probability formula which implies (15). Fix a vector $\mathbf{a} = (a_2, a_3, ..., a_L) \in ((\mathbb{Z}/n\mathbb{Z})^{\times})^{L-1}$, and consider the probability

$$:= \operatorname{Prob}((p_1, p_2, \dots, p_L) \text{ satisfies } (1_L) \text{ and for all } i \in \{2, 3, \dots, L\}, p_i \equiv a_i \mod n)$$

and (for any subset $G \subseteq GL_2(\mathbb{Z}/n\mathbb{Z})$) the subset

$$G_{\text{ali-sequence}}^{L,a} := \{ (g_1, g_2, \dots, g_L) \in G_{\text{ali-sequence}}^L : \text{ for all } i \in \{2, 3, \dots, L\}, \det(g_i) = a_i \}.$$

In case L = 1, the vector $\mathbf{a} \in ((\mathbb{Z}/n\mathbb{Z})^{\times})^0$ is nonexistent, and as before we interpret the empty condition as $G_{\text{ali-sequence}}^{1, a} = G$. Also note the decomposition

(17)
$$G_{\text{ali-sequence}}^{L,a} = G_{\mathcal{N}=a_2} \times G_{\mathcal{N}=a_3}^{\det=a_2} \times G_{\mathcal{N}=a_4}^{\det=a_3} \times \cdots \times G_{\mathcal{N}=a_L}^{\det=a_{L-1}} \times G^{\det=a_L}$$

Finally, note that if $a_1 \neq a_2$, then $G_{\text{ali-sequence}}^{L, a_1} \cap G_{\text{ali-sequence}}^{L, a_2} = \emptyset$, and so we have a disjoint union

$$G_{\text{ali-sequence}}^{L} = \bigsqcup_{\boldsymbol{a} \in ((\mathbb{Z}/n\mathbb{Z})^{\times})^{L-1}} G_{\text{ali-sequence}}^{L, \boldsymbol{a}}.$$

For similar reasons, we have

$$\mathcal{P}_{(1_L)}(t) = \sum_{\boldsymbol{a} \in ((\mathbb{Z}/n\mathbb{Z})^{\times})^{L-1}} \mathcal{P}_{(1_L)}^{\boldsymbol{a}}(t).$$

Thus, (15) will follow from

(18)
$$\mathcal{P}^{\boldsymbol{a}}_{(1_L)}(t) \approx \frac{n^{L-1} \cdot \left| \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^{L, \boldsymbol{a}}_{\operatorname{ali-sequence}} \right|}{\left| \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^{L} \right|} \cdot \frac{1}{(\log t)^{L}},$$

which we will now derive by induction on L.

Base case: L = 1. Suppose that p_1 is a positive integer of size about t. One may interpret the prime number theorem as the probabilistic statement that

$$\mathcal{P}_{(1_{L=1})}(t) = \operatorname{Prob}(p_1 \text{ is prime}) \approx \frac{1}{\log t},$$

which is base case L = 1 of (18).

Induction step. Assume now that (18) holds for some fixed $L \ge 1$, and fix any vector $\boldsymbol{a} = (a_2, a_3, \dots, a_{L+1}) \in ((\mathbb{Z}/n\mathbb{Z})^{\times})^L$. Since the statement

 $(p_1, p_2, \dots, p_{L+1})$ satisfies (1_{L+1}) and for all $i \in \{2, 3, \dots, L+1\}$, $p_i \equiv a_i \mod n$

 $\mathcal{O}^{a}(t)$

is equivalent to

$$(p_1, p_2, ..., p_L)$$
 satisfies (1_L) and for all $i \in \{2, 3, ..., L\}$, $p_i \equiv a_i \mod n$,
 $p_{L+1} := p_L + 1 - a_{p_L}(E)$ is prime, and $p_{L+1} \equiv a_{L+1} \mod n$,

we see that

(19)
$$\mathcal{P}_{(1_{L+1})}^{(a_2, a_3, \dots, a_L, a_{L+1})}(t) = \mathcal{P}_{(1_L)}^{(a_2, a_3, \dots, a_L)}(t) \cdot \mathcal{P}(t),$$

where $\mathcal{P}(t)$ is the conditional probability that $p_{L+1} := p_L + 1 - a_{p_L}(E)$ is prime, and that $p_{L+1} \equiv a_{L+1} \mod n$, given that (1_L) holds. To estimate $\mathcal{P}(t)$, let us assume that (1_L) holds. First note that, by the Hasse bound $|a_p(E)| \le 2\sqrt{p}$, one has

$$p_{L+1} = p_1 + L - \sum_{i=1}^{L} a_{p_i}(E) \in \left[p_1 + L - 2L\sqrt{p_{\max}}, p_1 + L + 2L\sqrt{p_{\max}} \right],$$

where $p_{\text{max}} := \max\{p_i : i = 1, 2, ..., L\}$. By induction we have $p_{\text{max}} = t + O_L(\sqrt{t})$, and so $p_{L+1} \approx t$, with an error of $O_L(\sqrt{t})$. Now, if p_{L+1} were a positive integer of size about t selected independently of $(p_1, p_2, ..., p_L)$, then

(20) Prob
$$(p_{L+1} \text{ is prime and } p_{L+1} \equiv a_{L+1} \mod n) \approx \frac{1}{\varphi(n) \log t}$$

by the prime number theorem in arithmetic progressions. If the positive integer p_{L+1} were chosen randomly and independently of the previous primes, then the probability that $p_{L+1} \equiv a_{L+1} \mod n$ would be 1/n. However, p_{L+1} is not chosen independently of (p_1, p_2, \ldots, p_L) ; it is related to p_L by the formula $p_{L+1} = p_L + 1 - a_{p_L}(E)$. Thus, the congruence $p_{L+1} \equiv a_{L+1} \mod n$ is really the demand that

$$\operatorname{Frob}_{\mathbb{Q}(E[n])}(p_L) \in \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\mathcal{N}=a_{L+1}}$$

Since we assume that (1_L) holds, we know that $\operatorname{Frob}_{\mathbb{Q}(E[n])}(p_L) \in \operatorname{GL}_2(\mathbb{Z}/n\mathbb{Z})^{\det = a_L}$. It is thus natural to multiply (20) by the correction factor

$$\frac{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\mathcal{N}=a_{L+1}}^{\det=a_{L}}\right| / \left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^{\det=a_{L}}\right|}{1/n},$$

obtaining

(21)
$$\mathcal{P}(t) \approx \frac{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\mathcal{N}=a_{L+1}}^{\det=a_{L}}\right| / \left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^{\det=a_{L}}\right|}{1/n} \cdot \frac{1}{\varphi(n)\log t}$$
$$= \frac{n \left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\mathcal{N}=a_{L+1}}^{\det=a_{L}}\right|}{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})\right|} \cdot \frac{1}{\log t}.$$

By (17), we may rewrite (18) as

$$\mathcal{P}_{(1_L)}^{a}(t) \approx n^{L-1} \cdot \frac{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\mathcal{N}=a_2}\right|}{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})\right|} \cdot \left(\prod_{i=2}^{L-1} \frac{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\mathcal{N}=a_{i+1}}^{\det a_i}\right|}{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})\right|}\right) \\ \cdot \frac{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})\right|}{\left|\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})\right|} \cdot \frac{1}{(\log t)^L}$$

Plugging this expression and (21) into (19), and using the fact that

$$\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^{\det=a_L} = \left| \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^{\det=a_{L+1}} \right|$$

one concludes the induction step, completing the derivation of (18), and thus of (15).

Our analysis has motivated the following conjecture, wherein

$$\pi_{E,L}^{\text{ali-sequence}}(x) := \left| \left\{ p_1 \le x : \exists \text{ an aliquot sequence } (p_1, p_2, \dots, p_L) \text{ for } E \right\} \right|,$$
$$C_{E,L}^{\text{ali-sequence}} := \lim_{k \to \infty} \frac{n_k^{L-1} \cdot \left| \text{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})_{\text{ali-sequence}}^L \right|}{\left| \text{Gal}(\mathbb{Q}(E[n_k])/\mathbb{Q})^L \right|}.$$

Conjecture 3.1. Let *E* be an elliptic curve over \mathbb{Q} without complex multiplication and $L \ge 2$ a fixed integer. Then, as $x \longrightarrow \infty$, one has

$$\pi_{E,L}^{\text{ali-sequence}}(x) \sim C_{E,L}^{\text{ali-sequence}} \int_2^x \frac{1}{(\log t)^L} dt$$

Similarly to Proposition 2.6, one has

 $C_{E,L}^{\text{ali-sequence}} > 0 \iff \mathcal{G}_E$ has a (directed) walk of length L.

3.2. The conditional probability that $(p_1, p_2, ..., p_L)$ satisfies $(2'_L)$. We will now derive (16), completing the heuristic derivation of Conjecture 1.3. Suppose that $(p_1, p_2, ..., p_L)$ is an aliquot sequence of length *L* for *E*, i.e., that it satisfies (1_L) . What is the conditional probability that $\sum_{i=1}^{L} a_{p_i}(E) = r$? In the case L = 1, condition (1_L) is empty, and our question becomes identical to the Lang–Trotter conjecture for fixed Frobenius trace. In what follows, we will develop a probabilistic model in the same style as theirs.

Fixing a level *n*, the number $f_n(r, p) \ge 0$ will estimate the probability of the event that $\sum_{i=1}^{L} a_{p_i}(E) = r$, given that $(p = p_1, p_2, ..., p_L)$ is an aliquot sequence of length *L* for *E*. We will model the situation by assuming that the vector

(22)
$$\left(\operatorname{Frob}_{\mathbb{Q}(E[n])}(p_1), \operatorname{Frob}_{\mathbb{Q}(E[n])}(p_2), \dots \operatorname{Frob}_{\mathbb{Q}(E[n])}(p_L)\right) \in \operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})^L_{\operatorname{ali-sequence}}$$

is randomly distributed according to counting measure, and we will assume that the various $a_{p_i}(E)/(2\sqrt{p_i})$ are independent at infinity, i.e., that ϕ_L is the distribution

function for their sum. We will also assume independence of the random variables $\sum_{i=1}^{L} a_{p_i}(E)/(2\sqrt{p_i})$ and (22). Finally, in order to simplify our model, we will also regard all of the various primes p_i as having the same size, namely p. These considerations lead us to the following assumptions about the probabilities $f_n(r, p)$:

$$f_n(r, p) = 0 \quad \text{if } |r| > 2L\sqrt{p},$$

$$f_n(r, p) = \phi_L\left(\frac{r}{2\sqrt{p}}\right) \cdot \frac{n \left|\text{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\text{ali-sequence}}^{L,\sum \text{tr}=r}\right|}{\left|\text{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})_{\text{ali-sequence}}^{L}\right|} \cdot c_p \quad \text{if } |r| \le 2L\sqrt{p},$$

where c_p is some constant chosen so that $\sum_{r \in \mathbb{Z}} f_n(r, p) = 1$. Then, similarly to [Lang and Trotter 1976, pp. 31–32], one concludes that $c_p \sim \frac{1}{2\sqrt{p}}$, as $p \to \infty$. This leads to (16), completing the derivation of Conjecture 1.3.

4. Examples

We will now give some numerical evidence for Conjecture 1.3.

4.1. *Elliptic curves with* $C_{E,L} > 0$. Table 2 and Table 3 display some data for four elliptic curves. In each table, the column labeled "predicted" lists the approximate values of

$$C_{E,L} \int_{2}^{10^{13}} \frac{dt}{2\sqrt{t}(\log t)^{L}}$$

"actual" lists the values of $\pi_{E,L}(10^{13})$, and "% error" lists as a percentage the approximate values of

$$\frac{C_{E,L} \int_{2}^{10^{13}} \frac{dt}{2\sqrt{t}(\log t)^{L}} - \pi_{E,L}(10^{13})}{C_{E,L} \int_{2}^{10^{13}} \frac{dt}{2\sqrt{t}(\log t)^{L}}}$$

The first and third curves were already considered in [Silverman and Stange 2011], and are included here largely to show the contrast with the second curve. For each of these curves, a detailed list of the aliquot cycles with $p_1 \le 10^{13}$ may be found in an expanded version of this paper [Jones 2012].

E	predicted	actual	% error
$y^2 + y = x^3 - x$	318.98	332	-4.08%
$y^2 = x^3 + 6x - 2$	546.78	564	-2.97%
$y^2 + y = x^3 + x^2$	318.97	328	-2.83%
$y^2 + xy + y = x^3 - x^2$	318.95	331	-3.78%

Table 2. Data on $\pi_{E,2}(10^{13})$ for various *E*.

E	predicted	actual	% error
$y^2 + y = x^3 - x$	3.03	3	1.05%
$y^2 = x^3 + 6x - 2$	12.59	12	4.66%
$y^2 + y = x^3 + x^2$	3.04	2	34.10%
$y^2 + xy + y = x^3 - x^2$	3.02	4	-32.48%

Table 3. Data on $\pi_{E,3}(10^{13})$ for various *E*.

The four elliptic curves E under consideration satisfy

(23)
$$\left[\operatorname{GL}_2(\mathbb{Z}/n\mathbb{Z}):\operatorname{Gal}(\mathbb{Q}(E[n])/\mathbb{Q})\right] \le 2$$

for each $n \ge 1$ (see [Serre 1972, pp. 309–311; Lang and Trotter 1976, p. 51]). As shown in [Serre 1972, pp. 310–311], this is the smallest index that one can have for general *n* when the elliptic curve *E* is defined over \mathbb{Q} . We call any elliptic curve *E* satisfying (23) a *Serre curve*. Serre curves are thus elliptic curves for which Gal($\mathbb{Q}(E[n])/\mathbb{Q}$) is "as large as possible for all *n*," and it has been shown that, when ordered by height, almost all elliptic curves are Serre curves (see [Jones 2010; Radhakrishnan 2008]). One can show that for any Serre curve *E*, one has $C_{E,L} > 0$. In fact, if we define the constant C_L by

$$C_L := \frac{\phi_L(0)}{L} \cdot \lim_{k \to \infty} \frac{n_k^L \left| \operatorname{GL}_2(\mathbb{Z}/n_k \mathbb{Z})_{\text{ali-cycle}}^L \right|}{\left| \operatorname{GL}_2(\mathbb{Z}/n_k \mathbb{Z})^L \right|} = \frac{\phi_L(0)}{L} \cdot \prod_{l \text{ prime}} \frac{l^L \left| \operatorname{GL}_2(\mathbb{F}_l)_{\text{ali-cycle}}^L \right|}{\left| \operatorname{GL}_2(\mathbb{F}_l)^L \right|},$$

then for any Serre curve *E* one has that $C_{E,L} = C_L \cdot f_L(\Delta_{sf}(E))$, where $\Delta_{sf}(E)$ denotes the square-free part of the discriminant of any Weierstrass model of *E* and f_L is a positive function which approaches 1 as $|\Delta_{sf}(E)|$ approaches infinity. For L = 2 one has

$$C_{2} = \frac{\phi_{2}(0)}{2} \cdot \prod_{l \text{ prime}} \frac{l^{2} |\mathrm{GL}_{2}(\mathbb{F}_{l})_{\mathrm{ali-cycle}}^{2}|}{|\mathrm{GL}_{2}(\mathbb{F}_{l})^{2}|} = \frac{8}{3\pi^{2}} \cdot \prod_{l \text{ prime}} \frac{l^{2} (l^{4} - 2l^{3} - 2l^{2} + 3l + 3)}{[(l^{2} - 1)(l - 1)]^{2}}$$

\$\approx 0.077088124,

whereas for L = 3 one has

$$C_{3} = \frac{\phi_{3}(0)}{3} \prod_{l \text{ prime}} \frac{l^{3} |\mathrm{GL}_{2}(\mathbb{F}_{l})_{\mathrm{ali-cycle}}^{3}|}{|\mathrm{GL}_{2}(\mathbb{F}_{l})^{3}|}$$

= $\frac{\phi_{3}(0)}{3} \prod_{l \text{ prime}} \frac{l^{3} [l^{6} - 3l^{5} - 3l^{4} + 14l^{3} + (3 + \chi(l))l^{2} - (19 + 3\chi(l))l - 10 - 3\chi(l)]}{[(l^{2} - 1)(l - 1)]^{3}}$
 $\approx 0.019759298,$

E	$C_{E,2}$	$C_{E,3}$	$\Delta_{sf}(E)$
$y^2 + y = x^3 - x$	≈ 0.077093	≈ 0.019841	37
$y^2 = x^3 + 6x - 2$	≈ 0.132151	pprox 0.082365	-3
$y^2 + y = x^3 + x^2$	pprox 0.077091	pprox 0.019861	-43
$y^2 + xy + y = x^3 - x^2$	pprox 0.077088	pprox 0.019759	-53

Table 4. Values of $C_{E,2}$, $C_{E,3}$ and $\Delta_{sf}(E)$.

where $\chi(l) = \left(\frac{-3}{l}\right)$ denotes the character of conductor 3. Table 4 gives the values of $C_{E,2}$, $C_{E,3}$ and $\Delta_{sf}(E)$ for each of the four curves under consideration. The reason the second curve has a larger value of $C_{E,L}$ is that $|\Delta_{sf}(E)|$ is smaller for this curve than for the others.

4.2. An elliptic curve with $C_{E,L} = 0$. We will now discuss briefly the elliptic curve

(24)
$$E : y^2 = x^3 - 3x + 4$$

which was mentioned in the introduction, for which $\pi_{E,L}(x) \equiv 0$ and whose associated graph \mathscr{G}_E contains no closed walks at all. We will presently describe the Galois group Gal($\mathbb{Q}(E[4])/\mathbb{Q}$), which is an index 4 subgroup of GL₂($\mathbb{Z}/4\mathbb{Z}$). First, define the subgroup $H(4) \subseteq GL_2(\mathbb{Z}/4\mathbb{Z})$ by

$$H(4) := \left\{ \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ -1 & -1 \end{pmatrix}, \begin{pmatrix} -1 & -1 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} -1 & -1 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ -1 & -1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \right\}.$$

We then have

(25)
$$\operatorname{Gal}(\mathbb{Q}(E[4])/\mathbb{Q}) = H(4) \cdot \left(I + 2\left\{ \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \right\} \right).$$

(To see that the right-hand expression defines a subgroup of $GL_2(\mathbb{Z}/4\mathbb{Z})$, note that

$$\left\{ \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \right\} \subseteq M_{2 \times 2}(\mathbb{Z}/2\mathbb{Z})$$

is closed under addition and under $GL_2(\mathbb{Z}/2\mathbb{Z})$ -conjugation.)

Even though Gal($\mathbb{Q}(E[2])/\mathbb{Q}$) = GL₂($\mathbb{Z}/2\mathbb{Z}$), Gal($\mathbb{Q}(E[4])/\mathbb{Q}$) is a proper subgroup of GL₂($\mathbb{Z}/4\mathbb{Z}$), and so one has 4 | m_E . Furthermore, in this case the restriction map Gal($\mathbb{Q}(E[m_E])/\mathbb{Q}$) \rightarrow Gal($\mathbb{Q}(E[4])/\mathbb{Q}$) induces a graph morphism

(26)
$$\mathscr{G}_E = \mathscr{G}_E(m_E) \twoheadrightarrow \mathscr{G}_E(4),$$

which is surjective in the sense that it carries the vertex set $\mathcal{V}(m_E)$ onto $\mathcal{V}(4)$ and likewise carries $\mathscr{E}(m_E)$ onto $\mathscr{E}(4)$.

On the other hand, using (25), one finds that the directed graph $\mathcal{G}_E(4)$ is:

NATHAN JONES

Infinitely many primes p for which $|E(\mathbb{F}_p)|$ is prime. The non-CM case of a conjecture of Koblitz (see [Koblitz 1988] and also [Zywina 2011]) expresses (in our terminology) that for any non-CM elliptic curve E, the existence of a single directed edge in \mathcal{G}_E implies the existence of infinitely many primes p for which $|E(\mathbb{F}_p)|$ is prime. Taking E to be the elliptic curve given by (24), we see by the surjectivity of (26) together with (27) that \mathcal{G}_E contains at least one directed edge. Thus, assuming Koblitz's conjecture, there are infinitely many primes p for which $|E(\mathbb{F}_p)|$ is prime.

Finitely many aliquot cycles for E. Continuing with the example (24), by the surjectivity of (26) together with (27), we see that \mathscr{G}_E contains no closed walks at all. By Proposition 2.6, there are only finitely many aliquot cycles (p_1, p_2, \ldots, p_L) for *E*. This particular example may be explained as follows. Whenever $p_2 = |E(\mathbb{F}_{p_1})|$ for some prime p_1 , we see from (27) that $(tr(\operatorname{Frob}_{\mathbb{Q}(E[4])}(p_1)), \det(\operatorname{Frob}_{\mathbb{Q}(E[4])}(p_1))) = (-1, 1)$ (otherwise, $|E(\mathbb{F}_{p_1})|$ would be even). But then

 $(\operatorname{tr}(\operatorname{Frob}_{\mathbb{Q}(E[4])}(p_2)), \operatorname{det}(\operatorname{Frob}_{\mathbb{Q}(E[4])}(p_2))) \in \{(0, -1), (2, -1)\},\$

in which case $|E(\mathbb{F}_{p_2})|$ must be even. One deduces that *E* has no aliquot cycles of length $L \ge 2$, and indeed no aliquot *sequences* of length $L \ge 3$.

Remark 4.1. There is a modular curve X of level 4 and genus 0 with $|X(\mathbb{Q})| = \infty$, whose noncuspidal \mathbb{Q} -rational points correspond to elliptic curves E' for which $-\Delta_{E'}$ is a perfect square. For almost all such elliptic curves E', one may find an appropriate twist E of E' for which (25) holds, and thus for which $\lim_{x\to\infty} \pi_{E,L}(x) < \infty$ for $L \ge 2$. The elliptic curve (24) is one such example.

Acknowledgments

The author gratefully acknowledges A. C. Cojocaru, who first brought this question to my attention, and also J. Silverman for a stimulating discussion. Also many thanks to A. Sutherland, who provided help with the computations (a description of the software used therein may be found in [Kedlaya and Sutherland 2008]). Finally, thanks to the anonymous referee for a careful reading of the manuscript and several helpful comments.

References

- [Balog et al. 2011] A. Balog, A.-C. Cojocaru, and C. David, "Average twin prime conjecture for elliptic curves", *Amer. J. Math.* **133**:5 (2011), 1179–1229. MR 2012j:11118 Zbl 05969056
- [David and Pappalardi 1999] C. David and F. Pappalardi, "Average Frobenius distributions of elliptic curves", *Internat. Math. Res. Notices* 4 (1999), 165–183. MR 2000g:11045 Zbl 0934.11033

- [Gottschlich 2012] A. Gottschlich, "On positive integers *n* dividing the *n*th term of an elliptic divisibility sequence", *New York J. Math.* **18** (2012), 409–420. MR 2928585 Zbl 06098855
- [Jones 2010] N. Jones, "Almost all elliptic curves are Serre curves", *Trans. Amer. Math. Soc.* **362**:3 (2010), 1547–1570. MR 2011d:11130 Zbl 1204.11088
- [Jones 2012] N. Jones, "Elliptic aliquot cycles of fixed length", preprint, 2012. arXiv 1212.1010
- [Kedlaya and Sutherland 2008] K. S. Kedlaya and A. V. Sutherland, "Computing L-series of hyperelliptic curves", pp. 312–326 in *Algorithmic number theory*, edited by A. J. van der Poorten and A. Stein, Lecture Notes in Comput. Sci. **5011**, Springer, Berlin, 2008. MR 2010d:11070 Zbl 1232.11078
- [Koblitz 1988] N. Koblitz, "Primality of the number of points on an elliptic curve over a finite field", *Pacific J. Math.* **131**:1 (1988), 157–165. MR 89h:11023 Zbl 0608.10010
- [Lang and Trotter 1976] S. Lang and H. Trotter, *Frobenius distributions in* GL₂-*extensions: Distribution of Frobenius automorphisms in* GL₂-*extensions of the rational numbers*, Lecture Notes in Mathematics **504**, Springer, Berlin, 1976. MR 58 #27900 Zbl 0329.12015
- [Mazur 1972] B. Mazur, "Rational points of abelian varieties with values in towers of number fields", *Invent. Math.* 18 (1972), 183–266. MR 56 #3020 Zbl 0245.14015
- [Radhakrishnan 2008] V. Radhakrishnan, *An asymptotic formula for the number of non–Serre curves in a two-parameter family of elliptic curves*, Ph.D. thesis, University of Colorado at Boulder, 2008, Available at http://search.proquest.com/docview/304629348. MR 2711545
- [Serre 1968] J.-P. Serre, *Abelian l-adic representations and elliptic curves*, W. A. Benjamin, New York, 1968. MR 41 #8422 Zbl 0186.25701
- [Serre 1972] J.-P. Serre, "Propriétés galoisiennes des points d'ordre fini des courbes elliptiques", *Invent. Math.* **15**:4 (1972), 259–331. In French. MR 52 #8126 Zbl 0235.14012
- [Silverman and Stange 2011] J. H. Silverman and K. E. Stange, "Amicable pairs and aliquot cycles for elliptic curves", *Exp. Math.* **20**:3 (2011), 329–357. MR 2012g:11109
- [Taylor 2008] R. Taylor, "Automorphy for some *l*-adic lifts of automorphic mod *l* Galois representations, II", *Publ. Math. Inst. Hautes Études Sci.* 108 (2008), 183–239. MR 2010j:11085 Zbl 1169.11021
- [Zywina 2011] D. Zywina, "A refinement of Koblitz's conjecture", *Int. J. Number Theory* **7**:3 (2011), 739–769. MR 2012e:11107 Zbl 05913798

Received July 16, 2012. Revised December 10, 2012.

NATHAN JONES DEPARTMENT OF MATHEMATICS UNIVERSITY OF MISSISSIPPI HUME HALL 305 P.O. BOX 1848 UNIVERSITY, MS 38677 UNITED STATES

ncjones@olemiss.edu

ASYMPTOTIC L⁴ NORM OF POLYNOMIALS DERIVED FROM CHARACTERS

DANIEL J. KATZ

Littlewood investigated polynomials with coefficients in $\{-1, 1\}$ (Littlewood polynomials), to see how small their ratio of norms $||f||_4/||f||_2$ on the unit circle can become as deg $f \to \infty$. A small limit is equivalent to slow growth in the mean square autocorrelation of the associated binary sequences of coefficients of the polynomials. The autocorrelation problem for arrays and higher dimensional objects has also been studied; it is the natural generalization to multivariable polynomials. Here we find, for each n > 1, a family of *n*-variable Littlewood polynomials with lower asymptotic $||f||_4/||f||_2$ than any known hitherto. We discover these through a wide survey, infeasible with previous methods, of polynomials whose coefficients come from finite field characters. This is the first time that the lowest known asymptotic ratio of norms $||f||_4/||f||_2$ for multivariable polynomials $f(z_1, \ldots, z_n)$ is strictly less than what could be obtained by using products $f_1(z_1) \cdots f_n(z_n)$ of the best known univariate polynomials.

1. Introduction

1A. *History and main result.* Littlewood [1966; 1968] pioneered the study of the L^4 norm on the complex unit circle of polynomials whose coefficients lie in $\{-1, 1\}$, and in particular wanted to know how small their ratio of norms $||f||_4/||f||_2$ can become as deg $f \to \infty$. He suspected, based on calculations of Swinnerton-Dyer, that this ratio could be made to approach 1 asymptotically, but the smallest limiting ratio he could find was $\sqrt[4]{4/3}$ for the Rudin–Shapiro polynomials [Littlewood 1968]. The L^4 norm is of particular interest since it serves as a lower bound for the L^{∞} norm and is easier to calculate than most other L^r norms. Erdős [1957, Problem 22; 1962] had conjectured that $||f||_{\infty}/||f||_2$ is bounded away from 1 for nonconstant polynomials with complex coefficients of unit magnitude. This was disproved in [Kahane 1980], but the modified problem where we restrict to coefficients in $\{-1, 1\}$ remains open [Newman and Byrnes 1990], and would be solved if one could prove

The author was supported by funding from an NSERC grant awarded to Jonathan Jedwab. *MSC2010:* primary 11C08; secondary 11T24, 42A05, 11B83.

Keywords: L^4 norm, Littlewood polynomial, character polynomial, Fekete polynomial, character sum.

DANIEL J. KATZ

that $||f||_4/||f||_2$ is bounded away from 1 as deg $f \to \infty$. Polynomials in one or more variables with coefficients in $\{-1, 1\}$ and small $||f||_4/||f||_2$ are equivalent to binary sequences and arrays (that is, those that simply list the coefficients of the polynomials) with low mean square aperiodic autocorrelation. Such sequences and arrays are important in the theory of communications¹ [Golay 1977] and statistical physics [Bernasconi 1987]. Accordingly, we define a *Littlewood polynomial* in *n* variables to have the form

$$f(z_1,\ldots,z_n) = \sum_{j_1=0}^{s_1-1} \cdots \sum_{j_n=0}^{s_n-1} f_{j_1,\ldots,j_n} z_1^{j_1} \cdots z_n^{j_n},$$

with coefficient f_{j_1,\ldots,j_n} in $\{-1, 1\}$, and our L^r norm for $f(z_1,\ldots,z_n)$ is

$$\|f\|_r = \left(\frac{1}{(2\pi)^n} \int_0^{2\pi} \cdots \int_0^{2\pi} \left|f\left(\exp(i\theta_1), \ldots, \exp(i\theta_n)\right)\right|^r d\theta_1 \cdots d\theta_n\right)^{1/r}.$$

Note that $||f||_2^2 = s_1 \cdots s_n$ for our Littlewood polynomial.

For univariate Littlewood polynomials, the lowest asymptotic ratio of norms $||f||_4/||f||_2$ found by Littlewood himself [1968] was $\sqrt[4]{4/3}$ for the Rudin–Shapiro polynomials. Two decades later, this was improved to $\sqrt[4]{7/6}$ by Høholdt and Jensen [1988], using modifications of Fekete polynomials. Over two decades later still, in [Jedwab et al. 2013b], another modification was shown to yield further improvement:

Theorem 1.1 (Jedwab, Katz, Schmidt). There is a family of univariate Littlewood polynomials that, as deg $f \to \infty$, has $||f||_4/||f||_2 \to B_1$, the largest real root of $27x^{12} - 498x^8 + 1164x^4 - 722$, which is less than $\sqrt[4]{22/19}$.

Prior to this paper, for each *n*, the lowest known asymptotic $||f||_4/||f||_2$ for *n*-variable Littlewood polynomials

$$f(z_1, \ldots, z_n)$$
 (in the limit as $\deg_{z_1} f, \ldots, \deg_{z_n} f \to \infty$)

was simply the *n*-th power of the lowest known ratio for univariate polynomials, based on the fact that if $f(z_1, \ldots, z_n) = f_1(z_1) \cdots f_n(z_n)$, then $||f||_r =$ $||f_1||_r \cdots ||f_n||_r$. For bivariate Littlewood polynomials, Schmidt [2011] obtained an asymptotic $||f||_4/||f||_2$ of $\sqrt{7/6}$ in this way (via Høholdt and Jensen's univariate polynomials mentioned above), and foresaw the possibility that the asymptotic $||f||_4/||f||_2$ could be lowered to B_1^2 , contingent upon the conjecture that was later established as Theorem 1.1. In this paper, we show that one can do better than this product construction, even when based on the best univariate polynomials now known (those of Theorem 1.1).

¹In this milieu, results are expressed in terms of the *merit factor*, defined as $||f||_2^4/(||f||_4^4 - ||f||_2^4)$.

Theorem 1.2. For each n > 1, there is a family of n-variable Littlewood polynomials $f(z_1, \ldots, z_n)$ which, as $\deg_{z_1} f, \ldots, \deg_{z_n} f \to \infty$, has $||f||_4/||f||_2$ tending to a value strictly less than B_1^n .

The lowest known asymptotic $||f||_4/||f||_2$ for *n*-variable Littlewood polynomials is an algebraic number depending on *n*, and is specified precisely in Section 1C after we define in Section 1B the polynomials that are involved.²

1B. *Character polynomials.* The polynomials used in [Høholdt and Jensen 1988], [Jedwab et al. 2013b], and the current paper to break previous records for the lowest known asymptotic $||f||_4/||f||_2$ are all character polynomials, that is, polynomials whose coefficients are given by characters of finite fields. The L^4 norms of character polynomials have already been studied extensively [Høholdt and Jensen 1988; Jensen and Høholdt 1989; Jensen et al. 1991; Bömer and Antweiler 1993; Borwein 2002; Borwein and Choi 2000; 2002; Jedwab 2005; Høholdt 2006; Jedwab and Schmidt 2010; Schmidt 2011], but it took the new methods of this paper to discover and verify the properties of the polynomials of our Theorem 1.2.

The interrelation between the additive and multiplicative structures of finite fields endow character polynomials with their remarkable qualities: the coefficients of an *additive character polynomial* are obtained by applying an additive character of a finite field to its nonzero elements arranged multiplicatively (listed as successive powers of a primitive element), while the coefficients of a *multiplicative character polynomial* are obtained by applying a multiplicative character polynomial are obtained by applying a multiplicative character of a finite field to its elements arranged additively (as \mathbb{Z} -linear combinations of the generators, arrayed in a box whose dimensionality equals the number of generators). Thus an additive character polynomial has the form

(1)
$$f(z) = \sum_{j \in S} \psi(\alpha(j+t)) z^j,$$

where $\psi : \mathbb{F}_q \to \mathbb{C}$ is a nontrivial additive character, the *support S* is a set of the form $\{0, 1, \ldots, s - 1\}$, the *translation t* is an element of \mathbb{Z} , and the *arrangement* α is a group epimorphism from \mathbb{Z} to \mathbb{F}_q^* . A multiplicative character polynomial has the form

(2)
$$f(z_1, \ldots, z_e) = \sum_{j=(j_1, \ldots, j_e) \in S} \chi(\alpha(j+t)) \, z_1^{j_1} \cdots z_e^{j_e},$$

where *e* is a positive integer, χ is a nontrivial complex-valued multiplicative character of $\mathbb{F}_q = \mathbb{F}_{p^e}$ with *p* prime, the *support S* is $S_1 \times \cdots \times S_e$ with each S_k a set of

²Gulliver and Parker [2005] have also studied $||f||_4/||f||_2$ for multivariable Littlewood polynomials, but in a very different limit: they let the number of variables tend to infinity while keeping the degree in each variable less than or equal to one.

the form $\{0, 1, ..., s_k - 1\}$, while the *translation* $t = (t_1, ..., t_e)$ is in \mathbb{Z}^e , and the *arrangement* α is a group epimorphism from \mathbb{Z}^e to \mathbb{F}_q . We always extend nontrivial multiplicative characters to take 0 to 0.

We now define the Fekete polynomials and their modifications used in [Høholdt and Jensen 1988] and [Jedwab et al. 2013b]. For an odd prime p, the p-th Fekete polynomial is a multiplicative character polynomial using the quadratic character (Legendre symbol) over the prime field \mathbb{F}_p , support $S = \{0, 1, \dots, p-1\}$, translation t = 0, and arrangement $\alpha : \mathbb{Z} \to \mathbb{F}_p$ given by reduction modulo p. Fekete polynomials are themselves the subject of many fascinating studies linking number theory and analysis [Fekete and Pólya 1912; Pólya 1919; Montgomery 1980; Baker and Montgomery 1990; Conrey et al. 2000; Borwein et al. 2001; Borwein and Choi 2002].

The polynomials used in [Høholdt and Jensen 1988] to obtain asymptotic $||f||_4/||f||_2$ of $\sqrt[4]{7/6}$ have the same character, support, and arrangement, but the translations *t* are chosen such that $t/p \rightarrow 1/4$ as $p \rightarrow \infty$, and any coefficient of 0 (arising from the extended multiplicative character) is replaced with 1 to obtain Littlewood polynomials. To obtain asymptotic $||f||_4/||f||_2$ less than $\sqrt[4]{22/19}$, we used in [Jedwab et al. 2013b] a different limit for t/p, and allowed the support $S = \{0, 1, \ldots, s-1\}$ to be of size other than *p*, and in fact let s/p tend to a number slightly larger than 1 as $p \rightarrow \infty$.

The families of character polynomials used here are based on similar asymptotics: we say that a family $\{f_i\}_{i \in I}$ of additive character polynomials is *size-stable* to mean that if we write \mathbb{F}_{q_i} and S_i for the field and support of f_i , then $\{q_i : i \in I\}$ is infinite and $|S_i|/(q_i-1)$ tends to a positive real number σ (called the *limiting size*) as $q_i \to \infty$. Likewise, we say that a family of *e*-variable multiplicative character polynomials $\{f_i\}_{i \in I}$ is *size-stable* to mean that if we write $\mathbb{F}_{q_i} = \mathbb{F}_{p_i^e}$ and $S_i = S_{i,1} \times \cdots \times S_{i,e}$ for the field and support of f_i , then the set of primes $\{p_i : i \in I\}$ is infinite and for each $k \in \{1, \ldots, e\}$, the ratio $|S_{i,k}|/p_i$ tends to a positive real number σ_k as $q_i \to \infty$. We call $\sigma_1, \ldots, \sigma_e$ the *limiting sizes*. And we say that a family of *e*-variable multiplicative character polynomials $\{f_i\}_{i \in I}$ is *translation-stable* to mean that if we write $\mathbb{F}_{q_i} = \mathbb{F}_{p_i^e}$ and $t_i = (t_{i,1}, \ldots, t_{i,e})$ for the field and translation of f_i , then the set of primes $\{p_i : i \in I\}$ is infinite and for each $k \in \{1, \ldots, e\}$, the ratio $t_{i,k}/p_i$ tends to a real number τ_k as $q_i \to \infty$. We call τ_1, \ldots, τ_e the *limiting translations*.

1C. Subsidiary results. We discovered the polynomials of Theorem 1.2 via a survey, enabled by the methods presented in this paper, of the asymptotic L^4 norms of both additive and multiplicative character polynomials. Quadratic multiplicative characters behave differently than nonquadratic ones, so we treat them separately: we have *quadratic families* in which every character is quadratic, and *nonquadratic* families in which none is. We then have three theorems: one for additive characters

and two for the different types of multiplicative characters, and we express our limiting norms in terms of the function

(3)
$$\Omega(x, y) = \sum_{n \in \mathbb{Z}} \max(0, 1 - |xn - y|)^2,$$

which is defined and continuous on $\{(x, y) \in \mathbb{R}^2 : x \neq 0\}$.

Theorem 1.3. Let $\{f_i\}_{i \in I}$ be a size-stable family, with limiting size σ , of additive character polynomials over fields $\{\mathbb{F}_{q_i}\}_{i \in I}$.

(i) As $q_l \to \infty$, we have

$$\frac{\|f_{\iota}\|_{4}^{4}}{\|f_{\iota}\|_{2}^{4}} \to -\frac{2}{3}\sigma + 2\Omega\left(\frac{1}{\sigma}, 0\right).$$

(ii) This limit is globally minimized if and only if σ is the unique root in $(1, 1 + \frac{9}{64})$ of $x^3 - 12x + 12$.

Theorem 1.4. Let $\{f_i\}_{i \in I}$ be a size-stable family, with limiting sizes $\sigma_1, \ldots, \sigma_e$, of *e*-variable nonquadratic multiplicative character polynomials over fields $\{\mathbb{F}_{q_i}\}_{i \in I}$.

(i) As $q_l \to \infty$, we have

$$\frac{\|f_{\iota}\|_{4}^{4}}{\|f_{\iota}\|_{2}^{4}} \to -\frac{2^{e}}{3^{e}} \prod_{j=1}^{e} \sigma_{i} + 2 \prod_{j=1}^{e} \Omega\left(\frac{1}{\sigma_{j}}, 0\right).$$

(ii) This limit is globally minimized if and only if $\sigma_1, \ldots, \sigma_e$ all equal the unique root in $(1, 1+3^{e+1}/2^{2e+4})$ of

$$x^{3e} - \frac{3^e}{2^{e-3}}(x-1)(3x^2 - 4x + 2)^{e-1}.$$

Theorem 1.5. Let $\{f_i\}_{i \in I}$ be a size- and translation-stable family, with limiting sizes $\sigma_1, \ldots, \sigma_e$ and limiting translations τ_1, \ldots, τ_e , of e-variable quadratic multiplicative character polynomials over fields $\{\mathbb{F}_{q_i}\}_{i \in I}$.

(i) As $q_l \to \infty$, we have

$$\frac{\|f_{\iota}\|_{4}^{4}}{\|f_{\iota}\|_{2}^{4}} \to -\frac{2^{e+1}}{3^{e}} \prod_{j=1}^{e} \sigma_{i} + 2 \prod_{j=1}^{e} \Omega\left(\frac{1}{\sigma_{j}}, 0\right) + \prod_{j=1}^{e} \Omega\left(\frac{1}{\sigma_{j}}, 1 + \frac{2\tau_{j}}{\sigma_{j}}\right).$$

(ii) This limit is globally minimized if and only if $\sigma_1, \ldots, \sigma_e$ all equal the unique root in $(1, 1+3^{e+1}/2^{2e+3})$ of

$$x^{3e} - \frac{3^e}{2^{e-2}}(x-1)(3x^2 - 4x + 2)^{e-1} - \frac{3^e}{2^{2e}}(2x-1)^{2e-1},$$

and $\tau_j \in \{\frac{1}{4}(1-2\sigma_j) + n/2 : n \in \mathbb{Z}\}$ for each $j \in \{1, ..., e\}$.

These new theorems are much more general than the compositum of all previous results on the limiting ratio of L^4 to L^2 norm for character polynomials [Høholdt and Jensen 1988; Jensen and Høholdt 1989; Jensen et al. 1991; Borwein and Choi 2000; Schmidt 2011; Jedwab et al. 2013a; 2013b], and reveal for the first time the full functional form of the asymptotic ratio of norms as it depends on choice of character, limiting size, and limiting translation, thus enabling us to find multivariable Littlewood polynomials with lower asymptotic $||f||_4/||f||_2$ than any known hitherto.

For each $e \ge 1$, let A_e (resp., B_e) be the minimum asymptotic ratio of norms for a family of e-variable nonquadratic (resp., quadratic) character polynomials, as described in Theorem 1.4(ii) (for A_e) and Theorem 1.5(ii) (for B_e). Note that A_1 is also the minimum asymptotic ratio of norms achievable by a family of additive character polynomials as described in Theorem 1.3(ii). Rational approximations of $B_1^4, B_2^4, \ldots, B_5^4$ are obtained later in Lemma 7.1, and if desired, a computer may be used to obtain more accurate approximations of values of various A_e and B_e . For each $e \ge 1$, B_e is to date the lowest known asymptotic $||f||_4/||f||_2$ for a family of *e*-variable Littlewood³ polynomials $f(z_1, \ldots, z_e)$ in the limit as $\deg_{z_1} f, \ldots, \deg_{z_n} f \to \infty$. For e = 1, this recapitulates Corollary 3.2 of [Jedwab et al. 2013b], while for e > 1, the ratio obtained here is strictly lower than any found to date. Until now, the smallest known asymptotic ratio has been whatever can be obtained from the best univariate polynomials and the product construction $||f(z_1)\cdots f(z_e)||_r = ||f(z)||_r^e$, and so we are claiming that $B_e < B_1^e$ for every e > 1. This will give our main result, Theorem 1.2, but in fact we prove something more general: one always obtains a lower asymptotic ratio of norms with a single optimal family of quadratic character polynomials than one does using the product construction with two or more families of character polynomials (which could draw coefficients from $\{-1, 1\}$ or a larger set, depending on the characters involved).

Theorem 1.6. For each $e \ge 1$, let A_e (resp., B_e) be the minimum asymptotic ratio of L^4 to L^2 norm achievable by families of e-variable nonquadratic (resp., quadratic) multiplicative character polynomials, as described in Theorem 1.4(ii) (resp., Theorem 1.5(ii)).

Then $B_e < A_e$ for every $e \ge 1$ and $B_{e_1+e_2} < B_{e_1}B_{e_2}$ for every $e_1, e_2 \ge 1$.

1D. Organization of this paper. To prove Theorems 1.3–1.5, we first establish a general theorem for obtaining the L^4 norm of a polynomial from its Fourier interpolation in Section 3, after setting down notational conventions in Section 2. Our general theorem reduces the problem of computing L^4 norms of character

³Quadratic character polynomials are not always Littlewood because the extended quadratic character χ has $\chi(0) = 0$, so we replace each coefficient of 0 thus produced with a 1, and Corollary A.3 shows that this has no effect on the asymptotic ratio of norms.

polynomials to a pair of calculations (one for additive and one for multiplicative characters) involving Gauss sums, which are presented in Section 4. We use these in Section 5 to prove Theorems 1.3(i), 1.4(i), and 1.5(i). These respectively imply Theorems 1.3(ii), 1.4(ii), and 1.5(ii), but showing this demands delicate arguments which are sketched in Section 6. We prove Theorem 1.6 in Section 7. Some technical results used in Sections 5 and 6 are collected and proved in the Appendix.

2. Notation and conventions

For the rest of this paper p is a prime, and $q = p^e$ with e a positive integer. For any group Γ , we use $\widehat{\Gamma}$ to denote the group of characters from Γ into \mathbb{C} : thus $\widehat{\mathbb{F}}_q$ is the group of additive characters from \mathbb{F}_q to \mathbb{C} and $\widehat{\mathbb{F}}_q^*$ the group of multiplicative characters from \mathbb{F}_q^* to \mathbb{C} . We extend any nontrivial $\chi \in \widehat{\mathbb{F}}_q^*$ so that $\chi(0) = 0$.

To write the multiplicative character polynomial (2) compactly, we use the convention that if $j = (j_1, ..., j_e) \in \mathbb{Z}^e$, the notation z^j is a shorthand for $z_1^{j_1} \cdots z_e^{j_e}$. To make it easier to speak about supports of character polynomials (1) and (2), we call a finite set of consecutive integers a *segment*, and a finite Cartesian product of segments a *box*. If *S* is a subset of \mathbb{Z}^n and $t \in \mathbb{Z}^n$, then S + t is the translated subset $\{s + t : s \in S\}$.

3. L^4 norms via the Fourier transform

If Γ is a finite abelian group and $\{F_g\}_{g\in\Gamma}$ is a family of complex numbers, then for any $\eta \in \widehat{\Gamma}$, we have the Fourier transform

$$\hat{F}_{\eta} = \sum_{g \in \Gamma} F_g \eta(g),$$

with inverse

$$F_g = \frac{1}{|\Gamma|} \sum_{\eta \in \widehat{\Gamma}} \widehat{F}_\eta \overline{\eta(g)}.$$

We express the L^4 norm in terms of the Fourier interpolation.

Theorem 3.1. Let Γ be a finite abelian group, $\{F_g\}_{g\in\Gamma}$ a family of complex numbers, n a positive integer, and $\pi \in \text{Hom}(\mathbb{Z}^n, \Gamma)$. For any $\eta \in \widehat{\Gamma}$, let $\eta' \in \widehat{\mathbb{Z}^n}$ be $\eta \circ \pi$. If U is a finite subset of \mathbb{Z}^n and $F(z) = \sum_{u \in U} F_{\pi(u)} z^u \in \mathbb{C}[z_1, \ldots, z_n]$, then

$$\|F\|_{4}^{4} = \frac{1}{|\Gamma|^{5}} \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \sum_{\substack{\kappa,\lambda,\mu,\nu \in \widehat{\Gamma}}} \overline{\kappa'(a)\lambda'(b)}\mu'(c)\nu'(d)H(\kappa,\lambda,\mu,\nu),$$

where

$$H(\kappa,\lambda,\mu,\nu) = \sum_{\xi\in\widehat{\Gamma}} \hat{F}_{\xi\kappa} \hat{F}_{\xi\lambda} \hat{F}_{\xi\mu} \hat{F}_{\xi\nu}.$$

Proof. By the definition of the L^4 norm, we have

$$\|F\|_{4}^{4} = \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} F_{\pi(a)} F_{\pi(b)} \overline{F_{\pi(c)} F_{\pi(d)}},$$

and thus, using the inverse Fourier transform,

$$\|F\|_{4}^{4} = \frac{1}{|\Gamma|^{4}} \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \sum_{\substack{\kappa,\lambda,\mu,\nu \in \widehat{\Gamma}}} \hat{F}_{\kappa} \hat{F}_{\lambda} \overline{\hat{F}_{\mu} \hat{F}_{\nu} \kappa'(a) \lambda'(b)} \mu'(c) \nu'(d).$$

Since we are summing κ over all $\widehat{\Gamma}$, we can replace κ by $\xi \kappa$ for any given $\xi \in \widehat{\Gamma}$, and also do likewise with λ , μ , ν to obtain

$$\|F\|_{4}^{4} = \frac{1}{|\Gamma|^{4}} \sum_{\substack{a,b,c,d \in U\\a+b=c+d}} \sum_{\kappa,\lambda,\mu,\nu \in \widehat{\Gamma}} \hat{F}_{\xi\kappa} \hat{F}_{\xi\lambda} \overline{\hat{F}_{\xi\mu}} \hat{F}_{\xi\nu\kappa'(a)\lambda'(b)} \mu'(c)\nu'(d)$$

where we have omitted mention of the resulting factor of $\overline{\xi'(a)\xi'(b)\xi'(c)\xi'(d)}$, which equals 1 in view of the constraint in the first summation. Now sum ξ over $\widehat{\Gamma}$ and divide by $|\Gamma| = |\widehat{\Gamma}|$ to finish.

We apply this general theorem to additive and multiplicative character polynomials in two corollaries below. Such polynomials have Gauss sums as their Fourier coefficients, so for any $\psi \in \widehat{\mathbb{F}}_q$ and $\chi \in \widehat{\mathbb{F}}_q^*$, we define the *Gauss sum* associated with ψ and χ to be

$$G(\psi, \chi) = \sum_{a \in \mathbb{F}_q^*} \psi(a) \chi(a).$$

Corollary 3.2. If f(z) is an additive character polynomial with character $\psi \in \widehat{\mathbb{F}}_q$, support *S*, translation *t*, and arrangement α , then

$$\|f\|_4^4 = \frac{1}{(q-1)^5} \sum_{\substack{a,b,c,d \in S+t \\ a+b=c+d}} \sum_{\substack{\kappa,\lambda,\mu,\nu \in \widehat{\mathbb{F}_q^*}}} \overline{\kappa'(a)\lambda'(b)}\mu'(c)\nu'(d)H(\kappa,\lambda,\mu,\nu)$$

where for any $\eta \in \widehat{\mathbb{F}_q^*}$, we let $\eta' \in \widehat{\mathbb{Z}}$ be $\eta \circ \alpha$, and

$$H(\kappa,\lambda,\mu,\nu) = \sum_{\xi \in \widehat{\mathbb{F}_q^*}} G(\psi,\xi\kappa) G(\psi,\xi\lambda) \overline{G(\psi,\xi\mu)} G(\psi,\xi\nu).$$

Proof. Our additive character polynomial $f(z) = \sum_{s \in S} \psi(\alpha(s+t))z^s$ has the same L^r norms as $F(z) = z^t f(z) = \sum_{u \in S+t} \psi(\alpha(u))z^u$, so take $\Gamma = \mathbb{F}_q^*$, $F_g = \psi(g)$, $n = 1, \pi = \alpha$, and U = S + t in Theorem 3.1, and note that for $\eta \in \widehat{\mathbb{F}}_q^*$ we have $\widehat{F}_{\eta} = G(\psi, \eta)$.

Corollary 3.3. If f(z) is a multiplicative character polynomial with character $\chi \in \widehat{\mathbb{F}}_q^*$, support S, translation t, and arrangement α , then

$$\|f\|_4^4 = \frac{1}{q^5} \sum_{\substack{a,b,c,d \in S+t \\ a+b=c+d}} \sum_{\substack{\kappa,\lambda,\mu,\nu \in \widehat{\mathbb{F}}_q}} \overline{\kappa'(a)\lambda'(b)}\mu'(c)\nu'(d)H(\kappa,\lambda,\mu,\nu)$$

where for any $\eta \in \widehat{\mathbb{F}}_q$, we let $\eta' \in \widehat{\mathbb{Z}}^e$ be $\eta \circ \alpha$, and

$$H(\kappa,\lambda,\mu,\nu) = \sum_{\xi \in \widehat{\mathbb{F}}_q} G(\xi\kappa,\chi) G(\xi\lambda,\chi) \overline{G(\xi\mu,\chi)G(\xi\nu,\chi)}.$$

Proof. Our multiplicative character polynomial $f(z) = \sum_{s \in S} \chi(\alpha(s+t))z^s$ has the same L^r norms as $F(z) = z^t f(z) = \sum_{u \in S+t} \chi(\alpha(u))z^u$, so take $\Gamma = \mathbb{F}_q$, $F_g = \chi(g)$, $n = e, \pi = \alpha$, and U = S + t in Theorem 3.1, and note that for $\eta \in \widehat{\mathbb{F}}_q$ we have $\widehat{F}_{\eta} = G(\eta, \chi)$.

The key to L^4 norms is then the evaluation of the sums $H(\kappa, \lambda, \mu, \nu)$ in the above two corollaries, which we take up in the next section.

4. Two propositions on summations of Gauss sums

Here we estimate the values of the summations H that appear in Corollaries 3.2 and 3.3. We begin with some basic facts about Gauss sums, which are proved in Theorems 5.11 and 5.12 of [Lidl and Niederreiter 1997].

Lemma 4.1. If $\psi \in \widehat{\mathbb{F}}_q$ and $\chi \in \widehat{\mathbb{F}}_q^*$, then

- (i) $G(\psi, \chi) = q 1$ if both characters are trivial,
- (ii) $G(\psi, \chi) = 0$ if ψ is trivial and χ is not,
- (iii) $G(\psi, \chi) = -1$ if χ is trivial and ψ is not,
- (iv) $|G(\psi, \chi)| = \sqrt{q}$ if both characters are nontrivial, and
- (v) $\sum_{a \in \mathbb{F}_q^*} \psi(ba) \chi(a) = \overline{\chi}(b) G(\psi, \chi)$ for any $b \in \mathbb{F}_q^*$.

We first estimate the summation H appearing in Corollary 3.2.

Proposition 4.2. Let ψ be a nontrivial character in $\widehat{\mathbb{F}}_q$, and $\kappa, \lambda, \mu, \nu \in \widehat{\mathbb{F}}_a^*$. If

$$H = \sum_{\xi \in \widehat{\mathbb{F}_q^*}} G(\psi, \xi\kappa) G(\psi, \xi\lambda) \overline{G(\psi, \xi\mu)} G(\psi, \xi\nu),$$
$$M = \begin{cases} (q-1)^3 & \text{if } \{\kappa, \lambda\} = \{\mu, \nu\}, \\ 0 & \text{otherwise,} \end{cases}$$

then $|H - M| \leq (q - 1)q\sqrt{q}$.

Proof. First we consider the case where $\{\kappa, \lambda\} = \{\mu, \nu\}$, wherein

$$H = \sum_{\xi \in \widehat{\mathbb{F}_q^*}} |G(\psi, \xi \kappa)|^2 |G(\psi, \xi \lambda)|^2.$$

One can work out from parts (iii) and (iv) of Lemma 4.1 that $H = (q-2)q^2 + 1$ if $\kappa = \lambda = \mu = \nu$ and $H = (q-3)q^2 + 2q$ otherwise. Thus H - M = (q-1)(q-2) or 1 - q.

Now we consider the case where $\{\kappa, \lambda\} \neq \{\mu, \nu\}$, wherein

$$H = \sum_{\xi \in \widehat{\mathbb{F}_q^*}} \sum_{\substack{w, x, y, z \in \mathbb{F}_q^* \\ w, x, y, z \in \mathbb{F}_q^*}} \psi(w + x - y - z)\xi(wxy^{-1}z^{-1})\kappa(w)\lambda(x)\overline{\mu(y)\nu(z)}$$
$$= (q-1)\sum_{\substack{w, x, y, z \in \mathbb{F}_q^* \\ wx = yz}} \psi(w + x - y - z)\kappa(w)\lambda(x)\overline{\mu(y)\nu(z)}.$$

Now reparametrize the sum with w = uy and z = ux to obtain

$$H = (q-1) \sum_{u,x,y \in \mathbb{F}_q^*} \overline{\psi}((u-1)x) \psi((u-1)y) \kappa \overline{\nu}(u) \lambda \overline{\nu}(x) \kappa \overline{\mu}(y),$$

and since $\{\kappa, \lambda\} \neq \{\mu, \nu\}$, we can restrict to $u \neq 1$ without changing the value of the summation. Then Lemma 4.1(v) tells us that when we sum over x and y, we obtain

$$H = (q-1)G(\psi, \kappa \overline{\mu})G(\overline{\psi}, \lambda \overline{\nu}) \sum_{u \neq 0, 1} \overline{\kappa \lambda} \mu \nu(u-1)\kappa \overline{\nu}(u).$$

Now $\overline{\kappa\lambda\mu\nu}$ and $\overline{\kappa\nu}$ cannot both be the trivial character since $\{\kappa, \lambda\} \neq \{\mu, \nu\}$. If $\overline{\kappa\lambda\mu\nu}$ is trivial, then the sum over u is -1; if $\overline{\kappa\nu}$ is trivial, the sum is $-\overline{\kappa\lambda\mu\nu}(-1)$; otherwise, let ω be a generator of $\widehat{\mathbb{F}_q^*}$ and we can write the sum over u as

$$\sum_{u\neq 0,1}\omega((u-1)^a u^b)$$

for some nonzero $a, b \in \mathbb{Z}/(q-1)\mathbb{Z}$, and use the Weil bound [Weil 1948; Lidl and Niederreiter 1997, Theorem 5.41] to see that this sum is bounded in magnitude by \sqrt{q} . We can use this fact, along with Lemma 4.1(iii), (iv), to see that $|H| \le (q-1)q\sqrt{q}$.

Similarly, we estimate the summation *H* appearing in Corollary 3.3.

Proposition 4.3. Let χ be a nontrivial character in $\widehat{\mathbb{F}}_q^*$, and $\kappa, \lambda, \mu, \nu \in \widehat{\mathbb{F}}_q$. If

$$H = \sum_{\xi \in \widehat{\mathbb{F}}_q} G(\xi \kappa, \chi) G(\xi \lambda, \chi) \overline{G(\xi \mu, \chi) G(\xi \nu, \chi)}$$

and

$$M = \begin{cases} q^3 & \text{if } \{\kappa, \lambda\} = \{\mu, \nu\}, \\ q^3 & \text{if } \kappa = \lambda, \, \mu = \nu, \, \text{and } \chi \text{ is the quadratic character}, \\ 0 & \text{otherwise}, \end{cases}$$

then $|H - M| \leq 3q^2 \sqrt{q}$.

Proof. Let ϵ be the canonical additive character over \mathbb{F}_q . Then for any $\eta \in \widehat{\mathbb{F}}_q$, there is a unique $y \in \mathbb{F}_q$ such that $\eta(z) = \epsilon(yz)$ for all $z \in \mathbb{F}_q$. Let a, b, c, d be chosen so that $\kappa(z) = \epsilon(az), \lambda(z) = \epsilon(bz), \mu(z) = \epsilon(cz)$, and $\nu(z) = \epsilon(dz)$ for all $z \in \mathbb{F}_q$. Furthermore, we shall parametrize the sum of ξ over $\widehat{\mathbb{F}}_q$ in the definition of H by a sum over $x \in \mathbb{F}_q$, and replace $\xi(z)$ with $\epsilon(xz)$ wherever it occurs. Thus, in view Lemma 4.1(v) and (ii), we have

$$H = |G(\epsilon, \chi)|^4 \sum_{x \in \mathbb{F}_q} \overline{\chi}((x+a)(x+b))\chi((x+c)(x+d)),$$

and $|G(\epsilon, \chi)| = \sqrt{q}$ by Lemma 4.1(iv). Let *m* be the order of χ . Then

$$H = q^2 \sum_{x \in \mathbb{F}_q} \chi((x+a)^{m-1}(x+b)^{m-1}(x+c)(x+d)).$$

The magnitude of the Weil sum over x is bounded by $3\sqrt{q}$ unless the polynomial $(x+a)^{m-1}(x+b)^{m-1}(x+c)(x+d)$ is an *m*-th power in $\mathbb{F}_q[x]$. (See [Weil 1948; Lidl and Niederreiter 1997, Theorem 5.41].) It is an *m*-th power only if $\{a, b\} = \{c, d\}$ or if m = 2, a = b, and c = d, in which cases the Weil sum is either q - 1 (if a = b = c = d) or q - 2 (if there are two distinct roots).

5. Asymptotic L^4 norm

We prove Theorems 1.3(i), 1.4(i), and 1.5(i) in this section, by using the propositions from the previous section with Corollaries 3.2 and 3.3.

Proof of Theorems 1.4(i) and 1.5(i). Let χ be a nontrivial character in $\widehat{\mathbb{F}}_q^*$, let α be an epimorphism from \mathbb{Z}^e to \mathbb{F}_q , let $t \in \mathbb{Z}^e$, and let $S = S_1 \times \cdots \times S_e$ be a box where each S_j is a nonempty segment of the form $\{0, 1, \ldots, s_j - 1\}$. Recall from Section 2 our notational convention that z^s is shorthand for $z_1^{s_1} \cdots z_e^{s_e}$ when $s = (s_1, \ldots, s_e) \in \mathbb{Z}^e$. Let f(z) be the multiplicative character polynomial $\sum_{s \in S} \chi(\alpha(s + t))z^s$. We shall calculate $||f||_4^4$ first, and then investigate what happens asymptotically to this quantity in the limits considered in Theorems 1.4(i) and 1.5(i). By Corollary 3.3, we have

(4)
$$\|f\|_4^4 = \frac{1}{q^5} \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \sum_{\kappa,\lambda,\mu,\nu \in \widehat{\mathbb{F}}_q} \overline{\kappa'(a)\lambda'(b)}\mu'(c)\nu'(d)H(\kappa,\lambda,\mu,\nu),$$

where we let U = S + t, and for any $\eta \in \widehat{\mathbb{F}}_q$, we let $\eta' = \eta \circ \alpha$, and

$$H(\kappa,\lambda,\mu,\nu) = \sum_{\xi \in \widehat{\mathbb{F}}_q} G(\xi\kappa,\chi) G(\xi\lambda,\chi) \overline{G(\xi\mu,\chi)G(\xi\nu,\chi)}$$

By Proposition 4.3, we can write $H(\kappa, \lambda, \mu, \nu) = M(\kappa, \lambda, \mu, \nu) + N(\kappa, \lambda, \mu, \nu)$ with

$$M(\kappa, \lambda, \mu, \nu) = \begin{cases} q^3 & \text{if } \{\kappa, \lambda\} = \{\mu, \nu\}, \\ q^3 & \text{if } \kappa = \lambda, \mu = \nu, \text{ and } \chi \text{ is the quadratic character,} \\ 0 & \text{otherwise,} \end{cases}$$

and

(5)
$$|N(\kappa,\lambda,\mu,\nu)| \le 3q^2\sqrt{q},$$

for all κ , λ , μ , $\nu \in \widehat{\mathbb{F}_q^*}$.

If χ is nonquadratic, when we write out separately the contributions from M and N to (4), we get $||f||_4^4 = A + B - D + E$, where

$$\begin{split} A &= \frac{1}{q^2} \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \sum_{\substack{\kappa,\lambda \in \widehat{\mathbb{F}}_q}} \kappa'(c-a)\lambda'(d-b), \\ B &= \frac{1}{q^2} \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \sum_{\substack{\kappa,\lambda \in \widehat{\mathbb{F}}_q}} \kappa'(d-a)\lambda'(c-b), \\ D &= \frac{1}{q^2} \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \sum_{\substack{\kappa \in \widehat{\mathbb{F}}_q}} 1, \\ E &= \frac{1}{q^5} \sum_{\substack{\kappa,\lambda,\mu,\nu \in \widehat{\mathbb{F}}_q}} N(\kappa,\lambda,\mu,\nu) \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \overline{\kappa'(a)\lambda'(b)}\mu'(c)\nu'(d). \end{split}$$

Here *A* accounts for the value of *M* when $(\kappa, \lambda) = (\mu, \nu)$, and *B* accounts for the value of *M* when $(\kappa, \lambda) = (\nu, \mu)$, while *D* corrects for the double counting by *A* and *B* of the case $\kappa = \lambda = \mu = \nu$.

Note that A = B, and that A counts the number of $(a, b, c, d) \in U^4$ with $c-a=b-d \in \ker \alpha$. If we write $a = (a_1, \ldots, a_e), b = (b_1, \ldots, b_e), c = (c_1, \ldots, c_e)$, and $d = (d_1, \ldots, d_e)$, then $c-a \in \ker \alpha$ is equivalent to $c_1 - a_1 \equiv \cdots \equiv c_e - a_e \equiv 0$ (mod p), because α is an epimorphism from \mathbb{Z}^e to $\mathbb{F}_q = \mathbb{F}_{p^e}$ and so factors as $\alpha = \gamma \circ \beta$, with $\beta \colon \mathbb{Z}^e \to (\mathbb{Z}/p\mathbb{Z})^e$ coordinate-wise reduction modulo p and $\gamma \colon (\mathbb{Z}/p\mathbb{Z})^e \to \mathbb{F}_q$ a group isomorphism. Now $U = U_1 \times \cdots \times U_e$ with each $U_j = \{t_j, t_j + 1, \ldots, t_j + |S_j| - 1\}$, so for each $n \in \mathbb{Z}$, there are max $(0, |S_j| - p|n|)$ ways for $c_j - a_j$ to equal pn and the same number of ways for $b_j - d_j$ to equal pn.

So

$$A = B = \prod_{j=1}^{e} \sum_{n_j \in \mathbb{Z}} \max(0, |S_j| - p|n_j|)^2.$$

On the other hand, qD counts the number of $(a, b, c, d) \in U^4$ with c-a = b-d, so by the same argument we just used (with modulus 1 instead of p),

$$qD = \prod_{j=1}^{e} \sum_{n_j \in \mathbb{Z}} \max(0, |S_j| - |n_j|)^2,$$

from which we can compute

$$D = \prod_{j=1}^{e} \left(\frac{2|S_j|^3 + |S_j|}{3p} \right).$$

Now we bound *E* via two bounds: our bound (5) on *N*, and a technical result, Lemma A.1 in the Appendix, bounding the inner sum of *E*. We satisfy the condition on α demanded by this lemma, since $\alpha = \gamma \circ \beta$ with $\beta \colon \mathbb{Z}^e \to (\mathbb{Z}/p\mathbb{Z})^e$ coordinatewise reduction modulo *p* and $\gamma \colon (\mathbb{Z}/p\mathbb{Z})^e \to \mathbb{F}_q$ a group isomorphism. With these two bounds, we obtain

(6)
$$|E| \le 3 \cdot 64^e q \sqrt{q} \prod_{j=1}^e \max\left(1, \frac{|S_j|}{p}\right)^3 \prod_{j=1}^e (1 + \log p)^3.$$

Now we divide $||f||_4^4 = A + B - D + E$ by $||f||_2^4$ and consider the limit where each $|S_j|/p \to \sigma_j$ as $q \to \infty$, that is, consider what happens in a size-stable family of polynomials. Another technical result, Lemma A.2, shows that we can replace the denominator $||f||_2^4$ with $|S|^2$ without changing the limit. Then recall the definition (3) of Ω , and note that $A/|S|^2$ and $B/|S|^2$ tend to $\prod_{j=1}^e \Omega(1/\sigma_j, 0)$, that $D/|S|^2$ tends to $(\frac{2}{3})^e \prod_{j=1}^e \sigma_j$, and that $|E|/|S|^2$ tends to 0 in this limit.

If χ is quadratic, the proof is done in the same manner, except that there is now a contribution from *M* in the case where $\kappa = \lambda$ and $\mu = \nu$, and so we get $||f||_4^4 = A + B + C - 2D + E$, where *A*, *B*, *D*, and *E* are as defined above, and

$$C = \frac{1}{q^2} \sum_{\substack{a,b,c,d \in U \\ a+b=c+d}} \sum_{\substack{\kappa,\mu \in \widehat{\mathbb{F}}_q}} \kappa'(-a-b)\mu'(c+d).$$

Note that we subtract *D* twice now because *A*, *B*, and *C* count the case where $\kappa = \lambda = \mu = \nu$ three times. *C* counts the number of $(a, b, c, d) \in U^4$ with $a + b = c + d \in \ker \alpha$. Following the method we used to determine *A*, write $a = (a_1, \ldots, a_e), b = (b_1, \ldots, b_e), c = (c_1, \ldots, c_e), \text{ and } d = (d_1, \ldots, d_e), \text{ and note that } a + b \in \ker \alpha \text{ is equivalent to } a_1 + b_1 \equiv \cdots \equiv a_e + b_e \equiv 0 \pmod{p}.$

Since $U = U_1 \times \cdots \times U_e$ with each $U_j = \{t_j, t_j + 1, \dots, t_j + |S_j| - 1\}$, there are $\max(0, |S_j| - |np - (2t_j + |S_j| - 1)|)$ ways to obtain $a_j + b_j = np$ with $(a_j, b_j) \in U_j^2$, and the same number of ways to obtain $c_j + d_j = np$ with $(c_j, d_j) \in U_j^2$, so

$$C = \prod_{j=1}^{e} \sum_{n_j \in \mathbb{Z}} \max(0, |S_j| - |pn_j - |S_j| - 2t_j + 1|)^2,$$

and if we have both size- and translation-stability, then $|S_j|/p \to \sigma_j$ and $t_j/p \to \tau_j$ as $q \to \infty$, so that $C/|S|^2 \to \prod_{j=1}^e \Omega\left(1/\sigma_j, 1+2\tau_j/\sigma_j\right)$.

Proof of Theorem 1.3(i). The proof is the same, mutatis mutandis, as for the e = 1 case of Theorem 1.4(i), with the roles of \mathbb{F}_q and \mathbb{F}_q^* exchanged. Corollary 3.2 and Proposition 4.2 replace Corollary 3.3 and Proposition 4.3, and Lemma A.2 becomes unnecessary as $||f||_2^2$ for an additive character polynomial f is always precisely equal to the cardinality of the support of f. These, and other attendant minor changes resulting from the exchange of \mathbb{F}_q and \mathbb{F}_q^* , cause (5) to become $|N(\kappa, \lambda, \mu, \nu)| \leq (q-1)q\sqrt{q}$, and (6) to become

$$|E| \le 64q\sqrt{q} \max(1, |S|/(q-1))^3 (1 + \log(q-1))^3,$$

and any other printed instance of p or q should be replaced with q - 1.

6. Minimizing the asymptotic ratio of L^4 to L^2 norm

Here we prove Theorems 1.3(ii), 1.4(ii), and 1.5(ii) by finding the limiting sizes and (for quadratic multiplicative character polynomials) the limiting translations that globally minimize the ratio of the L^4 to L^2 norm.

Proof of Theorem 1.4(ii). In view of Theorem 1.4(i), we are trying to minimize the limiting ratio of norms, given by the function

$$K(x_1, \dots, x_e) = -\frac{2^e}{3^e} \prod_{j=1}^e x_i + 2 \prod_{j=1}^e \Phi(x_j),$$

for x_1, \ldots, x_e positive real numbers (the limiting sizes), where for positive x, we define

(7)
$$\Phi(x) = \Omega\left(\frac{1}{x}, 0\right) = \sum_{n \in \mathbb{Z}} \max\left(0, 1 - \frac{|n|}{x}\right)^2,$$

which is differentiable for $x \neq 0$ and is C^{∞} for $x \notin \mathbb{Z}$.

Step 1. We can assume that each $x_j > 1$ because otherwise the partial derivative of *K* with respect to x_j would be negative.

Step 2. We can assume that $(x_1, \ldots, x_e) \in (1, 3)^e$: indeed, Lemma A.6 shows that $K(x_1, \ldots, x_e) \ge \Phi(x_1) \cdots \Phi(x_e)$, and note that $\Phi(x)$ is increasing for x > 1, that $\Phi(1) = 1$, $\Phi(3) = \frac{19}{9} > 2$, and $K(1, \ldots, 1) = 2 - \left(\frac{2}{3}\right)^e < 2$. This proves that a global minimum exists and lies in $(1, 3)^e$: the closure of $(1, 3)^e$ is compact and *K* is continuous thereupon.

Step 3. Suppose $(\sigma_1, \ldots, \sigma_e)$ to be global minimizer of K. Then the partial derivatives of K must vanish there, whence for each $k \in \{1, \ldots, e\}$, we have $2u(\sigma_k) \prod_{j=1}^e U(\sigma_j) = 1$, where $u(x) = x \Phi'(x) / \Phi(x)$ and $U(x) = 3\Phi(x) / (2x)$ for x > 0.

Step 4. Then one can show that U(x) is strictly decreasing on [1, 3], with $U(3) = \frac{19}{18}$. Thus we must have $u(\sigma_k) \le \frac{1}{2} \cdot \left(\frac{18}{19}\right)^e < \frac{1}{2}$ for all $k \in \{1, \dots, e\}$. Then examination of u(x) shows that u(x) strictly increases from 0 to $\frac{1}{2}$ for $x \in [1, 2 - \sqrt{2/3}]$, and then $u(x) > \frac{1}{2}$ for $x \in (2 - \sqrt{2/3}, 3)$. This then forces $\sigma_1 = \cdots = \sigma_e < 2 - \sqrt{2/3} < \frac{6}{5}$. Step 5. Now $U(\sigma_1) > U\left(\frac{6}{5}\right) > \frac{9}{7}$, so this forces $u(\sigma_1) < \frac{1}{2} \cdot \left(\frac{7}{9}\right)^e \le \frac{7}{18}$, which in turn forces $\sigma_1 < \frac{8}{7}$. Then $U(\sigma_1) > U\left(\frac{8}{7}\right) > \frac{4}{3}$, so this forces $u(\sigma_1) < \frac{1}{2} \cdot \left(\frac{3}{4}\right)^e$. Since $u(x) \ge 8(x-1)/3$ for $x \in [1, \frac{8}{7}]$, this forces $\sigma_1 < 1 + 3^{e+1}/2^{2e+4}$.

Step 6. Now our problem is reduced to the single-variable minimization of

$$\Theta(x) = K(x, \dots, x) = -\left(\frac{2x}{3}\right)^e + 2\Phi(x)^e$$

on the interval $(1, 1+3^{e+1}/2^{2e+4})$. It is not hard to see that $d\Theta/dx$ vanishes if and only if $x^{3e} - (3^e/2^{e-3})(x-1)(3x^2-4x+2)^{e-1}$ vanishes. Meanwhile $d^2\Theta/dx^2 > 0$ on our interval: by computing its value and then dropping a nonnegative term, we can see that $d^2\Theta/dx^2$ is at least

$$-e(e-1)\frac{2^{e}x^{e-2}}{3^{e}} + 8e\frac{3-2x}{x^{4}}\Phi(x)^{e-1} \ge -e(e-1)\frac{19^{e}}{24^{e}} + 2e > 0.$$

This proves that there is a unique minimum: the unique root a_e of

$$x^{3e} - \frac{3^e}{2^{e-3}}(x-1)(3x^2 - 4x + 2)^{e-1}$$

lying in $(1, 1+3^{e+1}/2^{2e+4})$.

Proof of Theorem 1.3(ii). This is accomplished exactly as the e = 1 case of the proof of Theorem 1.4(ii) above, save that Lemma A.5 replaces Lemma A.6.

Proof of Theorem 1.5(ii). In view of Theorem 1.5 (i), we are trying to minimize the limiting ratio of norms, given by the function

(8)
$$-\frac{2^{e+1}}{3^{e}}\prod_{j=1}^{e}x_{j}+2\prod_{j=1}^{e}\Omega\left(\frac{1}{x_{j}},0\right)+\prod_{j=1}^{e}\Omega\left(\frac{1}{x_{j}},1+\frac{2y_{j}}{x_{j}}\right)$$

for x_1, \ldots, x_e positive real numbers (the limiting sizes) and y_1, \ldots, y_e arbitrary real numbers (the limiting translations).

Step 1. We invoke Lemma A.4(i) to see that we can confine our search to x_1, \ldots, x_e greater than or equal to $\frac{1}{2}$. For as long as $x_j \leq \frac{1}{2}$, the lemma shows that we can always arrange for y_j to be such that $\Omega(x_j^{-1}, 1 + 2x_j^{-1}y_j) = 0$, and we note that $\Omega(x_j^{-1}, 0) = 1$ for all $x_j \in (0, \frac{1}{2}]$. Thus we can increase x_j to $\frac{1}{2}$ to lower (8) through the term $-(2^{e+1}/3^e) \prod_{j=1}^e x_j$ while keeping the other terms constant.

Step 2. Now we invoke Lemma A.4(ii) to see that for fixed x_1, \ldots, x_e , we minimize the last term of (8) if and only if we arrange that

$$y_j \in \left\{\frac{1-2x_j}{4} + \frac{m}{2} : m \in \mathbb{Z}\right\}$$

for each $j \in \{1, ..., e\}$. The problem is thus reduced to the minimization of

$$\Lambda(x_1, \dots, x_e) = -\frac{2^{e+1}}{3^e} \prod_{j=1}^e x_j + 2 \prod_{j=1}^e \Phi(x_j) + \prod_{j=1}^e \Psi(x_j)$$

for x_1, \ldots, x_e positive real numbers, where $\Phi(x)$ is as defined in (7), and

$$\Psi(x) = \sum_{n \in \mathbb{Z}} \max\left(0, 1 - \frac{|2n+1|}{2x}\right)^2$$

for x > 0. Note that Ψ is differentiable for $x \neq 0$ and is C^{∞} for $x \notin \mathbb{Z} + \frac{1}{2}$.

Step 3. We can assume that each $x_j > \frac{1}{2}$ because otherwise the partial derivative of Λ with respect to x_j would be negative.

Step 4. We can assume that $x_1, \ldots, x_e \in (\frac{1}{2}, 3)$: indeed, Lemma A.6 shows that $\Lambda(x_1, \ldots, x_e) \ge \Phi(x_1) \cdots \Phi(x_e)$, and note that $\Phi(x)$ is nondecreasing for $x > \frac{1}{2}$, that $\Phi(\frac{1}{2}) = 1$, $\Phi(3) = \frac{19}{9} > 2$, and $\Lambda(1, \ldots, 1) = 2 - 2(\frac{2}{3})^e + (\frac{1}{2})^e < 2$. This proves that a global minimum exists and lies in $(\frac{1}{2}, 3)^e$: the closure of $(\frac{1}{2}, 3)^e$ is compact and Λ is continuous thereupon.

Step 5. Suppose that $(\sigma_1, \ldots, \sigma_e)$ is a global minimizer of Λ . Then the partial derivatives of Λ must vanish there, whence

(9)
$$u(\sigma_k) \prod_{j=1}^e U(\sigma_j) + \frac{1}{2}v(\sigma_k) \prod_{j=1}^e V(\sigma_j) = 1,$$
$$x \Phi'(x) = 2\Phi(x) = x \Psi'(x) = 2\Psi(x)$$

where $u(x) = \frac{x\Phi'(x)}{\Phi(x)}$, $U(x) = \frac{3\Phi(x)}{2x}$, $v(x) = \frac{x\Psi'(x)}{\Psi(x)}$, and $V(x) = \frac{3\Psi(x)}{2x}$ for $x > \frac{1}{2}$.

Step 6. We can assume $\sigma_1, ..., \sigma_e \in (1, 3)$: see (9) and note that u(x) = 0 for $x \in (\frac{1}{2}, 1], \frac{1}{2}v(x)V(x) < 1$ for $x \in (\frac{1}{2}, 1]$, and V(x) < 1 for $x \in (\frac{1}{2}, 3)$.

Step 7. It is not difficult to show that U(x) strictly decreases and V(x) strictly increases on [1, 3] with $U(3) = \frac{19}{18}$ and $V(1) = \frac{3}{4}$, and that $0 \le u(x) < 1 \le v(x)$ for $x \in [1, 3]$. Thus (9) shows that we must have $u(\sigma_k) < \left(\frac{18}{19}\right)^e (1 - 3^e/2^{2e+1})$ for all *k*. This forces $u(\sigma_k) < \frac{7}{10}$ for all *k*, and examination of the function *u* shows that $u(x) \ge \frac{7}{10}$ for $x \in [\frac{5}{2}, 3]$, and so we must have $\sigma_k < \frac{5}{2}$ for all *k*.

Now one can repeat the argument on the interval $\begin{bmatrix} 1, \frac{5}{2} \end{bmatrix}$ to show that every $\sigma_k < 2$, then repeat it again on $\begin{bmatrix} 1, 2 \end{bmatrix}$ to show $\sigma_k < \frac{5}{4}$. Further repetitions give $\sigma_k < \frac{6}{5}$, $\sigma_k < \frac{13}{11}$, and $\sigma_k < \frac{7}{6}$. Since $U(x) > \frac{4}{3}$ while v(x), $V(x) \ge 0$ for $x \in (1, \frac{7}{6})$, we have $u(\sigma_k) < (\frac{3}{4})^e$ for all k, and since $u(x) \ge 8(x-1)/3$ for $x \in [1, \frac{7}{6}]$, this means that $\sigma_k < 1 + 3^{e+1}/2^{2e+3}$ for all k.

Step 8. So we have $1 < \sigma_k < \min(\frac{7}{6}, 1+3^{e+1}/2^{2e+3})$ for all k. Consider the products in (9): since each $\sigma_j \in (1, \frac{7}{6})$, we have

$$\left(\frac{4}{3}\right)^e < \prod_{j=1}^e U(\sigma_j) < \left(\frac{3}{2}\right)^e$$
, while $\left(\frac{3}{4}\right)^e < \prod_{j=1}^e V(\sigma_j) < \left(\frac{7}{8}\right)^e$.

We now claim that for a given $A \in \left[\left(\frac{4}{3}\right)^e, \left(\frac{3}{2}\right)^e\right]$ and $B \in \left[\left(\frac{3}{4}\right)^e, \left(\frac{7}{8}\right)^e\right]$, there is at most one solution $x \in \left(1, \frac{7}{6}\right)$ to $Au(x) + \frac{1}{2}Bv(x) = 1$, which will force $\sigma_1 = \cdots = \sigma_e$. For if we set $w(x) = Au(x) + \frac{1}{2}Bv(x)$, we can show that w'(x) > 0 for $x \in \left(1, \frac{7}{6}\right)$: on this interval, we have $u(x) = 4(x-1)/(3x^2-4x+2)$ v(x) = 2/(2x-1), and it is not difficult to show that $u'(x) > \frac{3}{2}$ and v'(x) > -4, so that

$$w'(x) > \frac{3}{2}A - 2B \ge \left(\frac{3}{2}\right)\left(\frac{4}{3}\right)^e - 2\left(\frac{7}{8}\right)^e > 0.$$

Step 9. Now our problem is reduced to the single-variable minimization of

$$T(x) = \Lambda(x, ..., x) = -2(2x/3)^{e} + 2\Phi(x)^{e} + \Psi(x)^{e}$$

for $x \in (1, 1 + 3^{e+1}/2^{2e+3})$. It is not hard to see that dT/dx vanishes if and only if

$$x^{3e} - \frac{3^e}{2^{e-2}}(x-1)(3x^2 - 4x + 2)^{e-1} - \frac{3^e}{2^{2e}}(2x-1)^{2e-3}$$

does. Meanwhile we claim that the second derivative of T is strictly positive on our interval: by dropping some nonnegative terms we see that

$$\frac{d^2T}{dx^2}(x) \ge -e(e-1)\frac{2^{e+1}x^{e-2}}{3^e} + e\frac{8(3-2x)}{x^4}\Phi(x)^{e-1} + e\frac{3-4x}{x^4}\Psi(x)^{e-1}.$$

Thus for e = 1, the second derivative is at least $(27 - 20x)/x^4$, which is strictly positive on our interval $(1, 1 + \frac{9}{32})$. For e = 2, we can use the fact that

$$0 \le \Psi(x) \le 1 \le \Phi(x)$$

on our interval $(1, 1 + \frac{27}{128})$ to show that the second derivative is at least

$$-\frac{16}{9} + \frac{54 - 40x}{x^4}$$

which is strictly positive on our interval. Finally, if $e \ge 3$, we have

$$1 + \frac{3^{e+1}}{2^{2e+3}} < \frac{7}{6}$$

and so on our interval $(1, 1+3^{e+1}/2^{2e+3})$ we have $8(3-2x)/x^4 > \frac{20}{7}$, $\Phi(x) \ge 1$, $(3-4x)/x^4 \ge -1$, and $0 \le \Psi(x) \le 1$, so that

$$\frac{d^2T}{dx^2}(x) \ge -\frac{8}{9} \left(\frac{7}{9}\right)^{e-2} e(e-1) + \frac{13}{7}e > 0.$$

This proves that there is a unique global minimum for this single-variable problem: the unique root b_e of

$$x^{3e} - \frac{3^e}{2^{e-2}}(x-1)(3x^2 - 4x + 2)^{e-1} - \frac{3^e}{2^{2e}}(2x-1)^{2e-1}$$

lying in $(1, 1 + 3^{e+1}/2^{2e+3})$. Thus we have found that global minima are obtained precisely when $\sigma_1 = \cdots = \sigma_e = b_e$ and $\tau_j \in \{\frac{1}{4}(1 - 2b_e) + m/2 : m \in \mathbb{Z}\}$ for each $j \in \{1, \ldots, e\}$.

7. Proof of Theorem 1.6

For A_e and B_e as defined in Theorem 1.6, we first show that $B_e < A_e$ for each $e \ge 1$. Given the minimizing conditions described in Theorems 1.4(ii) and 1.5(ii), it suffices to show that

$$-\frac{2^{e+1}}{3^e}x^e + 2\Omega\left(\frac{1}{x}, 0\right)^e + \Omega\left(\frac{1}{x}, \frac{1}{2x}\right)^e < -\frac{2^e}{3^e}x^e + 2\Omega\left(\frac{1}{x}, 0\right)^e.$$

for $x \in [1, \frac{3}{2}]$. This follows if $\Omega(1/x, 1/2x) < 2x/3$, or using the definition (3) of Ω , if $4x^3 - 12x^2 + 12x - 3 > 0$ for $x \in [1, \frac{3}{2}]$, which is routine to show.

Now we show that $B_{e_1+e_2} < B_{e_1}B_{e_2}$ for any $e_1, e_2 \ge 1$. We use a technical Lemma 7.1 below, which provides bounds on the B_e . It shows that if $e_1 \ge 5$ or $e_2 \ge 5$, then $B_{e_1}B_{e_2} > \sqrt[4]{2} > B_{e_1+e_2}$. So we may confine ourselves to the case where $1 \le e_1 \le e_2 \le 4$. If we define $B_0 = 1$, then the bounds in Lemma 7.1 also show us that

$$\frac{B_1}{B_0} > \frac{B_2}{B_1} > \frac{B_3}{B_2} > \frac{B_4}{B_3} > \frac{B_5}{B_4}.$$

Thus we note that $B_{e_1}B_{e_2} > B_{e_1-1}B_{e_2+1}$. If $e_1 + e_2 > 5$, we can repeat this argument to show that $B_{e_1}B_{e_2}$ is greater than $B_{e_1+e_2-5}B_5$, which we have already shown to exceed $B_{e_1+e_2}$. On the other hand, if $e_1 + e_2 \le 5$, repetition of the same argument produces $B_{e_1}B_{e_2} > B_0B_{e_1+e_2} = B_{e_1+e_2}$.

Lemma 7.1. For each $e \ge 1$, let B_e be the minimum asymptotic ratio of L^4 to L^2 norm achievable by a family of e-variable quadratic multiplicative character polynomials as described in Theorem 1.5(ii). Then

- (i) $\sqrt[4]{103/89} < B_1 < \sqrt[4]{22/19}$,
- (ii) $\sqrt[4]{86/65} < B_2 < \sqrt[4]{75/56}$,
- (iii) $\sqrt[4]{142/95} < B_3 < \sqrt[4]{116/77}$,
- (iv) $\sqrt[4]{100/61} < B_4 < \sqrt[4]{107/65}$,
- (v) $\sqrt[4]{7/4} < B_5 < \sqrt[4]{128/73}$, and
- (vi) $\sqrt[4]{7/4} < B_e < \sqrt[4]{2}$ for all $e \ge 6$.

Proof. By Section 6 of the proof of Theorem 1.5(ii) (see page 389), for each $e \ge 1$ the quantity B_e^4 is the minimum of the function

$$T(x) = -\frac{2^{e+1}}{3^e}x^e + 2\Omega\left(\frac{1}{x}, 0\right)^e + \Omega\left(\frac{1}{x}, \frac{1}{2x}\right)^e$$

on the interval $(1, 1+3^{e+1}/2^{2e+3})$, upon which the second derivative of *T* is shown to be positive. Thus if we can find x_1, x_2, x_3 in this interval with $x_1 < x_2 < x_3$ and $T(x_2) < \min(T(x_1), T(x_3))$, then we will have shown that the minimizing value of *x* lies in the interval (x_1, x_3) . Then we can use $B_e^4 \le T(x_2)$ for our upper bound and, by the monotonicity of $\Omega(1/x, 0)$ and $\Omega(1/x, 1/(2x))$, we can use

$$B_e^4 > -\frac{2^{e+1}}{3^e} x_3^e + 2\Omega\left(\frac{1}{x_1}, 0\right)^e + \Omega\left(\frac{1}{x_1}, \frac{1}{2x_1}\right)^e$$

as a lower bound. We use this technique to prove bounds in (i)–(v). The calculations done by hand are tedious, so here we simply state choices of the triple (x_1, x_2, x_3) that establish stricter bounds than the ones we claim above:

for B_1 use $\left(\frac{55}{52}, \frac{128}{121}, \frac{73}{69}\right)$, for B_2 use $\left(\frac{18}{17}, \frac{17}{16}, \frac{16}{15}\right)$, for B_3 use $\left(\frac{21}{20}, \frac{20}{19}, \frac{19}{18}\right)$, for B_4 use $\left(\frac{26}{25}, \frac{25}{24}, \frac{24}{23}\right)$, and for B_5 use $\left(\frac{36}{35}, \frac{35}{34}, \frac{34}{33}\right)$.

Henceforth assume that $e \ge 6$, and let b_e be the unique value in the interval $(1, 1 + 3^{e+1}/2^{2e+3})$ such that $T(b_e) = B_e^4$. Now note that $\Omega(1/x, 0) \ge 1$ and $\Omega(1/x, 1/(2x)) \ge \frac{1}{2}$ for $x \ge 1$, and that $2b_e/3 < (2^{15} + 3^7)/(3 \cdot 2^{14}) < 1$, so that

$$B_6^4 > -2\left(\frac{2^{15}+3^7}{3\cdot 2^{14}}\right)^6 + 2 + \frac{1}{2^6} > \frac{7}{4},$$

and if $e \ge 7$, then

$$B_e^4 > -2\left(\frac{2^{15}+3^7}{3\cdot 2^{14}}\right)^7 + 2 > \frac{7}{4}.$$

This proves our lower bound on B_e when $e \ge 6$.

To prove our upper bound on B_e when $e \ge 6$, write $b_e = 1 + c_e$ with $0 < c_e < 3^{e+1}/2^{2e+3}$, and use the definition (3) of Ω and the fact that $1/(1 + c_e) > 1 - c_e$ to estimate $B_e = T(b_e)$ as

$$B_e^4 < -\frac{2^{e+1}}{3^e} + 2(1+2c_e^2)^e + \frac{1}{2^e}(1+c_e)^{2e}$$

$$\leq -\frac{2^{e+1}}{3^e} + 2 + 4ec_e^2 + 2^{e+3}c_e^4 + \frac{1}{2^e}(1+c_e)^{2e},$$

where the second inequality follows from a crude approximation with the binomial expansion. Now note that

$$2^{e+3}c_e^4 < \frac{3^{4e+4}}{2^{7e+9}} < \frac{2^e}{6\cdot 3^e} \quad \text{and} \quad 4ec_e^2 < e\frac{3^{2e+2}}{2^{4e+4}} < \frac{5\cdot 2^e}{4\cdot 3^e},$$

so that

$$B_e^4 < 2 - \frac{7}{12} \cdot \frac{2^e}{3^e} + \frac{1}{2^e} (1 + c_e)^{2e},$$

which will imply $B_e^4 < 2$ if we can show that $(3(1 + c_e)^2/4)^e \le \frac{7}{12}$. Given that $c_e \le 3^7/2^{15} < \frac{1}{10}$, the quantity being raised to the *e*-th power on the left hand side is less than 1, so it suffices to show that $(3(\frac{11}{10})^2/4)^6 \le \frac{7}{12}$.

Appendix A: Proofs of auxiliary results

Here we collect, for the sake of completeness, technical results used in our proofs. The first is a bound on a character sum used in the proofs of Theorems 1.3(i), 1.4(i), and 1.5(i) in Section 5.

Lemma A.1. Let Γ be a finite abelian group, n a positive integer, and

 $\pi_1,\ldots,\pi_n\in\operatorname{Hom}(\mathbb{Z},\Gamma)$

such that $\operatorname{im} \pi_1 + \cdots + \operatorname{im} \pi_n$ is the internal direct sum of $\operatorname{im} \pi_1, \ldots, \operatorname{im} \pi_n$ in Γ . Let $\pi \in \operatorname{Hom}(\mathbb{Z}^n, \Gamma)$ with $\pi(u_1, \ldots, u_n) = \pi_1(u_1) + \cdots + \pi_n(u_n)$ for all $(u_1, \ldots, u_n) \in \mathbb{Z}^n$, and let $U = U_1 \times \cdots \times U_n$ be a box in \mathbb{Z}^n . Then

$$T = \sum_{\substack{\kappa,\lambda,\mu,\nu\in\widehat{\Gamma} \\ a+b=c+d}} \left| \sum_{\substack{a,b,c,d\in U \\ a+b=c+d}} \overline{\kappa(\pi(a))\lambda(\pi(b))} \mu(\pi(c))\nu(\pi(d)) \right|$$

is no greater than

$$64^{n}|\Gamma|^{4}\prod_{j=1}^{n}\max\left(1,\frac{|U_{j}|}{|\mathrm{im}\,\pi_{j}|}\right)^{3}\prod_{j=1}^{n}(1+\log|\mathrm{im}\,\pi_{j}|)^{3}$$

Proof. Write $K = \bigoplus_{j=1}^{n} \operatorname{im} \pi_j$, so that $\widehat{K} = \prod_{j=1}^{n} \operatorname{im} \pi_j$. Each character of K extends to $[\Gamma : K]$ characters of Γ , and for any $\eta \in \operatorname{im} \pi_j$, let $\eta' \in \widehat{\mathbb{Z}}$ be $\eta \circ \pi_j$,

so that

$$T = [\Gamma:K]^4 \prod_{j=1}^n \sum_{\substack{\kappa_j, \lambda_j, \mu_j, \nu_j \in \widehat{\operatorname{im}\pi_j}}} \left| \sum_{\substack{a_j, b_j, c_j, d_j \in U_j \\ a_j + b_j = c_j + d_j}} \overline{\kappa'_j(a_j)\lambda'_j(b_j)} \mu'_j(c_j)\nu'_j(d_j) \right|,$$

and so it suffices to prove the bound when n = 1 and π is surjective. In this case Γ is a finite cyclic group, which we identify with $\mathbb{Z}/m\mathbb{Z}$ by identifying $\pi(1)$ with $1 \in \mathbb{Z}/m\mathbb{Z}$. Then we set $\epsilon_a = \exp(2\pi i a/m)$ for every $a \in \mathbb{Z}/m\mathbb{Z}$, and note that $\widehat{\Gamma}$ is the set of maps $a \mapsto \epsilon_{xa}$ with $x \in \mathbb{Z}/m\mathbb{Z}$. Thus,

$$T = \sum_{\substack{w, x, y, z \in \mathbb{Z}/m\mathbb{Z}}} \left| \sum_{\substack{a, b, c, d \in U \\ a+b=c+d}} \epsilon_{-wa-xb+yc+zd} \right|.$$

U is a set of consecutive integers in \mathbb{Z} , and note that translation of *U* does not influence the magnitude of the inner sum in *T*, so without loss of generality, we assume that $U = \{0, 1, ..., |U| - 1\}$. Then reparametrize the outer sum in *T* with x' = w - x, y' = y - w, z' = z - w, and *w* to obtain

$$T = m \sum_{\substack{x', y', z' \in \mathbb{Z}/m\mathbb{Z}}} \left| \sum_{\substack{a, b, c, d \in U \\ a+b=c+d}} \epsilon_{-x'b+y'c+z'd} \right|,$$

which is not more than $64m \max(m, |U|)^3 (1 + \log m)^3$ by [Jedwab et al. 2013b, Lemma 2.2].

The next result is used in the proofs of Section 5 to understand the asymptotic behavior of the L^2 norm for multiplicative character polynomials.

Lemma A.2. Let $\{f_i\}_{i \in I}$ be a size-stable family of e-variable multiplicative character polynomials with \mathbb{F}_{q_i} , S_i , t_i , and α_i the field, support, translation, and arrangement, respectively, of f_i for each $i \in I$. Then there is a Q and an N such that for all $i \in I$ with $q_i \ge Q$, we have $|S_i \cap (\ker \alpha_i - t_i)| \le N$. Thus $|S_i \cap \ker(\alpha_i - t_i)|/|S_i| \to 0$ and $||f_i||_2^2/|S_i| \to 1$ as $q_i \to \infty$.

Proof. Suppose that the limiting sizes for our size-stable family $\{f_i\}_{i \in I}$ of multiplicative character polynomials are $\sigma_1, \ldots, \sigma_e$. For each $\iota \in I$, let χ_ι be the character of f_ι , so that $f_\iota(z) = \sum_{s \in S_\iota} \chi_\iota(\alpha_\iota(s + t_\iota)) z^s$, and let $p_\iota = q_\iota^{1/e}$, which is the characteristic of the field \mathbb{F}_{q_ι} of f_ι . Since α_ι is an epimorphism, its restriction to each $p_\iota \times \cdots \times p_\iota$ cubical box in \mathbb{Z}^e is a bijection to \mathbb{F}_{q_ι} , and by the definition of size-stability, there is some Q such that for every $q_\iota \ge Q$, the support $S_\iota = S_{\iota,1} \times \cdots \times S_{\iota,e}$ can be covered with $N = \prod_{j=1}^e (\lfloor \sigma_j \rfloor + 1)$ such cubes, each of which contains one point of ker $\alpha_\iota - t_\iota$, so $|S_\iota \cap (\ker \alpha_\iota - t_\iota)| \le N$. Since the family is size-stable, $|S_\iota| \to \infty$ as $q_\iota \to \infty$. The squared L^2 norm of a polynomial is the sum of the squared magnitudes of its coefficients, and $\chi_\iota(\alpha_\iota(s + t_\iota)) = 0$ for $s \in S_\iota \cap (\ker \alpha_\iota - t_\iota)$ while

 $|\chi_t(\alpha_t(s+t_t))| = 1$ for all other $s \in S_t$. Thus $||f_t||_2^2 = |S_t \setminus (S_t \cap (\ker \alpha_t - t_t))|$, and so $||f_t||_2^2/|S_t| \to 1$ as $q_t \to \infty$.

Recall from footnote 3 of Section 1C that we sometimes wish to obtain a Littlewood polynomial from a quadratic character polynomial $f(z) = \sum_{s \in S} \chi(\alpha(s+t))z^s$, but f may have some coefficients equal to 0 because an extended nontrivial multiplicative character χ has $\chi(0) = 0$. More generally, if χ is a nontrivial multiplicative character (not necessarily quadratic), we may wish to obtain from f a polynomial with coefficients of unit magnitude. So we replace the zero coefficient for each z^s such that $s \in S \cap (\ker \alpha - t)$ with a coefficient of unit magnitude. We may choose each replacement coefficient independently of the others, and any polynomial gresulting from such replacements is called a *unimodularization* of f. The following corollary to Lemma A.2 shows that unimodularizing all the polynomials in a sizestable family of multiplicative character polynomials does not affect asymptotic ratio $||f||_4/||f||_2$.

Corollary A.3. Let $\{f_i\}_{i \in I}$ be a size-stable family of multiplicative character polynomials over fields $\{\mathbb{F}_{q_i}\}_{i \in I}$, and let g_i be a unimodularization of f_i for each $i \in I$. If r is a real number with $r \ge 2$ or if $r = \infty$, then $\|f_i\|_r / \|g_i\|_r \to 1$ as $q_i \to \infty$.

Proof. If $u \in \mathbb{C}$ with |u| = 1 and if $s = (s_1, \ldots, s_e) \in \mathbb{Z}^e$, then L^r norm of $uz^s = uz_1^{s_1} \cdots z_e^{s_e}$ is 1. By Lemma A.2 there is an N and a Q such that whenever $q_t \geq Q$, the two polynomials f_t and g_t differ by the sum of N or fewer such monomials, and so by the L^r triangle inequality, the difference between $||f_t||_r$ and $||g_t||_r$ cannot be greater in magnitude than N. Now $||g_t||_r \geq ||g_t||_2 = \sqrt{|S_t|}$ by monotonicity of L^r norms and the fact that the squared L^2 norm of a polynomial is the sum of the squared magnitudes of its coefficients, and $|S_t| \to \infty$ as $q_t \to \infty$ for a size-stable family.

The next result is used in Section 6 to find the limiting translations that globally minimize the asymptotic ratio of L^4 to L^2 norm for quadratic character polynomials.

Lemma A.4. Let x be a fixed nonzero value in \mathbb{R} and let y vary over \mathbb{R} .

- (i) If $|x| \ge 2$, the function $\Omega(x, y)$, considered as a function of y, achieves a global minimum value of 0 for $y \in \bigcup_{m \in \mathbb{Z}} [m|x|+1, (m+1)|x|-1]$ and for no other value of y.
- (ii) If $0 < |x| \le 2$, the function $\Omega(x, y)$, considered as a function of y, achieves a global minimum value of

$$\Omega\left(x,\frac{x}{2}\right) = \sum_{n \in \mathbb{Z}} \max\left(0, 1 - \left|\left(n + \frac{1}{2}\right)x\right|\right)^2$$

for $y \in \{x(m + \frac{1}{2}) : m \in \mathbb{Z}\}$ and for no other value of y.

Proof. For part (i), note that all the terms of $\Omega(x, y)$ are nonnegative. Since $\Omega(x, y) = \Omega(-x, y)$, we may assume without loss of generality that $x \ge 2$, and then the term max $(0, 1-|xn-y|)^2$ is nonvanishing if and only if (y-1)/x < n < (y+1)/x. Thus we can obtain a global minimum value of 0 if we can arrange that no integer lie in the interval ((y-1)/x, (y+1)/x). If *m* is the greatest integer lying below this interval (so that $y \ge mx + 1$), then for the next integer m + 1 to lie above the interval, it is necessary and sufficient that $y \le (m+1)x - 1$.

For part (ii), it is clear from the definition (3) of Ω that $\Omega(-x, y) = \Omega(x, y)$, $\Omega(x, -y) = \Omega(x, y)$, and $\Omega(x, y) = \Omega(x, y+x)$. So without loss of generality we may restrict our attention to the case where $0 < x \le 2$ and $0 \le y \le x/2$. In this case

$$\Omega(x, y) = \sum_{\lceil (y-1)/x \rceil \le n \le 0} (y-1-nx)^2 + \sum_{0 < n \le \lfloor (y+1)/x \rfloor} (y+1-nx)^2,$$

and we reparametrize the sums to obtain

$$\Omega(x, y) = \sum_{0 \le n \le \lfloor (1-y)/x \rfloor} (y - 1 + nx)^2 + \sum_{0 \le n \le \lfloor (1+y-x)/x \rfloor} (y - x + 1 - nx)^2,$$

and calculate

$$\frac{\partial}{\partial y}\Omega(x, y) = \sum_{0 \le n \le \lfloor (1-y)/x \rfloor} 2(y-1+nx) + \sum_{0 \le n \le \lfloor (1+y-x)/x \rfloor} 2(y-x+1-nx)$$
$$= 2\left\lfloor \frac{y+1}{x} \right\rfloor (2y-x) + \sum_{\lfloor (1+y-x)/x \rfloor < n \le \lfloor (1-y)/x \rfloor} 2(y-1+nx),$$

because (1-y)/x is greater than or equal to (1+y-x)/x, and note for the remainder of this proof that their difference is at most 1. Since $0 \le y \le x/2 \le 1$, we can see that both terms in the last expression for our partial derivative are nonpositive, with the summation over *n* strictly negative if y < x - 1, and the other term is strictly negative if $x - 1 \le y < x/2$. Thus our partial derivative is strictly negative for $0 \le y < x/2$, and so for our ranges of *x* and *y* values, the unique minimum is obtained when y = x/2.

The last two results are used in Section 6 to show that a large limiting size will make the ratio of L^4 to L^2 norm large.

Lemma A.5. If $\{f_i\}_{i \in I}$ is a size-stable family, with limiting size σ , of additive character polynomials over fields $\{\mathbb{F}_{q_i}\}_{i \in I}$, then

$$\liminf_{q_l \to \infty} \frac{\|f_l\|_4^4}{\|f_l\|_2^4} \ge \Omega\left(\frac{1}{\sigma}, 0\right).$$

Proof. For each $\iota \in I$, let f_{ι} have character ψ_{ι} , support S_{ι} , translation t_{i} , and arrangement α_{ι} , so that $f_{\iota}(z) = \sum_{s \in S_{\iota}} \psi_{\iota}(\alpha_{\iota}(s + t_{i}))z^{s}$. When we confine the values

of z to the complex unit circle, we have

$$\overline{f_{\iota}(z)} = \sum_{s \in S_{\iota}} \overline{\psi_{\iota}(\alpha_{\iota}(s+t_{i}))} z^{-s}.$$

We can consider $f_t(z)$ and $\overline{f_t(z)}$ formally as elements of $\mathbb{C}[z, z^{-1}]$, and view $||f_t||_4^4$ as the sum of the squared magnitudes of the coefficients of $f_t(z)\overline{f_t(z)}$. The coefficient of z^s in $f_t(z)\overline{f_t(z)}$ is

$$\sum_{\substack{u,v\in S_i\\u-v=s}}\psi_i(\alpha_i(u+t_i))\overline{\psi_i(\alpha_i(v+t_i))}.$$

Since α_i is an epimorphism from \mathbb{Z} to $\mathbb{F}_{q_i}^*$, we see that $\psi_t(\alpha_t(u+t_i)) = \psi_t(\alpha_t(v+t_i))$ whenever $u \equiv v \pmod{q_i - 1}$. Thus if $s \equiv 0 \pmod{q_i - 1}$, the coefficient of z^s in $f_t(z)\overline{f_t(z)}$ is equal to $|S_t \cap (s + S_t)|$. Since $||f_t||_4^4$ is the sum of the squared magnitudes of the coefficients of $f_t(z)\overline{f_t(z)}$ while $||f_t||_2^2$ is the coefficient of z^0 of the same, we have

$$\frac{\|f_{\iota}\|_{4}^{4}}{\|f_{\iota}\|_{2}^{4}} \geq \frac{\sum_{n \in \mathbb{Z}} |S_{\iota} \cap (n(q_{\iota} - 1) + S_{\iota})|^{2}}{|S_{\iota}|^{2}},$$

and then we note that $|S_{\iota} \cap (n(q_{\iota}-1)+S_{\iota})| = \max(0, |S_{\iota}|-|n(q_{\iota}-1)|)$ and apply the size-stability limit $|S_{\iota}|/(q_{\iota}-1) \to \sigma$ as $q_{\iota} \to \infty$.

Lemma A.6. If $\{f_i\}_{i \in I}$ is a size-stable family, with limiting sizes $\sigma_1, \ldots, \sigma_e$, of *e*-variable multiplicative character polynomials over fields $\{\mathbb{F}_{q_i}\}_{i \in I}$, then

$$\liminf_{q_l \to \infty} \frac{\|f_l\|_4^4}{\|f_l\|_2^4} \ge \prod_{j=1}^e \Omega\left(\frac{1}{\sigma_j}, 0\right).$$

Proof. For each $\iota \in I$, let f_{ι} have character χ_{ι} , support $S_{\iota} \subseteq \mathbb{Z}^{e}$, translation $t_{i} \in \mathbb{Z}^{e}$, and arrangement α_{ι} , so that $f_{\iota}(z_{1}, \ldots, z_{e}) = \sum_{s \in S_{\iota}} \chi_{\iota}(\alpha_{\iota}(s + t_{i}))z^{s}$, where we write z^{s} for $z_{1}^{s_{1}} \cdots z_{e}^{s_{e}}$ when $s = (s_{1}, \ldots, s_{e})$. Our proof runs the same as that of the previous lemma for additive character polynomials once we replace ψ_{ι} with χ_{ι} , but we must take care of the fact that $\chi_{\iota}(\alpha_{\iota}(s + t_{i})) = 0$ when $s \in -t_{i} + \ker \alpha_{\iota}$; otherwise, the coefficients are of unit magnitude. And of course the polynomials are in *e* variables and the coefficients have periodicity *p* in each direction. Thus if we define $V_{\iota} = S_{\iota} \setminus (-t_{i} + \ker \alpha_{\iota})$ we have

$$\frac{\|f_{\iota}\|_{4}^{4}}{\|f_{\iota}\|_{2}^{4}} \geq \frac{\sum_{n \in \mathbb{Z}^{e}} |V_{\iota} \cap (np_{\iota} + V_{\iota})|^{2}}{|V_{\iota}|^{2}},$$

but Lemma A.2 can be used to show that the ratio $|V_l \cap (np_l + V_l)|/|V_l|$ has the same limit as $|S_l \cap (np_l + S_l)|/|S_l|$ as $q_l \to \infty$.

Acknowledgements

The author thanks Jonathan Jedwab and Kai-Uwe Schmidt for helpful suggestions on the presentation of these results.

References

- [Baker and Montgomery 1990] R. C. Baker and H. L. Montgomery, "Oscillations of quadratic *L*-functions", pp. 23–40 in *Analytic number theory* (Allerton Park, IL, 1989), edited by B. C. Berndt et al., Progr. Math. **85**, Birkhäuser, Boston, MA, 1990. MR 91k:11071 Zbl 0718.11039
- [Bernasconi 1987] J. Bernasconi, "Low autocorrelation binary sequences: Statistical mechanics and configuration space analysis", *J. Phys. France* **48**:4 (1987), 559–567.
- [Bömer and Antweiler 1993] L. Bömer and M. Antweiler, "Optimizing the aperiodic merit factor of binary arrays", *Signal Process.* **30**:1 (1993), 1–13. Zbl 0769.93022
- [Borwein 2002] P. Borwein, *Computational excursions in analysis and number theory*, CMS Books in Mathematics **10**, Springer, New York, 2002. MR 2003m:11045 Zbl 1020.12001
- [Borwein and Choi 2000] P. Borwein and K.-K. S. Choi, "Merit factors of character polynomials", *J. London Math. Soc.* (2) **61**:3 (2000), 706–720. MR 2002d:11103 Zbl 1011.11059
- [Borwein and Choi 2002] P. Borwein and K.-K. S. Choi, "Explicit merit factor formulae for Fekete and Turyn polynomials", *Trans. Amer. Math. Soc.* **354**:1 (2002), 219–234. MR 2002i:11065 Zbl 1010.11017
- [Borwein et al. 2001] P. Borwein, K.-K. S. Choi, and S. Yazdani, "An extremal property of Fekete polynomials", *Proc. Amer. Math. Soc.* **129**:1 (2001), 19–27. MR 2001j:11061 Zbl 0987.11010
- [Conrey et al. 2000] B. Conrey, A. Granville, B. Poonen, and K. Soundararajan, "Zeros of Fekete polynomials", Ann. Inst. Fourier (Grenoble) 50:3 (2000), 865–889. MR 2001h:11108 Zbl 1007.11053
- [Erdős 1957] P. Erdős, "Some unsolved problems", *Michigan Math. J.* **4** (1957), 291–300. MR 20 #5157 Zbl 0081.00102
- [Erdős 1962] P. Erdős, "An inequality for the maximum of trigonometric polynomials", *Ann. Polon. Math.* **12** (1962), 151–154. MR 25 #5330 Zbl 0106.27702
- [Fekete and Pólya 1912] M. Fekete and G. Pólya, "Über ein Problem von Laguerre", *Rend. Circ. Mat. Palermo* **34**:1 (1912), 89–120.
- [Golay 1977] M. J. E. Golay, "The merit factor of long low autocorrelation binary sequences", *IEEE Trans. Inform. Theory* **23**:1 (1977), 43–51.
- [Gulliver and Parker 2005] T. Gulliver and M. Parker, "The multivariate merit factor of a Boolean function", in *Proceedings of the IEEE ITSOC Information Theory Workshop 2005 on Coding and Complexity* (Rotorua, New Zealand, 2005), IEEE, New York, 2005.
- [Høholdt 2006] T. Høholdt, "The merit factor problem for binary sequences", pp. 51–59 in *Applied algebra, algebraic algorithms and error-correcting codes*, edited by M. Fossorier et al., Lecture Notes in Comput. Sci. **3857**, Springer, Berlin, 2006. MR 2007c:94080 Zbl 1125.94008
- [Høholdt and Jensen 1988] T. Høholdt and H. E. Jensen, "Determination of the merit factor of Legendre sequences", *IEEE Trans. Inform. Theory* **34**:1 (1988), 161–164. Zbl 0652.40006
- [Jedwab 2005] J. Jedwab, "A survey of the merit factor problem for binary sequences", pp. 19–21 in *Sequences and their applications* (SETA 2004), edited by T. Helleseth et al., Lecture Notes in Computer Science **3486**, Springer, Berlin, 2005.

DANIEL J. KATZ

- [Jedwab and Schmidt 2010] J. Jedwab and K.-U. Schmidt, "Appended *m*-sequences with merit factor greater than 3.34", pp. 204–216 in *Sequences and their applications – SETA 2010*, edited by C. Carlet and A. Pott, Lecture Notes in Comput. Sci. 6338, Springer, Berlin, 2010. MR 2012m:94245 Zbl pre05781409
- [Jedwab et al. 2013a] J. Jedwab, D. J. Katz, and K.-U. Schmidt, "Advances in the merit factor problem for binary sequences", *J. Combin. Theory Ser. A* **120**:4 (2013), 882–906. MR 3022619 Zbl 06154122
- [Jedwab et al. 2013b] J. Jedwab, D. J. Katz, and K.-U. Schmidt, "Littlewood polynomials with small L^4 norm", *Adv. Math.* **241** (2013), 127–136. MR 3053707
- [Jensen and Høholdt 1989] H. E. Jensen and T. Høholdt, "Binary sequences with good correlation properties", pp. 306–320 in *Applied algebra, algebraic algorithms and error-correcting codes* (Menorca, 1987), edited by L. Huguet and A. Poli, Lecture Notes in Comput. Sci. **356**, Springer, Berlin, 1989. MR 90h:94009 Zbl 0675.94010
- [Jensen et al. 1991] J. M. Jensen, H. E. Jensen, and T. Høholdt, "The merit factor of binary sequences related to difference sets", *IEEE Trans. Inform. Theory* **37**:3, part 1 (1991), 617–626. MR 92j:94009 Zbl 0731.94011
- [Kahane 1980] J.-P. Kahane, "Sur les polynômes à coefficients unimodulaires", *Bull. London Math. Soc.* **12**:5 (1980), 321–342. MR 82a:30003 Zbl 0443.30005
- [Lidl and Niederreiter 1997] R. Lidl and H. Niederreiter, *Finite fields*, 2nd ed., Encyclopedia of Mathematics and its Applications 20, Cambridge University Press, 1997. MR 97i:11115 Zbl 0866.11069
- [Littlewood 1966] J. E. Littlewood, "On polynomials $\sum^{n} \pm z^{m}$, $\sum^{n} e^{\alpha_{m}i} z^{m}$, $z = e^{\theta_{i}}$ ", J. London *Math. Soc.* **41** (1966), 367–376. MR 33 #4237 Zbl 0142.32603
- [Littlewood 1968] J. E. Littlewood, *Some problems in real and complex analysis*, D. C. Heath and Co. Raytheon Education Co., Lexington, MA, 1968. MR 39 #5777 Zbl 0185.11502
- [Montgomery 1980] H. L. Montgomery, "An exponential polynomial formed with the Legendre symbol", *Acta Arith.* **37** (1980), 375–380. MR 82a:10041 Zbl 0369.10024
- [Newman and Byrnes 1990] D. J. Newman and J. S. Byrnes, "The L^4 norm of a polynomial with coefficients ± 1 ", *Amer. Math. Monthly* **97**:1 (1990), 42–45. MR 91d:30006
- [Pólya 1919] G. Pólya, "Verschiedene Bemerkungen zur Zahlentheorie", *Jahresber. Dtsch. Math.-Ver.* **28** (1919), 31–40. JFM 47.0882.06
- [Schmidt 2011] K.-U. Schmidt, "The merit factor of binary arrays derived from the quadratic character", *Adv. Math. Commun.* **5**:4 (2011), 589–607. MR 2012k:94119 Zbl 1238.94025
- [Weil 1948] A. Weil, "On some exponential sums", *Proc. Nat. Acad. Sci. U. S. A.* **34** (1948), 204–207. MR 10,234e Zbl 0032.26102

Received June 13, 2012.

DANIEL J. KATZ DEPARTMENT OF MATHEMATICS SIMON FRASER UNIVERSITY 8888 UNIVERSITY DRIVE BURNABY BC V5A 1S6 CANADA dkatz@sfu.ca
DEGREE-THREE SPIN HURWITZ NUMBERS

JUNHO LEE

Gunningham (2012) calculated all spin Hurwitz numbers in terms of combinatorics of the Sergeev algebra. Here we use a spin curve degeneration to obtain a recursion formula for degree-three spin Hurwitz numbers.

Let *D* be a complex curve of genus *h* and *N* be a theta characteristic on *D*, that is, $N^2 = K_D$. The pair (D, N) is called a *spin curve* of genus *h* with parity $p \equiv h^0(N) \pmod{2}$. For i = 1, ..., k, let $m^i = (m_1^i, ..., m_{\ell_i}^i)$ be an odd partition of d > 0, namely, all components m_j^i are odd. Fix *k* points $q^1, ..., q^k$ in *D* and consider degree-*d* maps $f : C \to D$ from possibly disconnected domains *C* of Euler characteristic χ that are ramified only over the fixed points q^i with ramification data m^i . Observe that the Riemann–Hurwitz formula shows

(0-1)
$$2d(1-h) - \chi + \sum_{i=1}^{k} (\ell(m^{i}) - d) = 0,$$

where $\ell(m^i) = \ell_i$ is the length of m^i . By the Hurwitz formula, the twisted line bundle

(0-2)
$$L_f = f^* N \otimes \mathbb{O}\left(\sum_{i,j} \frac{1}{2}(m_j^i - 1)x_j^i\right)$$

is a theta characteristic on *C* where $f^{-1}(q^i) = \{x_j^i\}_{1 \le j \le \ell_i}$ and *f* has multiplicity m_i^i at x_j^i . We define the parity p(f) of a map *f* by

$$(0-3) p(f) \equiv h^0(L_f) \pmod{2}.$$

Given odd partitions m^1, \ldots, m^k of d, the spin Hurwitz number of genus h and parity p is defined as a (weighted) sum of (ramified) covers f satisfying (0-1) with sign determined by the parity p(f):

(0-4)
$$H^{h,p}_{m^1,\dots,m^k} = \sum_f \frac{(-1)^{p(f)}}{|\operatorname{Aut}(f)|}$$

MSC2010: 14N35, 53D45.

Keywords: local Gromov-Witten invariants, spin curve, Schiffer variation.

Eskin, Okounkov, and Pandharipande [Eskin et al. 2008] calculated the genus h = 1 and odd parity spin Hurwitz numbers in terms of characters of the Sergeev group. Gunningham [2012] calculated all spin Hurwitz numbers in terms of combinatorics of the Sergeev algebra.

The trivial partition (1^d) of *d* is a partition whose components are all 1. If $m^k = (1^d)$, *f* has no ramification points over the fixed point q^k and hence we have

(0-5)
$$H^{h,p}_{m^1,\dots,m^{k-1},(1^d)} = H^{h,p}_{m^1,\dots,m^{k-1}}$$

When all partitions $m^i = (1^d)$, denote the spin Hurwitz numbers (0-4) by $H_d^{h,p}$. These are dimension-zero local GW invariants $GT_d^{\text{loc},h,p}$ of spin curve (D, N) that give all dimension-zero GW invariants of Kähler surfaces with a smooth canonical divisor; see [Kiem and Li 2007; 2011; Lee and Parker 2007; Maulik and Pandharipande 2008]. For notational simplicity, we set $H_{(3)^0}^{h,p} = H_3^{h,p}$ and for $k \ge 1$ write

$$H^{h, p}_{(3)^k}$$

for the spin Hurwitz numbers $H_{(3),...,(3)}^{h,p}$ with the same k partitions (3). Since there are two odd partitions (1³) and (3) of d = 3, by (0-5) it suffices to compute $H_{(3)^k}^{h,p}$ for $k \ge 0$. The aim of this paper is to use a spin curve degeneration to obtain the following recursion formula.

Theorem 0.1. If $h = h_1 + h_2$ and $p \equiv p_1 + p_2 \pmod{2}$, then, for $k_1 + k_2 = k$,

(0-6)
$$H_{(3)^{k}}^{h,p} = 3! H_{(3)^{k_{1}}}^{h_{1},p_{1}} \cdot H_{(3)^{k_{2}}}^{h_{2},p_{2}} + 3 H_{(3)^{k_{1}+1}}^{h_{1},p_{1}} \cdot H_{(3)^{k_{2}+1}}^{h_{2},p_{2}}$$

One can use Theorem 0.1 and the result of [Eskin et al. 2008] to explicitly compute the spin Hurwitz numbers of degree d = 3. In Proposition 7.1, we show that

(0-7)
$$H_{(3)^k}^{h,\pm} = 3^{2h-2}[(-1)^k 2^{k+h-1} \pm 1],$$

where + and - denote the even and odd parities. When the degree *d* is 1 or 2, the dimension-zero local GW invariants are given by the formulas

$$GT_1^{\text{loc},h,\pm} = \pm 1$$
 and $GT_2^{\text{loc},h,\pm} = \pm 2^{h-1};$

see Lemma 2.6 of [Lee 2013]. Since $GT_d^{\text{loc},h,p} = H_d^{h,p}$ as mentioned above, formula (0-7) shows $GT_3^{\text{loc},h,\pm} = 3^{2h-2}(2^{h-1}\pm 1).$

This calculation is, in fact, the main motivation for the paper.

In Section 1, we express the degree-*d* spin Hurwitz numbers (0-4) in terms of relative GW moduli spaces. We can then apply a degeneration method for a family of curves $\mathfrak{D} \to \Delta$ where the central fiber D_0 is a nodal curve and the general fiber

 D_{λ} ($\lambda \neq 0$) is a smooth curve. Section 2 describes the relative moduli space \mathcal{M}_0 of maps f into the nodal curve D_0 . In Section 3, we show that the union over $\lambda \in \Delta$ of relative moduli spaces \mathcal{M}_{λ} of maps into D_{λ} consists of connected components $\mathscr{X}_{m,f} \to \Delta$ containing $f \in \mathcal{M}_0$. Here m is the ramification data of f over nodes of D_0 such that $d - \ell(m)$ is even.

The (ordinary) Hurwitz numbers are sums of (ramified) maps modulo automorphism without sign. One can easily obtain a recursion formula for Hurwitz numbers by counting maps in the general fiber of $\mathscr{Z}_{m,f} \to \Delta$. For spin Hurwitz numbers, one needs to calculate parities of maps induced from a fixed spin structure on the family of curves \mathfrak{D} .

The novelty of our approach is to apply a Schiffer variation for the parity calculation. The space $\mathscr{X}_{m,f}$ is, in general, not smooth. In Section 4, we construct a smooth model for $\mathscr{X}_{m,f}$ by Schiffer variation. In Section 5, we use the smooth model to twist the pullback of the spin structure on \mathfrak{D} . When the degree *d* equals 3, the partition *m* is odd, either (1³) or (3). In this case, a suitable twisting immediately yields a required parity calculation. We prove Theorem 0.1 in Section 6 and formula (0-7) in Section 7.

For higher degree $d \ge 4$, the partition *m* may not be odd! A new parity calculation is needed. In [Lee and Parker 2012], we generalized the recursion formula (0-6) for higher-degree spin Hurwitz numbers by employing additional geometric analysis arguments for parity calculations.

1. Dimension zero relative GW moduli spaces

In this section, we express the spin Hurwitz numbers (0-4) in terms of dimensionzero relative GW moduli spaces. We follow the definitions of [Ionel and Parker 2003] for the relative GW theory.

Let *D* be a smooth curve of genus *h* and let $V = \{q^1, \ldots, q^k\}$ be a fixed set of points on *D*. Given partitions m^1, \ldots, m^k of *d*, a degree-*d* holomorphic map $f: C \to D$ from a possibly disconnected curve *C* is called *V*-regular with contact vectors m^1, \ldots, m^k if $f^{-1}(V)$ consists of $\sum \ell(m^i)$ contact marked points x_j^i $(1 \le j \le \ell(m^i))$ with $f(x_j^i) = q^i$ such that *f* has ramification index (or multiplicity) m_j^i at x_j^i . Two *V*-regular maps $(f, C; \{x_j^i\})$ and $(\tilde{f}, \tilde{C}; \{\tilde{x}_j^i\})$ are equivalent if they are isomorphic, that is, there is a biholomorphism $\sigma : C \to \tilde{C}$ with $\tilde{f} \circ \sigma = f$ and $\sigma(x_i^i) = \tilde{x}_j^i$ for all *i*, *j*. The relative moduli space

(1-1)
$$\mathcal{M}^{V}_{\chi,m^{1},\ldots,m^{k}}(D,d)$$

consists of equivalence classes of *V*-regular maps $(f, C; \{x_j^i\})$ with the Euler characteristic $\chi(C) = \chi$ and with contact vectors m^1, \ldots, m^k . Since no confusion can

JUNHO LEE

arise, we regard a point in the space (1-1) as a V-regular map $(f, C; \{x_j^i\})$. For simplicity, we often write a V-regular map $(f, C; \{x_j^i\})$ simply as f.

The (formal) complex dimension of the space (1-1) is given by the left side of the Riemann–Hurwitz formula (0-1):

(1-2)
$$2d(1-h) - \chi - \sum_{i=1}^{k} (d - \ell(m^{i})).$$

Suppose this dimension is zero. Then, for each *V*-regular map $(f, C; \{x_j^i\})$ in (1-1), forgetting the contact marked points x_j^i gives a (ramified) cover *f* that is ramified only over fixed points q^i and satisfies (0-1). The automorphism group Aut(f) of a (ramified) cover *f* consists of automorphisms $\sigma \in Aut(C)$ with $f \circ \sigma = f$. The automorphism group Aut(f, V) of a *V*-regular map $(f, C; \{x_j^i\})$ consists of automorphisms $\sigma \in Aut(f)$ with $\sigma(x_j^i) = x_j^i$ for all *i*, *j*.

For a partition *m* of *d*, let Aut(*m*) be the subgroup of symmetric group $S_{\ell(m)}$ permuting equal parts of the partition *m*.

Lemma 1.1. Let m^1, \ldots, m^k be as above and suppose the dimension (1-2) is zero.

- (a) If $m^i = (1^d)$ for some $1 \le i \le k$, Aut(f, V) is trivial for all f in (1-1).
- (b) If m^1, \ldots, m^k are all odd partitions,

$$H_{m^1,...,m^k}^{h,p} = \frac{1}{\prod_{i=1}^k |\operatorname{Aut}(m^i)|} \sum \frac{(-1)^{p(f)}}{|\operatorname{Aut}(f,V)|}$$

where the sum is over all f in (1-1) and p(f) is the parity (0-3).

Proof. Let $(f, C; \{x_j^i\})$ be a *V*-regular map in (1-1) and $\sigma \in \text{Aut}(f, V)$. If $m^i = (1^d)$, the set of branch points *B* of *f* is a subset of $V \setminus \{q^i\}$ and the restriction of σ to $C \setminus f^{-1}(B)$ is a covering transformation that fixes contact marked points x_1^i, \ldots, x_d^i . Noting $f^{-1}(B)$ is finite, we conclude that σ is an identity map on *C*. This proves (a).

As mentioned above, forgetting contact marked points x_j^i gives a (ramified) cover f satisfying (0-1). Conversely, given a (ramified) cover f satisfying (0-1), one can mark a point over q^i with ramification index m_j^i as a contact marked point x_j^i . Such marking gives V-regular maps $(f, C; \{x_j^i\})$ in $\prod_{i=1}^k |\operatorname{Aut}(m^i)|$ ways. Observe that $(f, C; \{x_j^i\})$ and $(f, C; \{\sigma(x_j^i)\})$ are isomorphic for each $\sigma \in \operatorname{Aut}(f)$ and that $\operatorname{Aut}(f, V)$ is a normal subgroup of $\operatorname{Aut}(f)$. Consequently, the quotient group $G = \operatorname{Aut}(f) / \operatorname{Aut}(f, V)$ acts freely on the set of V-regular maps $(f, C; \{x_j^i\})$ obtained by the (ramified) cover f. Its orbits give $\prod_{i=1}^k |\operatorname{Aut}(m^i)| / |G|$ points (that is, equivalence classes of V-regular maps) in the space (1-1), each of which has the same automorphism group $\operatorname{Aut}(f, V)$. Now (b) follows from counting maps with the parity of map modulo automorphisms.

2. Maps into a nodal curve

Let $D_0 = D_1 \cup E \cup D_2$ be a connected nodal curve of (arithmetic) genus h with two nodes p^1 and p^2 such that, for $i = 1, 2, E = \mathbb{P}^1$ meets D_i at node p^i and D_i has genus h_i with $h_1 + h_2 = h$. In this section, we consider maps into D_0 that are relevant to our subsequent discussion.

Below, we fix d, h, χ , and odd partitions m^1, \ldots, m^k of d so that the Riemann–Hurwitz formula (0-1) holds, or equivalently, the dimension formula (1-2) is zero. For each partition m of d, consider the product space

$$\mathcal{P}_{m} = \mathcal{M}_{\chi_{1},(1^{d}),m^{1},\dots,m^{k_{1}},m}^{V_{1}}(D_{1},d) \times \mathcal{M}_{\chi_{0},m,(1^{d}),m}^{V_{0}}(E,d) \times \mathcal{M}_{\chi_{2},m,m^{k_{1}+1},\dots,m^{k},(1^{d})}^{V_{2}}(D_{2},d)$$

where

$$V_1 = \{q^{k+1}, q^1, \dots, q^{k_1}, p^1\}, \quad V_0 = \{p^1, q^{k+2}, p^2\}, \quad V_2 = \{p^2, q^{k_1+1}, \dots, q^k, q^{k+3}\}$$

and

(2-1)
$$\chi_1 + \chi_0 + \chi_2 - 4\ell(m) = \chi.$$

For simplicity, let \mathcal{M}_m^1 , \mathcal{M}_m^0 , and \mathcal{M}_m^2 denote the first, second, and third factors of \mathcal{P}_m .

Lemma 2.1. If $\mathcal{P}_m \neq \emptyset$, the spaces \mathcal{M}_m^1 , \mathcal{M}_m^0 , and \mathcal{M}_m^2 have dimension zero. Consequently, $\chi_0 = 2\ell(m)$ and $d - \ell(m)$ is even.

Proof. Each \mathcal{M}_m^i $(0 \le i \le 2)$ has nonnegative dimension by the Riemann–Hurwitz formula. The formula (2-1) and our assumption that the dimension (1-2) is zero thus imply that each \mathcal{M}_m^i has dimension zero. The dimension formulas for \mathcal{M}_m^0 and \mathcal{M}_m^i (i = 1, 2) then show that $\chi_0 = 2\ell(m)$ and $d - \ell(m)$ is even because $d - \ell(m^i) = \sum (m_j^i - 1)$ is even for all $1 \le i \le k$.

Let |A| denote the cardinality of a set A.

Lemma 2.2.
$$|\mathcal{M}_m^0| = \frac{d! |\operatorname{Aut}(m)|}{\prod m_j}$$

Proof. Let $f \in \mathcal{M}_m^0$. Since $\chi_0 = 2\ell(m)$, the domain of f is a disjoint union of smooth rational curves E_j for $1 \le j \le \ell(m)$, and each restriction $f_j = f|_{E_j}$ has exactly one contact marked point over p^i (i = 1, 2) with multiplicity m_j , so f_j has degree m_j .

Consequently, forgetting contact marked points of maps in \mathcal{M}_m^0 gives exactly one map (as a cover) with automorphism group of order $|\operatorname{Aut}(m)| \prod m_j$. Here the factor $|\operatorname{Aut}(m)|$ appears because we can relabel maps f_j in $|\operatorname{Aut}(m)|$ ways and the factor $\prod m_j$ appears because each restriction map f_j (as a cover) has an automorphism group of order m_j . We then argue as in the proof of Lemma 1.1.

JUNHO LEE

For each $(f_1, f_0, f_2) \in \mathcal{P}_m$, by identifying contact marked points over $p^i \in D^i \cap E$ (i = 1, 2), one can glue the domains of f_i and f_0 to obtain a map $f : C \to D_0$ with $\chi(C) = \chi$. For notational convenience, we often write the glued map f as $f = (f_1, f_0, f_2)$. Denote by

$$(2-2) \mathcal{M}_{m,0}$$

the space of such glued maps $f = (f_1, f_0, f_2)$. Contact marked points are labeled, but nodal points of C are not labeled. Thus, we have the following.

Lemma 2.3. \mathcal{P}_m is a cover of $\mathcal{M}_{m,0}$ of degree $|\operatorname{Aut}(m)|^2$.

3. Limiting and gluing

Following [Ionel and Parker 2004], this section describes limiting and gluing arguments under a degeneration of target curves. Let $D_0 = D_1 \cup E \cup D_2$ be the nodal curve with fixed points q^1, \ldots, q^{k+3} as in Section 2. In Section 4, we construct a family of curves together with k + 3 sections:

Here the total space \mathfrak{D} is a smooth complex surface, $\Delta \subset \mathbb{C}$ is a disk with parameter λ , the central fiber is D_0 , the general fiber D_{λ} ($\lambda \neq 0$) is a smooth curve of genus h, and $Q^i(0) = q^i$ for $1 \le i \le k+3$. By Gromov's convergence theorem, a sequence of holomorphic maps into D_{λ} with $\lambda \to 0$ has a map into D_0 as a limit. For notational simplicity, for $\lambda \neq 0$ we set

(3-2)
$$\mathcal{M}_{\lambda} = \mathcal{M}_{\chi, m^1, \dots, m^{k+3}}^{V_{\lambda}}(D_{\lambda}, d), \text{ where } V_{\lambda} = \{Q^1(\lambda), \dots, Q^{k+3}(\lambda)\},$$

and denote the set of limits of sequences of maps in \mathcal{M}_{λ} as $\lambda \to 0$ by

$$\lim_{\lambda \to 0} \mathcal{M}_{\lambda}.$$

Lemma 3.1 shows that limit maps in (3-3) lie in the union of spaces (2-2), namely,

(3-4)
$$\lim_{\lambda \to 0} \mathcal{M}_{\lambda} \subset \bigcup_{m} \mathcal{M}_{m,0}$$

where the union is over all partitions *m* of *d* with $d - \ell(m)$ even.

Conversely, by the gluing theorem of [Ionel and Parker 2004], the domain of each map in $\mathcal{M}_{m,0}$ can be smoothed to produce maps in \mathcal{M}_{λ} for small $|\lambda|$. Shrinking Δ if necessary, for $\lambda \in \Delta$, one can assign to each $f_{\lambda} \in \mathcal{M}_{\lambda}$ a partition *m* of *d* by (3-4). Let $\mathcal{M}_{m,\lambda}$ be the set of all pairs (f_{λ}, m) . For each $f \in \mathcal{M}_{m,0}$, let

$$(3-5) \qquad \qquad \mathscr{Z}_{m,f} \to \Delta$$

be the connected component of $\bigcup_{\lambda \in \Delta} \mathcal{M}_{m,\lambda} \to \Delta$ that contains f, and let

$$(3-6) \qquad \qquad \mathfrak{L}_{m,f,\lambda}$$

denote the fiber of (3-5) over $\lambda \in \Delta$. It follows that, for $\lambda \neq 0$,

(3-7)
$$\mathcal{M}_{\lambda} = \bigsqcup_{f \in \mathcal{M}_{m,0}} \mathscr{X}_{m,f,\lambda}.$$

For $f = (f_1, f_0, f_2) \in \mathcal{M}_{m,0}$ where $m = (m_1, \dots, m_\ell)$, let y_j^i be the node mapped to p^i at which f_i and f_0 have multiplicity m_j . The gluing theorem shows that one can smooth each node y_j^i in m_j ways to produce $(\prod m_j)^2$ maps in $\mathcal{Z}_{m,f,\lambda}$, so

(3-8)
$$|\mathscr{Z}_{m,f,\lambda}| = \left(\prod m_j\right)^2 \quad (\lambda \neq 0).$$

In order to prove (3-4), we use the following fact on stable maps. An irreducible component of a stable holomorphic map f is a ghost component if its image is a point. Write the domain of f as $C^g \cup C$ where C^g is a connected curve whose irreducible components are all ghost components. Then the stability of f implies that

$$\chi(C^g) - \ell^g - n \le -1$$

where $\ell^g = |C^g \cap C|$ and *n* is the number of marked points on C^g .

Lemma 3.1. Let \mathcal{M}_r and $\mathcal{M}_{m,0}$ be as above. Then we have

$$\lim_{\lambda\to 0}\mathcal{M}_{\lambda}\subset \bigcup_m\mathcal{M}_{m,0}$$

where the union is over all partitions m of d with $d - \ell(m)$ even.

Proof. Let f be a limit map in (3-3). The domain C of f can be written as

(3-10)
$$C = C_1 \cup C_0 \cup C_2 \cup \left(\bigcup_{i=1}^{k+3} C_i^g\right) \cup C^g \cup \widetilde{C}^g$$

where C_0 maps to E, C_1 and C_2 map to D_1 and D_2 , C_i^g is the union of all ghost components over q^i , where i = 1, ..., k+3, C^g is the union of all ghost components over points in $D_0 \setminus (V_1 \cup V_0 \cup V_2)$, and \widetilde{C}^g is the union of all ghost components over $\{p^1, p^2\}$. Let $f_j = f|_{C_j}$ for j = 0, 1, 2. Observe that f_j is V_j -regular because C_j has no ghost components. Let \widehat{m}^i be a contact vector over q^i , \widetilde{m}^1 and \widetilde{m}^2 be contact vectors of f_1 and f_2 over p^1 and p^2 , and $\widetilde{m}^{0;1}$ and $\widetilde{m}^{0;2}$ be contact vectors of f_0 over p^1 and p^2 . The Riemann–Hurwitz formulas for f_0 , f_1 , and f_2 give

(3-11)
$$\sum_{j=0}^{2} \chi(C_j) \le 2d(1-h) + \sum_{i=1}^{k+3} (\ell(\widehat{m}^i) - d) + \sum_{i=1}^{2} (\ell(\widetilde{m}^i) + \ell(\widetilde{m}^{0;i})).$$

JUNHO LEE

For i = 1, ..., k + 3, let $\ell_i = |C_1 \cup C_0 \cup C_2 \cap C_i^g|$ and let n_i be the number of marked points on C_i^g . Since all marked points are limits of marked points, we have

(3-12)
$$\ell(\widehat{m}^i) = \ell(m^i) - n_i + \ell_i.$$

For j = 0, 1, 2, let $\tilde{\ell}_j = |C_j \cap \tilde{C}^g|$. Counting the number of nodes mapped to p^1 and p^2 shows

(3-13)
$$\sum_{i=1}^{2} (\ell(\widetilde{m}^{i}) - \widetilde{\ell}_{i}) = \sum_{i=1}^{2} |C_{i} \cap C_{0}| = \sum_{i=1}^{2} \ell(\widetilde{m}^{0;i}) - \widetilde{\ell}_{0}.$$

Let $\ell^{g} = |C_{1} \cup C_{0} \cup C_{2} \cap C^{g}|$. Since $\chi(C) = \chi$, by (3-10) and (3-13) we have

(3-14)
$$\chi = \sum_{j=0}^{2} \chi(C_j) + \sum_{i=1}^{k+3} (\chi(C_i^g) - 2\ell_i) + \chi(C^g) - 2\ell^g + \chi(\widetilde{C}^g) - \widetilde{\ell} - \sum_{i=1}^{2} (\ell(\widetilde{m}^i) + \ell(\widetilde{m}^{0;i})),$$

where $\tilde{\ell} = \tilde{\ell}_0 + \tilde{\ell}_1 + \tilde{\ell}_2$. By our assumption that formula (0-1) holds, it follows from (3-11), (3-12), and (3-14) that

(3-15)
$$\chi \leq \chi + \sum_{i=1}^{k+3} (\chi(C_i^g) - \ell_i - n_i) + \chi(C^g) - 2\ell^g + \chi(\widetilde{C}^g) - \widetilde{\ell}.$$

Noting that C^g and \widetilde{C}^g have no marked points, by (3-9) and (3-15), we conclude that the domain C of f has no ghost components. Consequently,

- f_i is V_i -regular for j = 0, 1, 2,
- $\widetilde{m}^i = \widetilde{m}^{0;i}$ for i = 1, 2 (see Lemma 3.3 of [Ionel and Parker 2004]) and $\widehat{m}^i = m^i$ for i = 1, ..., k + 3.

In particular, the equality in (3-11) holds; otherwise we have a strict inequality in (3-15). So, we have $\chi(C_0) = \ell(\widetilde{m}^1) + \ell(\widetilde{m}^2)$. But $\chi(C_0) \le 2 \min\{\ell(\widetilde{m}^1), \ell(\widetilde{m}^2)\}$. It follows that

- C_0 has $\ell(\widetilde{m}^1) = \ell(\widetilde{m}^2)$ connected components E_j with $\chi(E_j) = 2$ for all j,
- $\widetilde{m}_{i}^{1} = \deg(f_{0}|_{E_{i}}) = \widetilde{m}_{i}^{2}$ for all *j*, that is, $\widetilde{m}^{1} = \widetilde{m}^{2}$.

It follows that the Euler characteristics of C_0 , C_1 , and C_2 satisfy (2-1) by (3-14). Therefore, $f \in \mathcal{M}_{m,0}$ for $m = \tilde{m}^1 = \tilde{m}^2$ and $d - \ell(m)$ is even by Lemma 2.1.

4. Smooth model by Schiffer variation

A *Schiffer variation* of a nodal curve (compare [Arbarello et al. 2011, p. 184]) is obtained by gluing deformations $uv = \lambda$ near nodes with the trivial deformation

away from nodes. In this section, we use the method of Schiffer variation to construct a smooth model for the space $\mathscr{Z}_{m,f}$ in (3-5), which has several branches intersecting at f unless m is trivial.

In this section, we fix an odd partition $m = (n^{\ell})$, that is, $m = (m_1, \ldots, m_{\ell})$ with

(4-1)
$$m_1 = \cdots = m_\ell = n$$
, where $n = d/\ell$ is odd.

Let $f = (f_1, f_0, f_2)$ be a map in $\mathcal{M}_{m,0}$ in (2-2). As described in Section 2, the central fiber of $\rho : \mathfrak{D} \to \Delta$ is the nodal curve $D_0 = D_1 \cup E \cup D_2$ with two nodes $p^1 \in D_1 \cap E$ and $p^2 \in D_2 \cap E$ where $E = \mathbb{P}^1$. The domain of f is a nodal curve

$$C = C_1 \cup C_0 \cup C_2$$
, where $C_0 = \bigcup_{j=1}^{\ell} E_{\ell}$,

with 2ℓ nodes, such that, for i = 1, 2 and $j = 1, \dots, \ell$,

- $f^{-1}(p^i)$ consists of the ℓ nodes $y_i^i \in C_i \cap E_j$,
- C_i is smooth and $f|_{C_i} = f_i$ has ramification index $m_j = n$ at the node y_i^i ,
- $E_j = \mathbb{P}^1$ and $f|_{E_j} = f_0|_{E_j} : E_j \to E$ has ramification index $m_j = n$ at the node y_i^i .

The following is the main result of this section.

Proposition 4.1. Let f be as above. Then, for each vector $\zeta = (\zeta_1^1, \zeta_1^2, \dots, \zeta_{\ell}^1, \zeta_{\ell}^2)$, where ζ_j^i is an n^{th} root of unity, there are a family of curves $\varphi_{\zeta} : \mathscr{C}_{\zeta} \to \Delta$, with smooth total space \mathscr{C}_{ζ} , over a disk Δ (with parameter s) and a holomorphic map $\mathscr{F}_{\zeta} : \mathscr{C}_{\zeta} \to \mathfrak{D}$ satisfying:

- (a) the central fiber $C_{\zeta,0} = C$ and the restriction map $\mathscr{F}_{\zeta}|_{C} = f$;
- (b) the general fiber $C_{\zeta,s}$ ($s \neq 0$) is smooth and, for $\lambda = s^n \neq 0$,

(4-2)
$$\bigcup_{\zeta} \{f_{\zeta,s}\} = \mathscr{X}_{m,f,\lambda}.$$

where the union is over all ζ , $f_{\zeta,s} = \mathscr{F}_{\zeta}|_{C_{\zeta,s}}$ and $\mathscr{X}_{m,f,\lambda}$ is the space (3-6).

Proof. The proof consists of four steps.

Step 1. We first show how to construct the family of curves $\rho : \mathfrak{D} \to \Delta$ with k+3 sections. For i = 1, 2, a neighborhood of the node $p^i \in D_i \cap E$ can be regarded as the union $U^i \cup V^i$ of the two disks

$$U^{i} = \{u^{i} \in \mathbb{C} : |u^{i}| < 1\} \subset D_{i} \text{ and } V^{i} = \{v^{i} \in \mathbb{C} : |v^{i}| < 1\} \subset E_{i}$$

with their origins identified. We may assume that the fixed points q^1, \ldots, q^{k+3} in D_0 described in (2-1) lie outside these sets. Consider the regions

JUNHO LEE

$$A^{i} = \left\{ (u^{i}, v^{i}, \lambda) \in U^{i} \times V^{i} \times \Delta : u^{i} v^{i} = \lambda \right\},\$$
$$B = \bigcup_{i=1}^{2} G^{i} \cup \left[\left(D_{0} \setminus \bigcup_{i=1}^{2} (U^{i} \cup V^{i}) \right) \times \Delta \right],\$$

where

$$G^{i} = \left\{ (u^{i}, \lambda) \in U^{i} \times \Delta : |u^{i}| > \sqrt{|\lambda|} \right\} \cup \left\{ (v^{i}, \lambda) \in V^{i} \times \Delta : |v^{i}| > \sqrt{|\lambda|} \right\}.$$

We obtain a smooth complex surface \mathfrak{D} by gluing A^1 , A^2 , and B_0 using the maps

(4-3)
$$G^i \to A^i$$
 defined by $(u^i, \lambda) \to \left(u^i, \frac{\lambda}{u^i}, \lambda\right)$ and $(v^i, \lambda) \to \left(\frac{\lambda}{v^i}, v^i, \lambda\right)$.

Let $\rho : \mathfrak{D} \to \Delta$ be the projection to the last factor and define k + 3 sections Q^i of ρ by

$$Q^i(\lambda) = (q^i, \lambda).$$

Step 2. We can similarly construct a family of curves over a 2ℓ -dimensional polydisk:

(4-4)
$$\varphi_{2\ell} : \mathscr{X} \to \Delta_{2\ell} = \{ t = (t_1^1, t_1^2, \dots, t_\ell^1, t_\ell^2) \in \mathbb{C}^{2\ell} : |t_j^i| < 1 \}.$$

For each node $y_i^i \in C_i \cap E_j$, choose a neighborhood obtained from two disks

$$U_{j}^{i} = \{u_{j}^{i} \in \mathbb{C} : |u_{j}^{i}| < 1\} \subset C_{i} \text{ and } V_{j}^{i} = \{v_{j}^{i} \in \mathbb{C} : |v_{j}^{i}| < 1\} \subset E_{j}$$

by identifying the origins. Consider the regions

$$A_j^i = \left\{ (u_j^i, u_j^i, t) \in U_j^i \times V_j^i \times \Delta_{2\ell} : u_j^i v_j^i = t_j^i \right\},\$$

$$B_{2\ell} = \bigcup_{i,j} G_j^i \cup \left[\left(C \setminus \bigcup_{i,j} (U_j^i \cup V_j^i) \right) \times \Delta_{2\ell} \right],\$$

where

$$G_{j}^{i} = \left\{ (u_{j}^{i}, t) \in U_{j}^{i} \times \Delta_{2\ell} : |u_{j}^{i}| > \sqrt{|t_{j}^{i}|} \right\} \cup \left\{ (v_{j}^{i}, t) \in V_{j}^{i} \times \Delta_{2\ell} : |v_{j}^{i}| > \sqrt{|t_{j}^{i}|} \right\}.$$

We can then obtain a smooth complex manifold \mathscr{X} of dimension $2\ell + 1$ by gluing $\bigcup A_i^i$ and $B_{2\ell}$ with the maps

(4-5)
$$G_j^i \to A_j^i$$
 defined by $(u_j^i, t) \to \left(u_j^i, \frac{t_j^i}{u_j^i}, t\right)$ and $(v_j^i, t) \to \left(\frac{t_j^i}{v_j^i}, v_j^i, t\right)$.

Let $\varphi_{2\ell} : \mathscr{X} \to \Delta$ be the projection to the factor *t*.

Step 3. Since f_i and $f_0|_{E_j}$ have ramification index $m_j = n$ at y_j^i , we may assume (after coordinates change) that on U_j^i and V_j^i the map f can be written as

(4-6)
$$U_j^i \to U^i \text{ by } u_j^i \to (u_j^i)^n \text{ and } V_j^i \to V^i \text{ by } v_j^i \to (v_j^i)^n.$$

For each i, j, define a map

(4-7)
$$G_j^i \to G^i$$
 by $(u_j^i, t) \to ((u_j^i)^n, (t_j^i)^n)$ and $(u_j^i, t) \to ((v_j^i)^n, (t_j^i)^n)$.

On the other hand, for each i, j, we have a map

(4-8)
$$A_j^i \to A^i$$
 defined by $(u_j^i, v_j^i, t) \to ((u_j^i)^n, (v_j^i)^n, (t_j^i)^n).$

These two maps (4-7) and (4-8) are glued together under the maps (4-3) and (4-5). The glued map extends to a holomorphic map $f_t : \mathscr{X}_t \to D_\lambda$ if and only if

(4-9)
$$(t_1^1)^n = (t_1^2)^n = \dots = (t_\ell^1)^n = (t_\ell^2)^n = \lambda.$$

There are $n^{2\ell}$ solutions t of (4-9) and the extension map f_t is given by

$$(x,t) \to (f(x),\lambda) \text{ on } \mathscr{X}_t - \bigcup A_j^i.$$

Step 4. For each vector $\zeta = (\zeta_1^1, \zeta_1^2, \dots, \zeta_\ell^1, \zeta_\ell^2)$, where each ζ_j^i is an n^{th} root of unity, define

 $\delta_{\zeta}: \Delta \to \Delta_{2\ell}$ by $s \to (\zeta_1^1 s, \zeta_1^2 s, \zeta_2^1 s, \zeta_2^2 s, \ldots, \zeta_\ell^1 s, \zeta_\ell^2 s).$

The pullback $\delta_{\epsilon}^* \mathscr{X}$ gives a family of curves:

The central fiber is $C_{\zeta,0} = C$ and the general fiber $C_{\zeta,s}$ ($s \neq 0$) is smooth. A neighborhood of the node y_i^i of C in \mathscr{C}_{ζ} can be viewed as

(4-11)
$$\hat{A}_{j}^{i} = \left\{ (u_{j}^{i}, v_{j}^{i}, s) \in \mathbb{C}^{3} : |u_{j}^{i}| < 1, |v_{j}^{i}| < 1, u_{j}^{i}v_{j}^{i} = \zeta_{j}^{i}s \right\}.$$

It follows that the total space \mathscr{C}_{ζ} is a complex smooth surface. Noting $\delta_{\zeta}(s)$ is a solution of (4-9) for $\lambda = s^n$, we obtain a holomorphic map $\mathscr{F}_{\zeta} : \mathscr{C}_{\zeta} \to \mathfrak{D}$ given by

(4-12)
$$\begin{array}{c} (u_j^i, v_j^i, s) \to ((u_j^i)^n, (v_j^i)^n, s^n) \quad \text{on } \hat{A}_j^i, \\ (x, s) \to (f(x), s^n) \qquad \text{on } \mathscr{C}_{\zeta} - \bigcup \hat{A}_j^i. \end{array}$$

Since the restriction $\mathcal{F}_{\zeta}|_{C} = f$ by (4-6) and (4-12), it remains to show (4-2). By our choice of fixed points q^{i} on D_{0} , each contact marked point x_{j}^{i} of f lies in $\mathscr{C}_{\zeta} - \bigcup \hat{A}_{j}^{i}$. Thus, by (4-12), the pullback $\mathcal{F}_{\zeta}^{*}Q^{i}$ of the section Q^{i} of ρ gives a section X_{j}^{i} of φ_{ζ} given by $X_{j}^{i}(s) = (x_{j}^{i}, s)$. After marking the points $X_{j}^{i}(s)$ in $C_{\zeta,s}$, the restriction map

$$f_{\zeta,s} = \mathscr{F}_{\zeta}|_{C_{\zeta,s}} : C_{\zeta,s} \to D_{\lambda}, \text{ where } \lambda = s^n \neq 0,$$

JUNHO LEE

has contact marked points $X_j^i(s)$ over $Q^i(\lambda)$ with multiplicity m_j^i . This means that $f_{\zeta,s}$ lies in the space \mathcal{M}_{λ} in (3-2) for $\lambda = s^n$. Therefore, noting that $f_{\zeta,s} \to f$ as $s \to 0$ and that $|\mathscr{X}_{m,f,\lambda}| = n^{2\ell}$ by (3-8), we conclude (4-2).

5. Spin structure and parity

The aim of this section is to use a spin structure on a family of nodal curves [Cornalba 1989] to show the parity calculation in Proposition 5.4. Twisting a bundle as in (5-6) is a key idea for parity calculation.

We first introduce a spin structure on a family of nodal curves that is relevant to our discussion. We refer to [Cornalba 1989] for the definition of spin structure and more details. The relative dualizing sheaf ω_{ρ} of the family of curves $\rho : \mathfrak{D} \to \Delta$ in (3-1) is the canonical bundle $K_{\mathfrak{D}}$ on the total space \mathfrak{D} , since \mathfrak{D} is smooth and K_{Δ} is trivial. For each $\lambda \neq 0$, the restriction $K_{\mathfrak{D}}|_{D_{\lambda}}$ is the canonical bundle $K_{D_{\lambda}}$ on D_{λ} , and the restriction $K_{\mathfrak{D}}|_{D_0}$ is the dualizing sheaf ω_{D_0} of the nodal curve $D_0 = D_1 \cup E \cup D_2$. As described in Section 4, D_0 is locally given by $u^i v^i = 0$ near each node p^i in $D_i \cap E$ for i = 1, 2. Then the local generators of ω_{D_0} are du^i/u^i and dv^i/v^i with a relation $du^i/u^i + dv^i/v^i = 0$; see [Harris and Morrison 1998, p. 82]. This implies the restriction $\omega_{D_0}|_{D_i} = K_{D_i} \otimes \mathbb{O}(p^i)$. On the other hand, $1/u^i$ is a local defining function for the divisor -E on \mathfrak{D} near p^i . By restricting $1/u^i$ to D_i , one can see that $\mathbb{O}(-E)|_{D_i} = \mathbb{O}(-p^i)$. Consequently, for i = 1, 2,

(5-1)
$$K_{\mathfrak{D}}|_{D_i} \otimes \mathbb{O}(-E)|_{D_i} = \omega_{D_0}|_{D_i} \otimes \mathbb{O}(-p^1) = K_{D_i}.$$

From Cornalba's construction [1989, p. 570], there are a line bundle $\mathcal{N} \to \mathfrak{D}$ and a homomorphism $\Phi : \mathcal{N}^2 \to \omega_\rho = K_{\mathfrak{D}}$ satisfying the following.

- Φ vanishes identically on the exceptional component *E* and $\mathcal{N}|_E = \mathbb{O}_E(1)$.
- Since $\Phi|_E \equiv 0$, there is an induced homomorphism $\hat{\Phi} : \mathcal{N}^2 \to K_{\mathfrak{D}} \otimes \mathbb{O}(-E)$ such that Φ is the composition of $\hat{\Phi}$ with tensoring with η :

(5-2)
$$\Phi: \mathcal{N}^2 \xrightarrow{\hat{\Phi}} K_{\mathfrak{D}} \otimes \mathbb{O}(-E) \xrightarrow{\otimes \eta} K_{\mathfrak{D}},$$

where η is a section of $\mathbb{O}(E)$ with zero divisor *E*. Then, for i = 1, 2, the restriction

$$\hat{\Phi}|_{D_i} : (\mathcal{N}|_{D_i})^2 \to K_{\mathfrak{D}}|_{D_i} \otimes \mathbb{O}(-E)|_{D_i} = K_{D_i}$$

is an isomorphism so that the restriction $N_i = \mathcal{N}|_{D_i}$ is a theta characteristic on D_i .

• For each $\lambda \neq 0$, the restriction $\Phi|_{D_{\lambda}} : (\mathcal{N}|_{D_{\lambda}})^2 \to K_{D_{\lambda}}$ is an isomorphism so that the restriction $N_{\lambda} = \mathcal{N}|_{D_{\lambda}}$ is a theta characteristic on D_{λ} .

The pair (\mathcal{N}, Φ) is a spin structure on $\rho : \mathfrak{D} \to \Delta$ and the restriction $\mathcal{N}|_{D_0}$ is a theta characteristic on the nodal curve D_0 .

Remark 5.1. Atiyah [1971] and Mumford [1971] showed that the parity of a theta characteristic on a smooth curve is a deformation invariant. Cornalba [1989, Page 580] used the homomorphism Φ to extend Mumford's proof to the case of spin structure on a family of nodal curves. Thus, if p_1 , p_2 , and p are the parities of N_1 , N_2 , and N_λ ($\lambda \neq 0$), we have

$$p \equiv p_1 + p_2 \pmod{2}.$$

Let $\varphi_{\zeta} : \mathscr{C}_{\zeta} \to \Delta$ be the family of curves in Proposition 4.1. Recall that the central fiber of φ_{ζ} is $C = C_1 \cup C_0 \cup C_2$, where $C_0 = \bigsqcup_j E_j$ is a disjoint union of ℓ exceptional components E_j and $C_i \cap E_j = \{y_j^i\}$ for i = 1, 2 and $1 \le j \le \ell$. Similarly as for (5-1), by restricting local defining functions, we have

(5-3)
$$\mathbb{O}(\pm C_0)|_{C_i} = \mathbb{O}\left(\pm \sum_j y_j^i\right)$$
 $(i = 1, 2)$ and $\mathbb{O}(\pm C_0)|_{C_{\xi,s}} = \mathbb{O}$ $(s \neq 0).$

Since any fiber of φ_{ζ} is a principal divisor on \mathscr{C}_{ζ} , $\mathbb{O}(C) = \mathbb{O}$ and hence $\mathbb{O}(C_0) = \mathbb{O}(-C_1 - C_2)$. We also have

(5-4)
$$\mathbb{O}(\pm C_0)|_{E_j} = \mathbb{O}(\mp (C_1 + C_2))|_{E_j} = \mathbb{O}(\mp (y_j^1 + y_j^2)) = \mathbb{O}(\mp 2)(1 \le j \le \ell).$$

Let $f = (f_1, f_0, f_2)$ and $\mathscr{F}_{\zeta} : \mathscr{C}_{\zeta} \to \mathfrak{D}$ be the maps in Proposition 4.1. The ramification divisor $R_{\mathscr{F}_{\zeta}}$ of \mathscr{F}_{ζ} has local defining functions given by the Jacobian of \mathscr{F}_{ζ} , so (4-12) shows

(5-5)
$$R_{\mathcal{F}_{\zeta}} = \mathbb{O}(X_{\zeta} + (n-1)C) = \mathbb{O}(X_{\zeta}),$$

where $X_{\zeta} = \sum_{i,j} (m_j^i - 1) X_j^i$ and X_j^i is the section of φ_{ζ} defined in (4-12). Note that

- (i) the ramification divisor of $f_i = \mathcal{F}_{\zeta}|_{C_i}$ (i = 1, 2) is $R_{f_i} = X_{\zeta}|_{C_i} + \sum_i (n-1)y_i^i$;
- (ii) the ramification divisor of $f_{\zeta,s} = \mathcal{F}_{\zeta}|_{C_{\zeta,s}}$ $(s \neq 0)$ is $R_{f_{\zeta,s}} = X_{\zeta}|_{C_{\zeta,s}}$.

Now, noting *n* is odd, we twist the pullback bundle $\mathcal{F}_{\zeta}^* \mathcal{N}$ by setting

(5-6)
$$\mathscr{L}_{\zeta} = \mathscr{F}_{\zeta}^* \mathscr{N} \otimes \mathbb{O}\left(\frac{1}{2}X_{\zeta} + \frac{(n-1)}{2}C_0\right).$$

The lemma below shows that the twisted line \mathscr{L}_{ζ} restricts to a theta characteristic on each fiber of φ_{ζ} , including the central fiber *C*.

Lemma 5.2. Let \mathcal{L}_{ζ} be as above. Then:

- (a) $\mathscr{L}_{\zeta}|_{E_j} = \mathbb{O}(1)$ for $1 \le j \le \ell$.
- (b) $\mathscr{L}_{\zeta}|_{C_1} = L_{f_1}, \mathscr{L}_{\zeta}|_{C_2} = L_{f_2} and \mathscr{L}_{\zeta}|_{C_{\zeta,s}} = L_{f_{\zeta,s}} for s \neq 0$, where $L_{f_1}, L_{f_2}, L_{f_{\zeta,s}}$ are the theta characteristics on $C_1, C_2, C_{\zeta,s}$ defined by (0-2).

Proof. Part (a) follows from (5-4) and the fact that each restriction map $\mathscr{F}_{\zeta}|_{E_j}$ has degree *n*. Part (b) follows from (5-3), (i), and (ii).

Observe that the relative dualizing sheaf $\omega_{\varphi_{\zeta}}$ is the canonical bundle $K_{\mathscr{C}_{\zeta}}$ since \mathscr{C}_{ζ} is smooth. The Hurwitz formula and (5-5) thus imply that

(5-7)
$$\omega_{\varphi_{\zeta}} = K_{\mathscr{C}_{\zeta}} = \mathscr{F}_{\zeta}^* K_{\mathfrak{D}} \otimes \mathbb{O}(X_{\zeta}).$$

Define a homomorphism

(5-8)
$$\hat{\Psi}_{\zeta} : \mathscr{L}_{\zeta}^{2} = \mathscr{F}_{\zeta}^{*} \mathscr{N}^{2} \otimes \mathbb{O}(X_{\zeta} + (n-1)C_{0}) \rightarrow \mathscr{F}_{\zeta}^{*}(K_{\mathfrak{D}} \otimes \mathbb{O}(-E)) \otimes \mathbb{O}(X_{\zeta} + (n-1)C_{0})$$

by $\hat{\Psi}_{\zeta} = \mathscr{F}_{\zeta}^* \hat{\Phi} \otimes \text{Id}$, where $\hat{\Phi}$ is the induced homomorphism in (5-2). Noting that $\mathbb{O}(C) = \mathbb{O}$ and $\mathbb{O}(D_0) = \mathbb{O}$, by (4-12), we have

$$\mathcal{F}^*_{\zeta}\mathbb{O}(-E) = \mathcal{F}^*_{\zeta}\mathbb{O}(D_1 + D_2) = \mathbb{O}(n(C_1 + C_2)) = \mathbb{O}(-nC_0).$$

Together with (5-7), this implies that the right side of (5-8) is $K_{\mathscr{C}_{\zeta}} \otimes \mathbb{O}(-C_0)$. Now define a homomorphism $\Psi_{\zeta} : \mathscr{L}^2_{\zeta} \to K_{\mathscr{C}_{\zeta}}$ to be the composition

(5-9)
$$\Psi_{\zeta}: \mathscr{L}^2_{\zeta} \xrightarrow{\hat{\Psi}_{\zeta}} K_{\mathscr{C}_{\zeta}} \otimes \mathbb{O}(-C_0) \xrightarrow{\otimes \xi} K_{\mathscr{C}_{\zeta}}.$$

where ξ is a section of $\mathbb{O}(C_0)$ with zero divisor C_0 .

Lemma 5.3. $(\mathscr{L}_{\zeta}, \Psi_{\zeta})$ is a spin structure on $\varphi_{\zeta} : \mathscr{C}_{\zeta} \to \Delta$.

Proof. First, $\mathscr{L}_{\zeta}|_{E} = \mathbb{O}(1)$ by Lemma 5.2(a) and Ψ_{ζ} vanishes identically on each exceptional component E_{j} , since $\xi = 0$ on $C_{0} = \bigsqcup_{j} E_{j}$. Second, since $\hat{\Phi}|_{D_{i}}$ is an isomorphism, (5-3) and (i) show that, for i = 1, 2, the restriction

$$\hat{\Psi}|_{C_i} = f_i^*(\hat{\Phi}|_{D_i}) \otimes \mathrm{Id} : (\mathcal{L}_{\zeta}|_{C_i})^2 = f_i^* N_i^2 \otimes \mathbb{O}(R_{f_i}) \to f_i^* K_{D_i} \otimes \mathbb{O}(R_{f_i}) = K_{C_i}$$

is an isomorphism. Lastly, let $\lambda = s^n \neq 0$. Since $\Phi|_{D_{\lambda}}$ is an isomorphism, so is $\hat{\Phi}|_{D_{\lambda}}$. Thus, by (5-3), (ii), and the facts $K_{\mathfrak{D}}|_{D_{\lambda}} = K_{D_{\lambda}}$ and $\mathbb{O}(-E)|_{D_{\lambda}} = \mathbb{O}$, the restriction

$$\hat{\Psi}_{\zeta}|_{C_{\zeta,s}} = f_{\zeta,s}^* \hat{\Phi}|_{D_{\lambda}} \otimes \mathrm{Id} : (\mathscr{L}_{\zeta}|_{C_{\zeta,s}})^2 = f_{\zeta,s}^* N_{\lambda}^2 \otimes \mathbb{O}(R_{f_{\zeta,s}}) \to f_{\zeta,s}^* K_{D_{\lambda}} \otimes \mathbb{O}(R_{f_{\zeta,s}}) = K_{C_{\zeta,s}}$$

is an isomorphism. This implies that the restriction

$$\Psi_{\zeta}|_{C_{\zeta,s}}: (\mathscr{L}_{\zeta}|_{C_{\zeta,s}})^2 \to K_{C_{\zeta}}|_{C_{\zeta,s}} = K_{C_{\zeta},s}$$

is also an isomorphism. Therefore, we conclude that $(\mathscr{L}_{\zeta}, \Psi_{\zeta})$ is a spin structure on φ_{ζ} .

The following is a key fact for the proof of Theorem 0.1.

Proposition 5.4. Let $f = (f_1, f_0, f_2)$ and $f_{\zeta,s}$ be maps in Proposition 4.1. Then, for all $s \neq 0$,

(5-10)
$$p(f_{\zeta,s}) \equiv p(f_1) + p(f_2) \pmod{2}.$$

Proof. Since $(\mathscr{L}_{\zeta}, \Psi_{\zeta})$ is a spin structure on φ_{ζ} , Cornalba's proof, mentioned in Remark 5.1, shows that, for all $s \neq 0$,

$$h^{0}(\mathscr{L}_{\zeta}|_{C_{\zeta,s}}) \equiv h^{0}(\mathscr{L}_{\zeta}|_{C_{1}}) + h^{0}(\mathscr{L}_{\zeta}|_{C_{2}}) \pmod{2}.$$

This and Lemma 5.2(b) prove (5-10).

6. Proof of Theorem 0.1

Proof. Fix a spin structure (\mathcal{N}, Φ) on $\rho : \mathfrak{D} \to \Delta$ given in Section 5. Consider the space $\mathcal{M}_{m,0}$ in (2-2) where *m* is a partition of d = 3. In this case, by Lemma 2.1, either $m = (1^3)$ or m = (3). Note that both of them satisfy (4-1). Fix $\lambda \neq 0$ and let $f = (f_1, f_0, f_2)$ be a map in $\mathcal{M}_{m,0}$. Then (4-2) and (5-10) show that, for all $f_{\mu} \in \mathfrak{L}_{m, f, \lambda},$

(6-1)
$$p(f_{\mu}) \equiv p(f_1) + p(f_2) \pmod{2}$$

Lemma 1.1 and (3-7) show that

(6-2)
$$H_{(3)^{k}}^{h,p} = H_{(3)^{k},(1^{3})^{3}}^{h,p} = \frac{1}{(3!)^{3}} \left(\sum_{f \in \mathcal{M}_{(1^{3}),0}} \sum_{f_{\mu} \in \mathcal{Z}_{(1^{3}),f,\lambda}} (-1)^{p(f_{\mu})} + \sum_{f \in \mathcal{M}_{(3),0}} \sum_{f_{\mu} \in \mathcal{Z}_{(3),f,\lambda}} (-1)^{p(f_{\mu})} \right).$$

By (3-8) and (6-1), (6-2) becomes

(6-3)
$$H_{(3)^k}^{h,p} = \sum_{\substack{f = (f_1, f_0, f_2) \in \mathcal{M}_{(1^3, 0)}}} \frac{(-1)^{p(f_1) + p(f_2)}}{(3!)^3} + \sum_{\substack{f = (f_1, f_0, f_2) \in \mathcal{M}_{(3), 0}}} \frac{3^2(-1)^{p(f_1) + p(f_2)}}{(3!)^3}.$$

It then follows from Lemma 2.3 and (6-3) that

$$\begin{split} H_{(3)^{k}}^{h,p} &= \sum_{(f_{1},f_{0},f_{2})\in\mathscr{P}_{(1^{3})}} \frac{(-1)^{p(f_{1})+p(f_{2})}}{(3!)^{5}} + \sum_{(f_{1},f_{0},f_{2})\in\mathscr{P}_{(3)}} \frac{3^{2}(-1)^{p(f_{1})+p(f_{2})}}{(3!)^{3}} \\ &= \frac{1}{(3!)^{3}} \sum_{f_{1}\in\mathscr{M}_{(1^{3})}^{1}} (-1)^{p(f_{1})} \sum_{f_{2}\in\mathscr{M}_{(1^{3})}^{2}} (-1)^{p(f_{2})} + \frac{3}{(3!)^{2}} \sum_{f_{1}\in\mathscr{M}_{(3)}^{1}} (-1)^{p(f_{1})} \sum_{f_{2}\in\mathscr{M}_{(3)}^{2}} (-1)^{p(f_{2})} \\ &= 3! H_{(3)^{k_{1}}}^{h_{1},p_{1}} \cdot H_{(3)^{k_{2}}}^{h_{2},p_{2}} + 3 H_{(3)^{k_{1}+1}}^{h_{1},p_{1}} \cdot H_{(3)^{k_{2}+1}}^{h_{2},p_{2}}; \end{split}$$

the second equality follows from Lemma 2.2 and the last from Lemma 1.1.

JUNHO LEE

7. Calculation

Proposition 7.1. $H_{(3)^k}^{h,\pm} = 3^{2h-2}[(-1)^k 2^{k+h-1} \pm 1].$

Proof. The proof consists of four steps.

Step 1. We first show the following facts which we use in the computation below.

Lemma 7.2. (a)
$$H_{(3)^0}^{0,+} = H_3^{0,+} = \frac{1}{3!}$$
, (b) $H_{(3)^3}^{0,+} = -\frac{1}{3}$, (c) $H_{(3)^0}^{1,+} = H_3^{1,+} = 2$.

Proof. Consider the dimension-zero space $\mathcal{M}^{V}_{\chi}(\mathbb{P}^{1}, 3)$ where $V = \emptyset$. The Euler characteristic $\chi = 6$ by (0-1), and hence the space contains only one map $f : C \to \mathbb{P}^{1}$ where *C* is a disjoint union of three rational curves and $|\operatorname{Aut}(f)| = 3!$. This shows (a). Let (f, C) be a map in the dimension-zero space $\mathcal{M}^{V}_{\chi,(3),(3),(3)}(\mathbb{P}^{1}, 3)$. Then *C* is a connected curve of genus one and the theta characteristic L_{f} on *C* defined by (0-2) is

$$L_f = \mathbb{O}(-2x_1 + x_2 + x_3) = \mathbb{O}(x_1 - 2x_2 + x_3) = \mathbb{O}(x_1 + x_2 - 2x_3),$$

where x_1, x_2 , and x_3 are ramification points of f. This implies $L_f^3 = \mathbb{O}$, and hence $L_f = \mathbb{O}$ because $L_f^2 = L_f^3 = \mathbb{O}$. We have p(f) = 1. Therefore,

$$H^{0,+}_{(3)^3} = -H^0_{(3)^3} = -\frac{1}{3},$$

where $H_{(3)^3}^0$ denotes the (ordinary) Hurwitz number, which is calculated by using the character formula; see [Okounkov and Pandharipande 2006, (0.10)]. By Proposition 9.2 of [Lee and Parker 2007], the spin Hurwitz numbers $H_d^{h,p}$ are the dimension-zero local invariants of spin curve that count maps from possibly disconnected domains. Let $GW_d^{h,p}$ denote the dimension-zero local invariants of spin curve that count maps from connected domains. Then $H_d^{h,p}$ and $GW_d^{h,p}$ are related as follows:

$$1 + \sum_{d>0} H_d^{h,p} t^d = \exp\bigg(\sum_{d>0} G W_d^{h,p} t^d\bigg).$$

Now (c) follows from $GW_1^{1,+} = 1$, $GW_2^{1,+} = 1/2$, and $GW_3^{1,+} = 4/3$; see Section 10 of [Lee and Parker 2007].

Step 2. In this step, we compute $H_{(3)^k}^{1,-}$. For a spin curve of genus one with trivial theta characteristic. It follows from formula (3.12) of [Eskin et al. 2008] that

(7-1)
$$H_{(3)^k}^{1,-} = 2^{-k} [(f_{(3)}(21))^k - (f_{(3)}(3))^k].$$

Here the *central character* $f_{(3)}$ can be written as

$$f_{(3)} = \frac{1}{3} \boldsymbol{p}_3 + a_2 \boldsymbol{p}_1^2 + a_1 \boldsymbol{p}_1 + a_0$$

for some $a_i \in \mathbb{Q}$ ($0 \le i \le 2$), and the *supersymmetric functions* p_1 and p_3 are defined

by

$$p_1(m) = d - \frac{1}{24}$$
 and $p_3(m) = \sum_j m_j^3 - \frac{1}{240}$

where $m = (m_1, \ldots, m_\ell)$ is a partition of d. For k = 0, 1, (7-1) shows

(7-2)
$$H_{(3)^0}^{1,-} = 0$$
 and $H_{(3)}^{1,-} = -3$.

Lemma 7.2(b), (7-2), and formula (0-6) give $H_{(3)^2}^{1,-} = 3H_{(3)}^{1,-} \cdot H_{(3)^3}^{0,+} = 3$. Together with (7-1) and (7-2), this yields $f_{(3)}(21) = -4$ and $f_{(3)}(3) = 2$. From this and (7-1) we have, for $k \ge 0$,

(7-3)
$$H_{(3)^k}^{1,-} = (-1)^k 2^k - 1$$

Step 3. In this step, we compute $H_{(3)^k}^{h,+}$ for h = 0, 1. For $k \ge 1$, (7-2) and formula (0-6) give $H_{(3)^{k-1}}^{1,-} = 3H_{(3)}^{1,-} \cdot H_{(3)^k}^{0,+} = -3^2 H_{(3)^k}^{0,+}$. Combining this with Lemma 7.2(a) we obtain, for $k \ge 0$,

(7-4)
$$H_{(3)^k}^{0,+} = -\frac{1}{3^2}((-1)^{k-1}2^{k-1}-1).$$

Lemma 7.2(c), (7-3), (7-4), and formula (0-6) show

$$\begin{split} H^{2,+}_{(3)^0} &= 3! H^{1,-}_{(3)^0} \cdot H^{1,-}_{(3)^0} + 3 H^{1,-}_{(3)} \cdot H^{1,-}_{(3)} = 27, \\ H^{2,+}_{(3)} &= 3! H^{1,-}_{(3)^0} \cdot H^{1,-}_{(3)} + 3 H^{1,-}_{(3)} \cdot H^{1,-}_{(3)^2} = -27, \\ H^{2,+}_{(3)^0} &= 3! H^{1,+}_{(3)^0} \cdot H^{1,+}_{(3)^0} + 3 H^{1,+}_{(3)} \cdot H^{1,+}_{(3)} = 24 + 3 H^{1,+}_{(3)} \cdot H^{1,+}_{(3)}, \\ H^{2,+}_{(3)} &= 3! H^{1,+}_{(3)^0} \cdot H^{1,+}_{(3)} + 3 H^{1,+}_{(3)} \cdot H^{1,+}_{(3)^2} = 12 H^{1,+}_{(3)} + 3 H^{1,+}_{(3)} \cdot H^{1,+}_{(3)^2}, \\ H^{1,+}_{(3)^2} &= 3! H^{1,+}_{(3)^0} \cdot H^{0,+}_{(3)^2} + 3 H^{1,+}_{(3)} \cdot H^{0,+}_{(3)^3} = 4 - H^{1,+}_{(3)}. \end{split}$$

It follows that $H_{(3)}^{1,+} = -1$. Hence, Lemma 7.2(c), (7-4), and formula (0-6) give

(7-5)
$$H_{(3)^k}^{1,+} = 3! H_{(3)^0}^{1,+} \cdot H_{(3)^k}^{0,+} + 3 H_{(3)}^{1,+} \cdot H_{(3)^{k+1}}^{0,+} = (-1)^k 2^k + 1$$

Step 4. It remains to compute $H_{(3)^k}^{h,p}$ for $h \ge 2$. The formula (0-6) gives

$$H_{(3)^{k}}^{h,p} = 3! H_{(3)^{0}}^{h-1,p} \cdot H_{(3)^{k}}^{1,+} + 3 H_{(3)}^{h-1,p} \cdot H_{(3)^{k+1}}^{1,+}$$

From this, we can deduce that, for $h \ge 2$,

$$(7-6) \quad \begin{pmatrix} H_{(3)^{k}}^{h,p} \\ H_{(3)^{k+1}}^{h,p} \end{pmatrix} = \begin{pmatrix} 3!H_{(3)^{k}}^{1,+} & 3H_{(3)^{k+1}}^{1,+} \\ 3!H_{(3)^{k+1}}^{1,+} & 3H_{(3)^{k+2}}^{1,+} \end{pmatrix} \begin{pmatrix} H_{(3)^{0}}^{h-1,p} \\ H_{(3)}^{h-1,p} \end{pmatrix}$$
$$= \begin{pmatrix} 3!H_{(3)^{k}}^{1,+} & 3H_{(3)^{k+1}}^{1,+} \\ 3!H_{(3)^{k+1}}^{1,+} & 3H_{(3)^{k+2}}^{1,+} \end{pmatrix} \begin{pmatrix} 3!H_{(3)^{0}}^{1,+} & 3H_{(3)}^{1,+} \\ 3!H_{(3)}^{1,+} & 3H_{(3)^{k+2}}^{1,+} \end{pmatrix} \begin{pmatrix} 3!H_{(3)^{0}}^{1,+} & 3H_{(3)^{2}}^{1,+} \\ 3!H_{(3)}^{1,+} & 3H_{(3)^{k+2}}^{1,+} \end{pmatrix} \begin{pmatrix} 3!H_{(3)^{0}}^{1,+} & 3H_{(3)^{2}}^{1,+} \\ 3!H_{(3)}^{1,+} & 3H_{(3)^{2}}^{1,+} \end{pmatrix}^{h-2} \begin{pmatrix} H_{(3)^{0}}^{1,p} \\ H_{(3)}^{1,p} \end{pmatrix}.$$

Therefore, (7-3), (7-5), and (7-6) complete the proof.

References

- [Arbarello et al. 2011] E. Arbarello, M. Cornalba, and P. A. Griffiths, *Geometry of algebraic curves*, vol. II, Grundlehren der Mathematischen Wissenschaften **268**, Springer, Heidelberg, 2011. MR 2012e:14059 Zbl 1235.14002
- [Atiyah 1971] M. F. Atiyah, "Riemann surfaces and spin structures", *Ann. Sci. École Norm. Sup.* (4) **4** (1971), 47–62. MR 44 #3350 Zbl 0212.56402
- [Cornalba 1989] M. Cornalba, "Moduli of curves and theta-characteristics", pp. 560–589 in *Lectures on Riemann surfaces* (Trieste, 1987), edited by M. Cornalba et al., World Scientific, Teaneck, NJ, 1989. MR 91m:14037 Zbl 0800.14011
- [Eskin et al. 2008] A. Eskin, A. Okounkov, and R. Pandharipande, "The theta characteristic of a branched covering", *Adv. Math.* **217**:3 (2008), 873–888. MR 2008k:14065 Zbl 1157.14014
- [Gunningham 2012] S. Gunningham, "Spin Hurwitz numbers and topological quantum field theory", preprint, 2012. arXiv 1201.1273
- [Harris and Morrison 1998] J. Harris and I. Morrison, *Moduli of curves*, Graduate Texts in Mathematics **187**, Springer, New York, 1998. MR 99g:14031 Zbl 0913.14005
- [Ionel and Parker 2003] E.-N. Ionel and T. H. Parker, "Relative Gromov–Witten invariants", *Ann. of Math.* (2) **157**:1 (2003), 45–96. MR 2004a:53112 Zbl 1039.53101
- [Ionel and Parker 2004] E.-N. Ionel and T. H. Parker, "The symplectic sum formula for Gromov–Witten invariants", *Ann. of Math.* (2) **159**:3 (2004), 935–1025. MR 2006b:53110 Zbl 1075.53092
- [Kiem and Li 2007] Y.-H. Kiem and J. Li, "Gromov–Witten invariants of varieties with holomorphic 2-forms", preprint, 2007. arXiv 0707.2986
- [Kiem and Li 2011] Y.-H. Kiem and J. Li, "Low degree GW invariants of spin surfaces", *Pure Appl. Math. Q.* **7**:4 (2011), 1449–1475. MR 2918169 Zbl 06107784
- [Lee 2013] J. Lee, "Sum formulas for local Gromov–Witten invariants of spin curves", *Trans. Amer. Math. Soc.* **365**:1 (2013), 459–490. MR 2984064
- [Lee and Parker 2007] J. Lee and T. H. Parker, "A structure theorem for the Gromov–Witten invariants of Kähler surfaces", J. Differential Geom. 77:3 (2007), 483–513. MR 2010b:53159 Zbl 1130.53059
- [Lee and Parker 2012] J. Lee and T. H. Parker, "Recursion formulas for spin Hurwitz numbers", preprint, 2012. arXiv 1212.1825
- [Maulik and Pandharipande 2008] D. Maulik and R. Pandharipande, "New calculations in Gromov–Witten theory", *Pure Appl. Math. Q.* **4**:2 (2008), 469–500. MR 2009d:14073 Zbl 1156.14042
- [Mumford 1971] D. Mumford, "Theta characteristics of an algebraic curve", *Ann. Sci. École Norm. Sup.* (4) **4** (1971), 181–192. MR 45 #1918 Zbl 0216.05904
- [Okounkov and Pandharipande 2006] A. Okounkov and R. Pandharipande, "Gromov–Witten theory, Hurwitz theory, and completed cycles", *Ann. of Math.* (2) **163**:2 (2006), 517–560. MR 2007b:14123 Zbl 1105.14076

DEGREE-THREE SPIN HURWITZ NUMBERS

Received May 7, 2012. Revised September 1, 2012.

JUNHO LEE DEPARTMENT OF MATHEMATICS UNIVERSITY OF CENTRAL FLORIDA ORLANDO, FL 32816 UNITED STATES

junlee@mail.ucf.edu

(ℤ₂)³-COLORINGS AND RIGHT-ANGLED HYPERBOLIC 3-MANIFOLDS

YOULIN LI AND JIMING MA

For a compact 3-manifold N with connected nonempty boundary, let Γ be an admissible trivalent graph in ∂N that decomposes ∂N into a set of disks. As an extension of small covers, from a $(\mathbb{Z}_2)^3$ -coloring λ on $\partial N - \Gamma$, one can get a closed 3-manifold M_{λ} that admits a locally standard $(\mathbb{Z}_2)^3$ -action.

Suppose N is irreducible and atoroidal: say, a handlebody. We give a combinatorial necessary and sufficient condition for a $(\mathbb{Z}_2)^3$ -colorable pair (N, Γ) to admit a right-angled hyperbolic structure, which naturally induces a hyperbolic structure on M_{λ} .

1. Introduction

In this note, we study polyhedral hyperbolic 3-manifolds admitting $(\mathbb{Z}_2)^3$ -colorings on their connected boundaries, which correspond to closed hyperbolic 3-manifolds admitting locally standard $(\mathbb{Z}_2)^3$ -actions.

 $(\mathbb{Z}_2)^3$ -colorings and locally standard $(\mathbb{Z}_2)^3$ -actions. Small covers, or Coxeter orbifolds, were studied in [Davis and Januszkiewicz 1991]. They are a class of manifolds which admit locally standard $(\mathbb{Z}_2)^n$ -actions, such that the orbit spaces are *n*-dimensional simple polyhedra. The algebraic and topological properties of a small cover are closely related to the combinatorics of the orbit polyhedron and the coloring on its boundary. For example, the (mod 2) Betti number β_i of a small cover *M* agrees with h_i , where $h = (h_0, h_1, \ldots, h_n)$ is the *h*-vector of the polyhedron.

Those manifolds admitting locally standard $(\mathbb{Z}_2)^n$ -actions form a wider class than small covers. In this paper, we focus on the 3-dimensional case.

MSC2010: primary 57M50, 57M60; secondary 52B70.

The first author was partially supported by NSFC 11001171. The second author was partially supported by NSFC 10901038 and Shanghai NSF 10ZR1403600.

Keywords: $(\mathbb{Z}_2)^3$ -action, hyperbolic structure with polyhedral boundary, 3-manifold.

A standard representation of the $(\mathbb{Z}_2)^3$ -action on \mathbb{R}^3 is the natural action defined by

(1-1)
$$e_1: (x_1, x_2, x_3) \mapsto (-x_1, x_2, x_3),$$

(1-2)
$$e_2: (x_1, x_2, x_3) \mapsto (x_1, -x_2, x_3)$$

(1-3) $e_3: (x_1, x_2, x_3) \mapsto (x_1, x_2, -x_3).$

The actions e_1 , e_2 and e_3 generate the group $(\mathbb{Z}_2)^3$. This action fixes the origin of \mathbb{R}^3 such that its orbit space is the positive cone

$$\mathbb{R}^3_{\geq 0} = \{ (x_1, x_2, x_3) \in \mathbb{R}^3 \mid x_i \ge 0 \}$$

Definition 1.1. An effective $(\mathbb{Z}_2)^3$ -action on a 3-dimensional closed manifold M is said to be *locally standard* if it locally looks like the standard representation of $(\mathbb{Z}_2)^3$ -action on \mathbb{R}^3 . More precisely, if for each point x in M, there is a $(\mathbb{Z}_2)^3$ -invariant neighborhood U_x of x such that U_x is equivariantly homeomorphic to an invariant open subset of the standard $(\mathbb{Z}_2)^3$ -action on \mathbb{R}^3 .

The orbit space of a locally standard $(\mathbb{Z}_2)^3$ -action on a 3-dimensional closed manifold M is a compact manifold N with corners. In other words, it is a 3dimensional compact manifold N with a graph Γ on ∂N . The graph Γ on ∂N induces a cell decomposition on ∂N . The vertices of Γ are the image of fixed points of the $(\mathbb{Z}_2)^3$ -action, the (open) edges of Γ are the image of fixed points of subgroups $(\mathbb{Z}_2)^2 < (\mathbb{Z}_2)^3$ and (open) components of $\partial N - \Gamma$ are the image of fixed points of subgroups $\mathbb{Z}_2 < (\mathbb{Z}_2)^3$.

Definition 1.2. Let *N* be a 3-dimensional manifold with nonempty boundary, and Γ a trivalent graph in ∂N that gives a cell decomposition of ∂N . A $(\mathbb{Z}_2)^3$ -coloring is a map $\lambda : \partial N - \Gamma \to (\mathbb{Z}_2)^3 - 0$ such that $\lambda(f_1), \lambda(f_2)$ and $\lambda(f_3)$ generate $(\mathbb{Z}_2)^3$ for each triple of faces f_1, f_2 and f_3 sharing a common vertex.

Associated to a locally standard $(\mathbb{Z}_2)^3$ -action on M, there is a canonical $(\mathbb{Z}_2)^3$ coloring λ on $\partial N - \Gamma$ which colors each face $f \in \partial N - \Gamma$ by the element $e \in (\mathbb{Z}_2)^3$ that fixes f. For an *i*-dimensional cell f in the cell decomposition, i = 0, 1, 2, we have a group $G_f \cong (\mathbb{Z}_2)^{3-i}$ which is generated by the colorings in the faces which are adjacent to f. The locally standard $(\mathbb{Z}_2)^3$ -action on M induces a principal $(\mathbb{Z}_2)^3$ -bundle over N.

Conversely, by Lemma 3.1 of [Lü and Masuda 2009], from a $(\mathbb{Z}_2)^3$ -coloring λ on $(\partial N, \Gamma)$ and a principal $(\mathbb{Z}_2)^3$ -bundle over N, we can get a unique closed 3-manifold M. In particular, from a $(\mathbb{Z}_2)^3$ -coloring λ and the trivial principal $(\mathbb{Z}_2)^3$ -bundle over N, we get a 3-manifold M_{λ} which depends only on the coloring λ . By this we take eight copies of N, $N \times \{\alpha\}$ for each $\alpha \in (\mathbb{Z}_2)^3$, and construct a quotient

space M_{λ} under the following gluing rule:

(1-4)
$$(x, \alpha_1) \sim (y, \alpha_2) \Leftrightarrow \begin{cases} x = y, \ \alpha_1 = \alpha_2, & \text{if } x \text{ lies in the interior of } N, \\ x = y, \ \alpha_1 \alpha_2^{-1} \in G_f, & \text{if } x \text{ lies in a cell } f. \end{cases}$$

Then it is easy to see M_{λ} is a closed 3-manifold. In this paper, we only consider closed 3-manifolds associated to $(\mathbb{Z}_2)^3$ -colorings and trivial principal $(\mathbb{Z}_2)^3$ -bundles over N.

A simple example is that if we consider a coloring of the four faces of a tetrahedron by e_1 , e_2 , e_3 , $e_1+e_2+e_3$, respectively, then from the above construction, we get the closed orientable 3-manifold RP^3 . A tetrahedron admits a unique right-angled spherical structure, and the spherical structures on eight copies of the tetrahedron, when glued together, give rise to the unique spherical structure on RP^3 . This point of view is applied in this paper.

There are many works on manifolds with locally standard $(\mathbb{Z}_2)^3$ -actions. For example, 3-dimensional small covers are studied in [Lü 2009; Lü and Yu 2011]. Six operations on small covers were defined in [Lü and Yu 2011], which topologically behave well, such that every 3-dimensional small cover is obtained from the two simple small covers RP^3 and $S^1 \times RP^2$ by a sequence of these operations. It should be noted that the operations in [Lü and Yu 2011] give many disks in a simple convex polygon P, which intersects the 1-skeleton of P in at most four points but which is not vertex-linking or edge-linking. So the preimage of these disks are essential spheres or essential tori in the small cover M in general, and hence M does not admit a geometric structure [Scott 1983].

Polyhedral hyperbolic 3-manifolds. Andreev [1971] (see also [Roeder et al. 2007]) gives a complete characterization of compact hyperbolic polyhedra in dimension 3 with nonobtuse angles. The boundary of a compact hyperbolic polyhedron inherits a natural cell decomposition. The 1-skeleton of the cell decomposition is a graph Γ on the boundary of the 3-ball, and a dihedral angle is also given on each edge of Γ from the hyperbolic structure. Andreev's theorem is given in terms of a set of conditions on the dihedral angles. Besides its beauty, Andreev's theorem is also essential in the proof of Thurston's geometrization theorem for Haken 3-manifolds. The natural question is, given a cell decomposition of the boundary of the 3-ball, and a weight $\alpha_e \in (0, \pi)$ attached to each edge *e* of the cell decomposition, whether there is a compact hyperbolic polyhedron in \mathbb{H}^3 realizing this cell decomposition whose dihedral angles coincide with the attached weights. This question is still open now.

A clever approach for working with compact hyperbolic polyhedra having arbitrary dihedral angles is to express necessary and sufficient conditions for existence of a given polyhedron in terms of its polar dual in the de Sitter space; see [Hodgson and Rivin 1993]. For a generalization of Andreev's result to ideal and hyper-ideal hyperbolic polyhedra, see [Rivin 1996; Bao and Bonahon 2002]. Hyperbolic structures on topologically more complicated 3-manifold *N* with boundary are also studied. See [Schlenker 2002; 2003; 2005; 2006, Fillastre and Izmestiev 2009; 2011; Guéritaud 2009].

Suppose N is a compact 3-manifold with connected and nonempty boundary. In this note, we consider the right-angled hyperbolic structures on N with compressible boundary.

Given a graph Γ in ∂N , we call it *admissible* if the lift $\tilde{\Gamma}$ of Γ in the universal cover of N, say \tilde{N} , gives a cell decomposition of $\partial \tilde{N}$ such that each of its 2-cells has a closure homeomorphic to a disk, and each pair of such two disks shares at most one edge in $\tilde{\Gamma}$. A *right-angled hyperbolic realization* of (N, Γ) is a complete compact hyperbolic manifold N^* with right-angled polyhedral boundary (i.e., modeled on the orthogonal intersection of two half-spaces with totally geodesic boundaries in H^3 and having finite volume), endowed with a homeomorphism to N that sends the nonsmooth points of N^* precisely to the points of Γ . The nonsmooth points of N^* will be called the *singular locus* of this structure. From the homeomorphism between N^* and N, Γ is also called the *singular locus* for this hyperbolic realization. We will call such a structure on N^* a *hyperbolic structure with right-angled polyhedral boundary* on N. Hence these kinds of hyperbolic structures on N^3 .

Similar to all results above, it is interesting to give a kind of characterization of hyperbolic structures with right-angled polyhedral boundary on N. Since all the dihedral angles are right-angled, an easy argument shows that the graph Γ defined above must be trivalent.

It is well known that most 3-manifolds are hyperbolic 3-manifolds [Thurston 1982]. So it is interesting to consider locally standard $(\mathbb{Z}_2)^3$ -actions on closed hyperbolic 3-manifolds. It is natural to ask which closed hyperbolic 3-manifold admits a locally standard $(\mathbb{Z}_2)^3$ -action. The orbit space of a locally standard $(\mathbb{Z}_2)^3$ -action is a compact manifold N with a coloring λ in ∂N . If (N, Γ) admits a right-angled hyperbolic structure, then it is easy to see that M is hyperbolic. A pair (N, Γ) admits a unique right-angled hyperbolic structure. However, it may admit many different colorings. Each coloring, together with a principle $(\mathbb{Z}_2)^3$ -bundle over N, gives a manifold with a locally standard $(\mathbb{Z}_2)^3$ -action. So these give many different hyperbolic manifolds of the same volume. [Inoue 2008] gives a very clear description of right-angled hyperbolic polyhedra from this point of view.

The most interesting case is that N is a handlebody, or the simplest one, a 3-ball.

Main result. Suppose Γ is an admissible graph in ∂N . For a vertex v of Γ , we take a small closed regular neighborhood B of v in N, then B intersects N – int B in a disk D_v . We call D_v a vertex-linking disk. It intersects Γ in three points. The preimage of a vertex-linking disk in M_λ is a sphere which bounds a 3-ball in M_λ for

a $(\mathbb{Z}_2)^3$ -coloring λ . For an edge e of Γ , we also take a small closed neighborhood B of e in N; then B intersects N – int B in a disk D_e . We call D_e an *edge-linking disk*. It intersects Γ in four points. The preimage of an edge-linking disk in M_λ is a torus (or a Klein bottle) which bounds a solid torus (or a solid Klein bottle) in M_λ . We say a properly embedded disk D in the 3-manifold N intersects Γ *efficiently* if ∂D and Γ are in general position and there is no bigon in $\partial N - (\partial D \cup \Gamma)$. In this note we always assume a disk D intersects with Γ efficiently.

Our main result is the following:

Theorem 1.3. Let N be an irreducible, atoroidal and compact 3-manifold with connected nonempty boundary, and Γ be an admissible trivalent graph in ∂N which gives a cell decomposition of ∂N , such that $(\partial N, \Gamma)$ admits a $(\mathbb{Z}_2)^3$ -coloring. Then (N, Γ) realizes a right-angled hyperbolic structure if and only if every properly embedded disk D in N has $|D \cap \Gamma| \ge 5$, except when D is a vertex-linking disk or an edge-linking disk. Moreover, the realization is unique up to isometry.

Remark 1.4. In practice, much attention has been paid to the right-angled hyperbolic structures on handlebodies. They are irreducible, atoroidal and compact 3-manifolds with connected nonempty boundaries. So Theorem 1.3 can be applied to the handlebody case.

Remark 1.5. There are two canonical ways to study polyhedral hyperbolic structure on 3-manifold *M*: Alexandrov's method and the variational method; see, for example, [Fillastre and Izmestiev 2011]. Our approach in this note uses the doubling trick. A $(\mathbb{Z}_2)^3$ -coloring helps us find a closed 3-manifold on which we can apply the geometrization theorem.

2. Preliminaries

If (N, Γ) admits a right-angled hyperbolic structure, then Γ is admissible, and each of its 2-dimensional faces is a right-angled hyperbolic *n*-polygon. So $n \ge 5$.

Definition 2.1. Let Γ^* be the dual graph of Γ in ∂N . A *k*-circuit is a simple closed curve *C* in Γ^* consisting of *k* successive edges of Γ^* which is contractible in ∂N . A circuit is *elementary* if it bounds a disk *D* in ∂N and there is exactly one vertex *V* of Γ that lies in *D*. A *k*-circuit is *prismatic* if the endpoints of all the edges of Γ which intersect *C* are distinct.

Obviously, there is a one-to-one correspondence between edges of Γ and those of Γ^* .

Lemma 2.2. Suppose Γ is admissible. If *C* is a 3-circuit which is not prismatic, then *C* is isotopic to the boundary of a vertex-linking disk. If Γ contains no prismatic 3-circuit, and *C* is a 4-circuit which is not prismatic, then *C* is isotopic to the boundary of an edge-linking disk.

Proof. The proof is similar to the proofs of Lemmas 1.2 and 1.3 of [Roeder et al. 2007], by which we have that, if *C* is a nonprismatic 3-circuit, then it is an elementary circuit. So it is isotopic to the boundary of a vertex-linking disk. If Γ contains no prismatic 3-circuit, then every nonprismatic 4-circuit *C* separates off exactly two vertices of Γ from the remaining vertices of Γ , which in turn implies that *C* is isotopic to the boundary of an edge-linking disk. Actually, the authors of [Roeder et al. 2007] proved this for any graph in S^2 . Since Γ is admissible, their arguments can be extended verbatim in the general case.

We give a proposition on the orientability of a 3-manifold M_{λ} with a locally standard $(\mathbb{Z}_2)^3$ -action and trivial principal $(\mathbb{Z}_2)^3$ -bundle.

Proposition 2.3. Suppose N is a compact orientable connected 3-manifold with connected boundary. Then, for a $(\mathbb{Z}_2)^3$ -coloring λ on $(\partial N, \Gamma)$, M_{λ} is orientable if and only if there is a basis $\{e_1, e_2, e_3\}$ of $(\mathbb{Z}_2)^3$, such that the image of λ is contained in $\{e_1, e_2, e_3, e_1 + e_2 + e_3\}$.

Proof. For small covers, this proposition has been proved in Theorem 1.7 of [Nakayama and Nishimura 2005]. Recall that M_{λ} is determined by the coloring λ and the trivial principal $(\mathbb{Z}_2)^3$ -bundle over N. So M_{λ} is obtained by gluing eight copies of (N, λ) , and M_{λ} is orientable if and only if $H_3(M_{\lambda}, \mathbb{Z}) = \mathbb{Z}$. To calculate $H_3(M_{\lambda}, \mathbb{Z})$, we only need to consider the 3-cells and 2-cells in a cellular decomposition of M_{λ} , which is induced by a cellular decomposition of (N, Γ) . Note that ∂N is connected, so the arguments of the proof of Theorem 1.7 of [Nakayama and Nishimura 2005] hold in our case word-by-word.

3. Proof of Theorem 1.3

Proof of the necessity part of Theorem 1.3. Suppose (N, Γ) realizes a right-angled hyperbolic structure. If $D \subset N$ is a properly embedded disk which intersects Γ efficiently, and is not vertex-linking or edge-linking, then by Gauss–Bonnet theorem, we have $|D \cap \Gamma| \ge 5$. So the necessity part of Theorem 1.3 follows.

Proof of the sufficiency part of Theorem 1.3 in the case that M_{λ} **is orientable.** Recall that a closed orientable 3-manifold *M* is *irreducible* if every embedded 2-sphere *S* in *M* bounds a 3-ball; otherwise *M* is *reducible*. An embedded 2-sphere *S* which does not bound a 3-ball in *M* is called *essential*. A closed irreducible orientable 3-manifold *M* is *atoroidal* if every embedded torus *T* in *M* bounds a solid torus; otherwise *M* is *toroidal*. An embedded torus *T* which does not bound a solid torus is *essential* in *M*. See [Hempel 1976] or [Jaco 1980].

We need the equivariant sphere theorem of Meeks, Simon, and Yau [Meeks et al. 1982], but the reformulation by Dunwoody [1985] is more convenient for us.

Theorem 3.1. Let G be a finite group that acts on a closed orientable 3-manifold M by homeomorphisms. Suppose M has a G-equivariant triangulation. If there exists an essential 2-sphere S in M, then there exists an essential 2-sphere S_1 in M which is in general position with respect to the triangulation, such that $g(S_1) = S_1$ or $g(S_1) \cap S_1 = \emptyset$ for every $g \in G$.

We also need the following equivariant torus theorem; see [Freedman et al. 1983; Jaco and Shalen 1979; Johannson 1979].

Theorem 3.2. Let G be a finite group which acts on a closed 3-manifold M by homeomorphisms. Suppose M is irreducible, orientable and contains an essential torus. Then either M is Seifert-fibered, or M contains a G-equivariant essential torus.

First, we show the following lemmas.

Lemma 3.3. M_{λ} is irreducible.

Proof. We give a triangulation \mathcal{T} of N, such that the graph Γ is contained in the 1-skeleton of \mathcal{T} . So the triangulation \mathcal{T} induces a triangulation of M_{λ} .

If M_{λ} is reducible, then by the equivariant sphere theorem, there is a $(\mathbb{Z}_2)^3$ equivariant sphere *S* which is essential in M_{λ} . We denote $S \cap N \times \{1\}$ by *A*, which is a compact surface with nonempty boundary if $A \neq \emptyset$. We may assume $A \neq \emptyset$, otherwise we can use the $(\mathbb{Z}_2)^3$ -action to find another $(\mathbb{Z}_2)^3$ -invariant sphere *S'* which has nonempty intersection with $N \times \{1\}$. Since *A* is obtained from *S* by the $(\mathbb{Z}_2)^3$ -action and *S* is connected, *A* is connected.

Since *S* is in general position with respect the triangulation of M_{λ} , *A* is in general position with respect to the triangulation of ∂N , in particular, with respect to Γ . So there is a cell decomposition of ∂A : for each face *f* of $\partial N - \Gamma$, $f \cap A$ is an edge in ∂A . Moreover, the coloring on $\partial N - \Gamma$ now induces a coloring on ∂A , which we denote by λ_A , and *S* is obtained from copies of *A* by the gluing rule from λ_A .

By Definition 1.2, the colorings on any two adjacent edges of ∂A are different. So we have a subgroup *G* of $(\mathbb{Z}_2)^3$ which has index 1 or 2 in $(\mathbb{Z}_2)^3$, such that for any $g \in G$ we have g(S) = S, and for any $h \in (\mathbb{Z}_2)^3 - G$ we have $h(S) \cap S = \emptyset$.

In other words, *S* is obtained by gluing 4 or 8 copies of *A*, and the edges in ∂A contribute a 4-valence graph in *S*. So $\chi(S) = m(\chi(A) - E/2 + E/4) = 2$, where *E* is the number of edges in ∂A , and m = 4 or 8. If m = 4, then E = 2 and $\chi(A) = 1$. So *A* is a disk with ∂A consisting of 2 edges. This is impossible by the assumption in Theorem 1.3. If m = 8, then E = 3 and $\chi(A) = 1$. So *A* is a disk with ∂A consists of 3 edges. Moreover, we have that $\partial A \cap f$ is connected for each face *f*. Suppose otherwise; i.e., suppose that $\partial A \cap f$ consists of at least two arcs. We have an edge *e* of Γ which intersects *A* such that the two sides of *e* both are in the face *f*. Then when we lift ∂N to the universal cover \tilde{N} . The closure of the lifting \tilde{f} of the face *f*

is not a disk, contradicting the assumption that Γ is admissible. So ∂A is a 3-circuit. Thus, by the assumption and Lemma 2.2, A is a vertex-linking disk in N. The preimage of a vertex-linking disk in M_{λ} is a sphere which bounds a 3-ball in M_{λ} . This contradicts the assumption that S is essential in M_{λ} . So M_{λ} is irreducible. \Box

Lemma 3.4. If M_{λ} is a toroidal Seifert manifold, then there is an essential torus in M_{λ} which is $(\mathbb{Z}_2)^3$ -equivariant.

Proof. Suppose e_1 , e_2 and e_3 are three orientation-reversing involutions which generate the $(\mathbb{Z}_2)^3$ -action. Since M_λ admits orientation-reversing involutions, according to Theorems 8.2 and 8.5 of [Neumann and Raymond 1978], M_λ is Seifert-fibered with Euler number 0; i.e., M_λ contains horizontal incompressible surfaces which are transversal to each fiber. In other words, M_λ is a surface bundle over S^1 with horizontal incompressible surfaces as surface fibers. We already proved in Lemma 3.3 that M_λ is irreducible, so the Euler characteristic of the base orbifold of M_λ is negative or zero. Thus M_λ admits the geometries $H^2 \times \mathbb{R}$ or E^3 . We refer the readers to [Scott 1983] for the details about these two geometries.

For each i = 1, 2, 3, Fix (e_i) contains no nonorientable closed surfaces since the nonorientable closed surfaces are one-sided in M_{λ} . According to [Meeks and Scott 1986], for each $i = 1, 2, 3, e_i$ is isotopic to an isometry. So Fix (e_i) consists of some totally geodesic, and hence incompressible, closed surfaces in M_{λ} .

If M_{λ} admits the $H^2 \times \mathbb{R}$ geometry, then it has unique Seifert fibration structure. So each homeomorphism sends regular fibers to regular fibers. Then, among Fix (e_1) , Fix (e_2) and Fix (e_3) , at least two of them, say Fix (e_1) and Fix (e_2) , consist of vertical essential tori, and e_3 keeps each regular fiber invariant and reverses its orientation. By the definition of $(\mathbb{Z}_2)^3$ -action, Fix $(e_1) \cap$ Fix $(e_2) \neq \emptyset$, and Fix (e_1) intersects Fix (e_2) transversely. Choose a torus component T in Fix $(e_1) = e_2(p) = e_2(p)$. So $e_2(T) \subset$ Fix (e_1) . Thus we have that either $e_2(T) = T$ or $e_2(T) \cap T = \emptyset$. By the assumption $T \cap$ Fix $(e_2) \neq \emptyset$, we have $e_2(T) = T$. Moreover, $e_3(T) = T$. Hence T is invariant by e_1 , e_2 and e_3 , and hence is invariant by each element of the group $(\mathbb{Z}_2)^3$. So it is an essential torus which is $(\mathbb{Z}_2)^3$ -equivariant.

If M_{λ} admits the E^3 geometry, then according to Theorems 8.2 and 8.5 of [Neumann and Raymond 1978], it is either the 3-torus T^3 or the Seifert manifold with invariant {0; (2, 1), (2, -1), (2, 1), (2, -1)}. For the former case, we choose a Seifert fibration structure which is fibred by all circles isotopic to the circles in Fix(e_1) \cap Fix(e_2). Then we can apply the same argument as above to obtain a $(\mathbb{Z}_2)^3$ -equivariant essential torus. For the latter case, we fix the Seifert fibration structure given before, and then all horizontal incompressible surfaces in M_{λ} are isotopic essential tori. This is because M_{λ} has a unique structure of surface bundles over S^1 , since its first Betti number is 1; see [Thurston 1986]. So we can still

assume that both $Fix(e_1)$ and $Fix(e_2)$ consist of vertical essential tori, and e_3 keeps each regular fiber invariant and reverses its orientation. The same argument in the previous paragraph still applies, and the same conclusion still holds.

Lemma 3.5. M_{λ} is atoroidal.

Proof. By Theorem 3.2 and Lemma 3.4, if M_{λ} is toroidal, then there is a $(\mathbb{Z}_2)^3$ -equivariant essential torus $T \subset M_{\lambda}$.

Similar to the sphere case in Lemma 3.3, we also give a triangulation \mathcal{T} of N. We denote $T \cap N \times \{1\}$ by A, which is nonempty, and is a compact connected surface with nonempty boundary. Also, there is a cell decomposition of ∂A induced from the triangulation of N.

Similar to the argument in the sphere case in Lemma 3.3, we have $\chi(T) =$ $m(\chi(A) - E/2 + E/4) = 0$, where E is the number of edges in ∂A , and m is an integer. So E = 4 and $\chi(A) = 1$. Thus A is a disk with ∂A consists of 4 edges. Moreover, $\partial A \cap f$ is connected for each face f. Suppose otherwise; i.e., suppose that $\partial A \cap f$ consists of at least two arcs. When we lift ∂N to the universal cover N, two of the four edges forming ∂A belong to the same face. If these two edges are adjacent in ∂A , then by the same argument as in the proof of Lemma 3.3, we obtain a contradiction. If these two edges are not adjacent in ∂A , then the lift of these two arcs in the universal cover are identified. So in the universal cover, there are two disks which share two distinct edges, contradicting the assumption that Γ is admissible. Thus ∂A is a 4-circuit. Therefore, by the assumption and Lemma 2.2, A is an edge-linking disk. The preimage of an edge-linking disk in M_{λ} is a torus (or a Klein bottle) which bounds a solid torus (or a solid Klein bottle, which is impossible since we assume M_{λ} is orientable in this subsection), so it is not essential. This contradicts the assumption that T is essential. So M_{λ} is atoroidal.

Lemma 3.6. M_{λ} is not a Seifert manifold.

Proof. Suppose M_{λ} is an orientable Seifert manifold with orientable base orbifold, and M_{λ} is neither a lens space nor S^3 . Here the lens spaces don't include S^3 or $S^2 \times S^1$. By Theorem 8.2 of [Neumann and Raymond 1978] and its proof, if there is an orientation-reversing involution on M_{λ} , then the Seifert invariant of M_{λ} is $\{g; (a_1, b_1), (a_1, -b_1), (a_2, b_2), (a_2, -b_2), \ldots, (a_t, b_t), (a_t, -b_t)\}$, where g is the genus of the base orbifold. Since M_{λ} is atoroidal, we have g = 0 and t = 1, and hence M_{λ} is a lens space. This is a contradiction.

Suppose M_{λ} is an orientable Seifert manifold with nonorientable base orbifold, and M_{λ} is not a lens space. By Theorem 8.5 of [Neumann and Raymond 1978] and its proof, if there is an orientation-reversing involution on M_{λ} , then the Seifert invariant of M_{λ} is {k; (a_1, b_1) , $(a_1, -b_1)$, (a_2, b_2) , $(a_2, -b_2)$, ..., (a_t, b_t) , $(a_t, -b_t)$ }, where k is the genus of the nonorientable base orbifold. If $t \ge 1$, then M_{λ} cannot be atoroidal. If t = 0, then M_{λ} is either reducible or toroidal. In both cases, we arrive at contradictions.

Suppose M_{λ} is a lens space. By the main result in [Kwun 1970], among all lens spaces, only RP^3 admits orientation-reversing involutions. Moreover, RP^3 admits exactly one orientation-reversing involution up to isotopies, and the set of fixed points of this involution is an RP^2 , which has Euler characteristic 1. However, according to the definition of locally standard $(\mathbb{Z}_2)^3$ -action, for any nontrivial element $e \in (\mathbb{Z}_2)^3$, its fixed point set Fix(e) is a union of k-polygons ($k \ge 5$ by our assumption), and each vertex in Fix(e) is adjacent to 4 edges. Let v be the number of vertices in Fix(e). Then the number of edges in Fix(e) is 2v, and the number of faces of Fix(e) is less than or equal to 4v/5. So the Euler characteristic of Fix(e)is negative — a contradiction.

Suppose M_{λ} is the 3-sphere S^3 . From the fact that the orientation-preserving mapping class group of S^3 is trivial, we know S^3 admits exactly one orientation-reversing involution up to isotopy, and the set of fixed points of this involution is an S^2 , which has Euler characteristic 2. Then similar to the argument in the previous paragraph, we get a contradiction. So the lemma follows.

By Lemmas 3.3, 3.5 and 3.6, M_{λ} is a closed, irreducible, and atoroidal manifold which is not Seifert-fibered. So by Perelman's proof of Thurston's geometrization theorem (see [Cao and Zhu 2006; Bessières et al. 2010; Kleiner and Lott 2008; Morgan and Tian 2007]), M_{λ} is a hyperbolic 3-manifold. By [Dinkelbach and Leeb 2009], every smooth action of a finite group on a hyperbolic 3-manifold is conjugate to an isometric action. Since each $e \in (\mathbb{Z}_2)^3$ is conjugate to an isometric involution, its fixed point set is a totally geodesic surface in M_{λ} . Since $(\mathbb{Z}_2)^3$ is an Abelian group, by elementary arguments for the isometric group of hyperbolic 3-space H^3 , all these totally geodesic surfaces intersect orthogonally. So the hyperbolic structure on M_{λ} induces a hyperbolic structure on (N, Γ) . Conversely, each right-angled hyperbolic structure on (N, Γ) induces a hyperbolic structure on M_{λ} . So the right-angled realization of (N, Γ) is unique. This ends the proof of Theorem 1.3 in the case that M_{λ} is orientable.

Proof of the sufficiency part of Theorem 1.3 in the case that M_{λ} **is nonorientable.** Let $\pi : \tilde{M}_{\lambda} \to M_{\lambda}$ be the orientable double cover of M_{λ} , and τ be the covering transformation of \tilde{M}_{λ} . Note that τ is orientation-reversing. By the lifting theorem, for each *i*, e_i lifts to an action, say \tilde{e}_i , on \tilde{M}_{λ} such that $\tilde{e}_i(x_0) = x_0$, where $x_0 \in \tilde{M}_{\lambda}$ projects to a vertex of Γ in ∂N .

We show that \tilde{e}_i and \tilde{e}_j commute, for $1 \le i, j \le 3$. It is easy to verify that $\tilde{e}_i \tilde{e}_j$ is the lift of $e_i e_j$, and $\tilde{e}_j \tilde{e}_i$ is the lift of $e_j e_i$. Since $e_i e_j = e_j e_i$, and $\tilde{e}_i \tilde{e}_j (x_0) = \tilde{e}_j \tilde{e}_i (x_0)$, by the unique lifting property, $\tilde{e}_i \tilde{e}_j = \tilde{e}_j \tilde{e}_i$. We also show that τ and \tilde{e}_i commute,

for $1 \le i \le 3$. It is easy to verify that both $\tau \tilde{e}_i$ and $\tilde{e}_i \tau$ are lifts of e_i . So either $\tau \tilde{e}_i = \tilde{e}_i \tau$ or $\tau \tau \tilde{e}_i = \tilde{e}_i \tau$. The latter is $\tilde{e}_i = \tilde{e}_i \tau$ in fact, which is impossible. So $\tau \tilde{e}_i = \tilde{e}_i \tau$. Therefore we have an action of $(\mathbb{Z}_2)^4$ on \tilde{M}_{λ} .

If \tilde{M}_{λ} is a toroidal Seifert manifold, then by Lemma 3.4, there is an essential vertical torus T in \tilde{M}_{λ} which is fixed by \tilde{e}_1 , and is invariant by \tilde{e}_i , for i = 2, 3. For any point $p \in T$, we have $\tilde{e}_1 \tau(p) = \tau \tilde{e}_1(p) = \tau(p)$. So $\tau(T) \subset \text{Fix}(\tilde{e}_1)$. Hence either $\tau(T) = T$ or $\tau(T) \cap T = \emptyset$. It is straightforward to verify that T is $(\mathbb{Z}_2)^4$ -equivariant.

Therefore, similar to the previous subsection, we can prove that \tilde{M}_{λ} is irreducible and atoroidal. Moreover, if \tilde{M}_{λ} is an atoroidal Seifert manifold, then it must be S^3 or RP^3 . The action of τ on \tilde{M}_{λ} has no fixed points. However, as stated in the previous subsection, any orientation-reversing involution on S^3 or RP^3 must have fixed points. We arrive at a contradiction.

So \tilde{M}_{λ} is hyperbolic. Similar to the arguments in the previous subsection, (N, Γ) admits a unique right-angled hyperbolic structure.

4. Examples

In this section we give three examples.

Example 4.1. The simplest way to construct a handlebody which admits rightangled hyperbolic structure is from the Löbell polyhedron L(n) for $n \ge 5$ (see, for example, [Inoue 2008]). A Löbell polyhedron L(n) admits a right-angled hyperbolic structure. Gluing two opposite *n*-gon faces of L(n), we get a solid torus admitting right-angled hyperbolic structures, and whose boundary consists of 2n octagons.

For instance, from L(5), which is a dodecahedron, we can get three solid tori, according to the twisting angle of gluing. All these solid tori satisfy Theorem 1.3. It is easy to see that they admit $(\mathbb{Z}_2)^3$ -colorings, but don't admit one which satisfies the orientability criterion in Proposition 2.3.

This kind of right-angled hyperbolic solid tori are "simple", by which we mean we can obtain a right-angled hyperbolic polyhedron by cutting along a totally geodesic right-angled n-polygon P from the solid tori, where P intersects the boundaries of the solid tori orthogonally.

Example 4.2. A hexagonal tessellation of \mathbb{R}^2 with a coloring is shown in Figure 1. We assume that the diameter of a hexagon is 1. We take a \mathbb{Z}^2 -action on \mathbb{R}^2 , such that its fundamental domain is a rectangle *R* whose vertical edges have length 4.5, and whose horizontal edges have length $3\sqrt{3}$. So there are six hexagons in each horizontal layer and each vertical layer. Gluing the boundaries of *R*, we get a torus *T*.

We can show that any solid torus bounded by T with coloring shown in Figure 1 satisfies the orientability criterion in Proposition 2.3 as well as the assumption of Theorem 1.3. So it admits a right-angled hyperbolic structure.



Figure 1

We fix a homeomorphism from T to the boundary of a solid torus J, so it is natural to ask whether the pair (J, Γ) admits a right-angled hyperbolic structure.

It is easy to see that any essential simple closed curve C in this T intersects Γ in at least five points, and any curve C which bounds a disk D in T intersects Γ in at least five points, unless that D is a vertex-linking disk or an edge-linking disk. So for any solid torus J which is bounded by T, (J, Γ) realizes a right-angled hyperbolic structure.

If the boundary of the unique essential disk in the solid torus J is the image of a horizontal line, then the hyperbolic solid torus can be decomposed into three copies of the Löbell polyhedron L(6) along three totally geodesic right-angled hexagons in the solid torus. The same claim holds if the boundary of the unique essential disk in the solid torus is the image of the straight lines which have angles $\pi/3$ or $2\pi/3$ with the horizontal lines.

Except in these three cases, the right-angled hyperbolic structure cannot be obtained by gluing two faces of a right-angled hyperbolic polyhedron by an isometry,



Figure 2

so it is not "simple". Suppose otherwise; then the totally geodesic right-angled *k*-polygon *P* which decomposes *J* is in general position with Γ , and so some faces of ∂J must be decomposed into a set of right-angled hyperbolic *n*-polygons by ∂P . Note that $n \ge 5$, so if ∂P enters a face *f* of $\partial J - \Gamma$, then it exits *f* from the opposite edge of *f* from where it enters. It is easy to see that ∂P is the image of the lines in \mathbb{R}^2 which have angles $0, \pi/3$ or $2\pi/3$ with the horizontal lines.

Example 4.3. The graph Γ decomposes the torus illustrated in Figure 2 [Chen 2009] into three hexagons, say f_1 , f_2 and f_3 . We color f_i by $e_i \in (\mathbb{Z}_2)^3$ for i = 1, 2, 3. There are two sets of disks in Theorem 1.3. The first one consists of boundary parallel disks. The second one consists of essential disks, i.e., not boundary parallel.

For any embedding of (T^2, Γ) of Figure 2 into a solid torus J, the boundary parallel disks satisfy the assumption of Theorem 1.3. So if we embed (T^2, Γ) into a solid torus J by a map f so that the unique essential disk D (up to isotopy) intersects Γ in at least 5 points, then by Theorem 1.3, we get a right-angled hyperbolic structure on $(J, f(\Gamma))$. Note that for a fixed embedding of $\Gamma \rightarrow T^2$, there are only finitely many isotopy classes of simple closed curves which intersect Γ in at most 4 points.

In general, if the pair $(\partial N, \Gamma)$ admits a $(\mathbb{Z}_2)^3$ -coloring and ∂N has genus at least one, then it may admit many colorings for a re-embedding of Γ into ∂N . This in turn induces many closed 3-manifolds from locally standard $(\mathbb{Z}_2)^3$ -actions.

Acknowledgements

The second author would like to thank Michel Boileau and Jean-Marc Schlenker for discussions on these topics. Both authors would like to thank Zhi Lü for helpful conversations. We also would like to thank the referee for many useful comments.

References

- [Andreev 1971] E. M. Andreev, "On convex polyhedra of finite volume in Lobačevski space", *Math. USSR Sbornik* **12**:3 (1971), 225–259. Zbl 0252.52005
- [Bao and Bonahon 2002] X. Bao and F. Bonahon, "Hyperideal polyhedra in hyperbolic 3-space", *Bull. Soc. Math. France* **130**:3 (2002), 457–491. MR 2003k:52007 Zbl 1033.52009
- [Bessières et al. 2010] L. Bessières, G. Besson, S. Maillot, M. Boileau, and J. Porti, *Geometrisation of 3-manifolds*, EMS Tracts in Mathematics 13, European Mathematical Society, Zürich, 2010. MR 2012d:57027 Zbl 1244.57003
- [Cao and Zhu 2006] H.-D. Cao and X.-P. Zhu, "A complete proof of the Poincaré and geometrization conjectures—application of the Hamilton–Perelman theory of the Ricci flow", *Asian J. Math.* 10:2 (2006), 165–492. MR 2008d:53090 Zbl 1200.53057
- [Chen 2009] B. Chen, A study of 2-torus topology, Ph.D. thesis, Fudan University, 2009. In Chinese.
- [Davis and Januszkiewicz 1991] M. W. Davis and T. Januszkiewicz, "Convex polytopes, Coxeter orbifolds and torus actions", *Duke Math. J.* **62**:2 (1991), 417–451. MR 92i:52012 Zbl 0733.52006
- [Dinkelbach and Leeb 2009] J. Dinkelbach and B. Leeb, "Equivariant Ricci flow with surgery and applications to finite group actions on geometric 3-manifolds", *Geom. Topol.* **13**:2 (2009), 1129–1173. MR 2011b:53158 Zbl 1181.57023
- [Dunwoody 1985] M. J. Dunwoody, "An equivariant sphere theorem", *Bull. London Math. Soc.* **17**:5 (1985), 437–448. MR 87f:57008 Zbl 0592.57005
- [Fillastre and Izmestiev 2009] F. Fillastre and I. Izmestiev, "Hyperbolic cusps with convex polyhedral boundary", *Geom. Topol.* **13**:1 (2009), 457–492. MR 2009i:57039 Zbl 1179.57026
- [Fillastre and Izmestiev 2011] F. Fillastre and I. Izmestiev, "Gauss images of hyperbolic cusps with convex polyhedral boundary", *Trans. Amer. Math. Soc.* **363**:10 (2011), 5481–5536. MR 2012b:57034 Zbl 1238.57018
- [Freedman et al. 1983] M. Freedman, J. Hass, and P. Scott, "Least area incompressible surfaces in 3-manifolds", *Invent. Math.* **71**:3 (1983), 609–642. MR 85e:57012 Zbl 0482.53045
- [Guéritaud 2009] F. Guéritaud, "Deforming ideal solid tori", preprint, 2009. arXiv 0911.3067
- [Hempel 1976] J. Hempel, 3-Manifolds, Ann. of Math. Studies 86, Princeton University Press, 1976. MR 54 #3702 Zbl 0345.57001
- [Hodgson and Rivin 1993] C. D. Hodgson and I. Rivin, "A characterization of compact convex polyhedra in hyperbolic 3-space", *Invent. Math.* 111:1 (1993), 77–111. MR 93j:52015 Zbl 0784.52013
- [Inoue 2008] T. Inoue, "Organizing volumes of right-angled hyperbolic polyhedra", *Algebr. Geom. Topol.* **8**:3 (2008), 1523–1565. MR 2009k:57025 Zbl 1146.52005
- [Jaco 1980] W. Jaco, *Lectures on three-manifold topology*, CBMS Regional Conference Series in Mathematics 43, American Mathematical Society, Providence, R.I., 1980. MR 81k:57009 Zbl 0433.57001
- [Jaco and Shalen 1979] W. H. Jaco and P. B. Shalen, *Seifert fibered spaces in 3-manifolds*, Mem. Amer. Math. Soc. **220**, Amer. Math. Soc., Providence, RI, 1979. MR 81c:57010 Zbl 0415.57005
- [Johannson 1979] K. Johannson, *Homotopy equivalences of 3-manifolds with boundaries*, Lecture Notes in Mathematics **761**, Springer, Berlin, 1979. MR 82c:57005 Zbl 0412.57007
- [Kleiner and Lott 2008] B. Kleiner and J. Lott, "Notes on Perelman's papers", *Geom. Topol.* **12**:5 (2008), 2587–2855. MR 2010h:53098 Zbl 1204.53033
- [Kwun 1970] K. W. Kwun, "Scarcity of orientation-reversing PL involutions of lens spaces", *Michigan Math. J.* 17 (1970), 355–358. MR 43 #5535 Zbl 0191.54901

- [Lü 2009] Z. Lü, "2-torus manifolds, cobordism and small covers", *Pacific J. Math.* 241:2 (2009), 285–308. MR 2010k:55008 Zbl 1181.57036
- [Lü and Masuda 2009] Z. Lü and M. Masuda, "Equivariant classification of 2-torus manifolds", *Collog. Math.* 115:2 (2009), 171–188. MR 2010k:52023 Zbl 1165.57023
- [Lü and Yu 2011] Z. Lü and L. Yu, "Topological types of 3-dimensional small covers", *Forum Math.* **23**:2 (2011), 245–284. MR 2012h:57038 Zbl 1222.52015
- [Meeks and Scott 1986] W. H. Meeks, III and P. Scott, "Finite group actions on 3-manifolds", *Invent. Math.* **86**:2 (1986), 287–346. MR 88b:57039 Zbl 0626.57006
- [Meeks et al. 1982] W. Meeks, III, L. Simon, and S. T. Yau, "Embedded minimal surfaces, exotic spheres, and manifolds with positive Ricci curvature", *Ann. of Math.* (2) **116**:3 (1982), 621–659. MR 84f:53053 Zbl 0521.53007
- [Morgan and Tian 2007] J. Morgan and G. Tian, *Ricci flow and the Poincaré conjecture*, Clay Mathematics Monographs **3**, American Mathematical Society, Providence, RI, 2007. MR 2008d:57020 Zbl 1179.57045
- [Mostow 1973] G. D. Mostow, *Strong rigidity of locally symmetric spaces*, Annals of Mathematics Studies **78**, Princeton University Press, 1973. MR 52 #5874 Zbl 0265.53039
- [Nakayama and Nishimura 2005] H. Nakayama and Y. Nishimura, "The orientability of small covers and coloring simple polytopes", *Osaka J. Math.* 42:1 (2005), 243–256. MR 2006a:57023 Zbl 1065.05041
- [Neumann and Raymond 1978] W. D. Neumann and F. Raymond, "Seifert manifolds, plumbing, μ -invariant and orientation reversing maps", pp. 163–196 in *Algebraic and geometric topology* (Santa Barbara, CA, 1977), edited by K. C. Millett, Lecture Notes in Math. **664**, Springer, Berlin, 1978. MR 80e:57008 Zbl 0401.57018
- [Rivin 1996] I. Rivin, "A characterization of ideal polyhedra in hyperbolic 3-space", *Ann. of Math.* (2) 143:1 (1996), 51–70. MR 96i:52008 Zbl 0874.52006
- [Roeder et al. 2007] R. K. W. Roeder, J. H. Hubbard, and W. D. Dunbar, "Andreev's theorem on hyperbolic polyhedra", *Ann. Inst. Fourier (Grenoble)* **57**:3 (2007), 825–882. MR 2008e:51011 Zbl 1127.51012
- [Schlenker 2002] J.-M. Schlenker, "Hyperbolic manifolds with polyhedral boundary", preprint, 2002. arXiv 0111136
- [Schlenker 2003] J.-M. Schlenker, "Hyperideal polyhedra in hyperbolic manifolds", preprint, 2003. arXiv 0212355
- [Schlenker 2005] J.-M. Schlenker, "Hyperideal circle patterns", *Math. Res. Lett.* **12**:1 (2005), 85–102. MR 2005k:52043 Zbl 1067.52016
- [Schlenker 2006] J.-M. Schlenker, "Hyperbolic manifolds with convex boundary", *Invent. Math.* **163**:1 (2006), 109–169. MR 2006m:57023 Zbl 1091.53019
- [Scott 1983] P. Scott, "The geometries of 3-manifolds", *Bull. London Math. Soc.* **15**:5 (1983), 401–487. MR 84m:57009 Zbl 0561.57001
- [Thurston 1982] W. P. Thurston, "Three-dimensional manifolds, Kleinian groups and hyperbolic geometry", *Bull. Amer. Math. Soc.* (*N.S.*) **6**:3 (1982), 357–381. MR 83h:57019 Zbl 0496.57005
- [Thurston 1986] W. P. Thurston, "A norm for the homology of 3-manifolds", *Mem. Amer. Math. Soc.* **59**:339 (1986), i–vi and 99–130. MR 88h:57014

Received March 31, 2012. Revised November 11, 2012.

YOULIN LI AND JIMING MA

YOULIN LI DEPARTMENT OF MATHEMATICS SHANGHAI JIAOTONG UNIVERSITY 200240 SHANGHAI CHINA liyoulin@sjtu.edu.cn

JIMING MA SCHOOL OF MATHEMATICAL SCIENCES FUDAN UNIVERSITY 200433 SHANGHAI CHINA majiming@fudan.edu.cn
REAL CLOSED SEPARATION THEOREMS AND APPLICATIONS TO GROUP ALGEBRAS

TIM NETZER AND ANDREAS THOM

Dedicated to Konrad Schmüdgen on the occasion of his 65th birthday

In this paper we prove a strong Hahn–Banach theorem: separation of disjoint convex sets by linear forms is possible without any further conditions if the target field \mathbb{R} is replaced by a more general real closed extension field. From this we deduce a general Positivstellensatz for *-algebras, involving representations over real closed fields. We investigate the class of group algebras in more detail. We show that the cone of sums of squares in the augmentation ideal has an interior point if and only if the first cohomology vanishes. For groups with Kazhdan's property (T), the result can be strengthened to interior points in the ℓ^1 -metric. We finally reprove some strong Positivstellensätze by Helton and Schmüdgen, using our separation method.

1. Introduction

In this article we combine techniques from real algebraic geometry, convex geometry, and the unitary representation theory of discrete groups to address various problems that arise in the emerging field of noncommutative real algebraic geometry [Schmüdgen 2009]. Classical results — like Artin's solution of Hilbert's 17th problem — strive for a characterization of natural notions of positivity in terms of algebraic certificates. For example, Artin proved that every polynomial in *n* variables that is positive at every point on \mathbb{R}^n must be a sum of squares of rational functions. Much later, Schmüdgen [1991] proved that a strictly positive polynomial on a compact semialgebraic set must be a sum of squares of polynomials plus defining inequalities. More recently, similar questions were asked in a noncommutative context. Typically, the setup involves a *-algebra *A* and a family of representations \mathcal{F} . The question is now: Is every self-adjoint element of *A* that is positive (semi)definite in every representation in \mathcal{F} necessarily of the form $\sum_i a_i^* a_i$ for some $a_i \in A$? It turned out — similar to the more classical commutative case — that the cone $\Sigma^2 A = \{\sum_i a_i^* a_i \mid a_i \in A\} \subset A$ is an interesting object of study in itself. Natural

MSC2010: 46H15.

Keywords: real closed fields, group rings, Kazhdan's property (T), sums of squares.

questions are: Is $\Sigma^2(A) \cap (-\Sigma^2 A) = \{0\}$? Is $\Sigma^2 A$ closed in a natural topology? Does it contain interior points, in the finest locally convex topology, say?

A point q is called an algebraic interior point of a cone C if the cone intersects each line through q in an open interval around q. A point is an algebraic interior point if and only if it is an interior point in the finest locally convex topology (see [Cimprič et al. 2011, Proposition 1.3], for example). The question for interior points of cones has the following motivation. If a cone C has an (algebraic) interior point q, then for every point $a \in C^{\vee\vee}$ from the double dual, one has $a + \epsilon q \in C$ for all $\epsilon > 0$ (see [loc. cit.] for a proof of this well known fact). Using the standard Gelfand–Naimark–Segal construction, this yields the following Positivstellensatz for unital *-algebras:

Theorem. Assume that q is an interior point of the cone $\Sigma^2 A$. If $a = a^* \in A$ is positive semidefinite in each *-representation of A, then $a + \epsilon q \in \Sigma^2 A$ for all $\epsilon > 0$.

Our first main result is a different Positivstellensatz (Theorem 3.12): we prove that each element from a real reduced unital *-algebra that is positive in every generalized representation is necessarily in $\Sigma^2 A$. The notion of a *generalized representation* involves an extension of the standard real and complex numbers to more general real- and algebraically closed fields.

A natural and vast class of examples of *-algebras is given by complex group algebras $\mathbb{C}[\Gamma]$ of discrete countable groups. We study the cones $\Sigma^2 \mathbb{C}[\Gamma]$ and $\Sigma^2 \omega(\Gamma)$ in more detail, where $\omega(\Gamma) \subset \mathbb{C}[\Gamma]$ denotes the augmentation ideal; see Section 4. The situation for $\omega(\Gamma)$ is much more complicated, as the study is closely related to questions about first cohomology with unitary coefficients. We prove that $\Sigma^2 \omega(\Gamma)$ has an interior point if and only if $H_1(\Gamma, \mathbb{C}) = 0$. The cone $\Sigma^2 \omega(\Gamma)$ has an interior point in the ℓ^1 -metric if Γ has Kazhdan's property (T), and the converse holds if $H_2(\Gamma, \mathbb{C}) = 0$ (see Section 5). In Section 6, we analyze the situation for free groups more closely and reprove theorems of Schmüdgen and Helton.

Along the way we prove some new and powerful separation theorems in Sections 2 and 3. The Hahn–Banach separation theorems for convex sets only apply if additional conditions on the involved sets are imposed; sets have to be closed or have to have nonempty interior, etc. We can remove *all* additional assumptions at the expense of enlarging the target \mathbb{R} to some real closed extension of the real numbers; see Theorem 2.1.

2. A real closed separation theorem for convex sets

Throughout, we will work with various real closed fields **R** and *always* assume that $\mathbb{R} \subset \mathbf{R}$. The following is a first general separation theorem for convex cones.

Theorem 2.1. Let V be an \mathbb{R} -vector space, $C \subset V$ a convex cone and $x \notin C$. Then there exist a real closed field **R** containing \mathbb{R} and an \mathbb{R} -linear functional $\varphi \colon V \to R$

such that

$$\varphi(x) < 0$$
 and $\varphi(y) \ge 0$ for all $y \in C$.

We can even ensure $\varphi(y) > 0$ for all $y \in C \setminus (C \cap -C)$. Also, **R** depends only on V, not on x or C.

Proof. Let us first assume that *V* is finite-dimensional. We construct a complete flag of subspaces $V = H_1 \supset H_2 \supset \cdots \supset H_n = C \cap -C$, starting with $V = H_1$, in the following way. By the standard separation theorem for convex sets (see for example Theorem 2.9 in [Barvinok 2002]), we choose a nontrivial \mathbb{R} -linear functional $\varphi_i : H_i \rightarrow \mathbb{R}$ such that $\varphi_i(y) \ge 0$ for all $y \in C \cap H_i$ and $\varphi_i(x) \le 0$ (if $x \in H_i$). We then define $H_{i+1} := H_i \cap \{\varphi_i = 0\}$ and iterate the process. We finally extend each φ_i in any way to *V*. Now let **R** be a proper real closed extension field of \mathbb{R} . Choose positive elements

$$1 = \epsilon_1 > \epsilon_2 > \cdots > \epsilon_{n-1} > 0$$

from **R** such that $k \cdot \epsilon_i < \epsilon_{i-1}$ for all $k \in \mathbb{R}$. For example, ϵ_2 can be any infinitesimal element with respect to \mathbb{R} , which exists since **R** is a proper extension of \mathbb{R} ; the following ϵ_i can be taken as powers of ϵ_2 .

Then define

$$\varphi := \epsilon_1 \varphi_1 + \dots + \epsilon_{n-1} \varphi_{n-1}.$$

One checks that φ has the desired properties. This proves the claim in the case of finite dimension.

In the general case, consider the set \mathscr{P} of all finite-dimensional subspaces H of V. For each $H \in \mathscr{P}$, choose $\varphi_H : H \to \mathbf{R}$, separating x from $C \cap H$ as desired (if $x \in H$). Extend φ_H in any way to V. Now let ω be an ultrafilter on \mathscr{P} , containing the sets $\{H \in \mathscr{P} \mid y \in H\}$ for all $y \in V$. Consider the linear functional $\varphi : V \to \mathbf{R}^{\omega}$, $\varphi(v) = (\varphi_H(v))_{H \in \mathscr{P}}$. Here \mathbf{R}^{ω} denotes the ultrapower of \mathbf{R} with respect to ω . One checks that φ separates x from C as desired, by the theorem of Los (see for example Theorem 2.2.9 in [Prestel and Delzell 2001]).

Remark 2.2. In the usual way, one can now also deduce that any two convex disjoint sets in a vector space can be separated as above with a real-closed valued affine functional.

It turns out that we can also extend functionals quite often if we allow for an extension of the real closed field.

Theorem 2.3. Let V be an \mathbb{R} -vector space, $C \subseteq V$ a convex cone, and $H \subseteq V$ a subspace. Assume $(C + H) \cap -(C + H) = H$. Then for any real closed extension field **R** of \mathbb{R} and any \mathbb{R} -linear functional $\varphi \colon H \to \mathbf{R}$ with $\varphi \ge 0$ on $C \cap H$, there is a real closed extension field \mathbf{R}' of \mathbf{R} and an \mathbb{R} -linear functional $\overline{\varphi} \colon V \to \mathbf{R}'$ with $\overline{\varphi} \ge 0$ on C and $\overline{\varphi} = \varphi$ on H. We can even ensure $\overline{\varphi}(y) > 0$ for all $y \in C \setminus H$.

Proof. We apply Theorem 2.1 to the convex cone C + H in V and obtain a real closed field $\widetilde{\mathbf{R}}$ and an \mathbb{R} -linear functional $\psi: V \to \widetilde{\mathbf{R}}$ with $\psi = 0$ on H and $\psi(y) > 0$ for $a \in C \setminus H$. By amalgamation of real closed fields, we can assume without loss of generality that $\mathbf{R} = \widetilde{\mathbf{R}}$. Finally, let \mathbf{R}' be a real closed extension field of \mathbf{R} that contains an element $\delta > \mathbf{R}$. Extend φ to an \mathbf{R} -valued functional on V and set

$$\overline{\varphi} := \varphi + \delta \cdot \psi.$$

It is clear that $\overline{\varphi}$ coincides with φ on H and also that $\overline{\varphi}(y) > 0$ for all $y \in C \setminus H$. \Box

We will improve upon the separation results in the case of certain *-algebras in the next section.

3. Completely positive separation

Throughout this section, let A be a \mathbb{C} -algebra with involution *, not necessarily unital. We consider the cone of sums of hermitian squares

$$\Sigma^2 A = \left\{ \sum_{i=1}^n a_i^* a_i \mid n \in \mathbb{N}, a_i \in A \right\}$$

contained in the real vector subspace of hermitian elements

$$A^{h} = \{ a \in A \mid a^{*} = a \}.$$

If $b \in A^h \setminus \Sigma^2 A$, we find an \mathbb{R} -linear functional $\varphi \colon A^h \to \mathbf{R}$, into some real closed extension field \mathbf{R} of \mathbb{R} , such that $\varphi(b) < 0$, $\varphi(a^*a) \ge 0$ for all $a \in A$, by Theorem 2.1. We can extend φ uniquely to a \mathbb{C} -linear functional $\varphi \colon A \to \mathbf{R}[i]$ fulfilling $\varphi(a^*) = \overline{\varphi(a)}$. We will denote the algebraically closed field $\mathbf{R}[i]$ by \mathbf{C} from now on.

The condition $\varphi(a^*a) \ge 0$ for all $a \in A$ is called *positivity* of φ . We would now like positive and real-closed valued functionals to fulfill the Cauchy–Schwarz inequality

$$|\varphi(a^*b)|^2 \le \varphi(a^*a)\varphi(b^*b)$$

for all $a, b \in A$. However, this is not true in general, as the next example shows.

Example 3.1. Let $A = \mathbb{C}[t]$ be the univariate polynomial ring; * is coefficientwise conjugation. The cone $\Sigma^2 A$ equals the cone of nonnegative real polynomials. Consider the functional

$$\varphi \colon \mathbb{R}[t] \to \mathbf{R}, \quad p \mapsto p(0) + \epsilon p''(0),$$

where $\epsilon \in \mathbf{R}$ is positive and infinitesimal with respect to \mathbb{R} . One checks that φ is positive, but for $a = 1 + t^2$ and b = 1, we have

$$|\varphi(a^*b)|^2 = 1 + 4\epsilon + 4\epsilon^2 > 1 + 4\epsilon = \varphi(a^*a)\varphi(b^*b).$$

Example 3.2. The last example can be modified to even fulfill $\varphi(a^*a) > 0$ if $a \neq 0$. Indeed, let $1 = \epsilon_0 > \epsilon_1 > \epsilon_2 > \cdots > 0$ be a sequence of elements from **R**, such that $\mathbb{R} \cdot \epsilon_i < \epsilon_{i-1}$ for all *i*. Then the linear mapping $p \mapsto \sum_{i=0}^{\infty} \epsilon_i \cdot p^{(2i)}(0)$ is well-defined and strictly positive in the desired sense. If we further assume $\mathbb{R} \cdot \epsilon_2 < \epsilon_1^2$, then the same argument as in Example 3.1 shows that the Cauchy–Schwarz inequality is not fulfilled.

Definition 3.3. A \mathbb{C} -linear functional $\varphi \colon A \to \mathbf{C}$ with $\varphi(a^*) = \overline{\varphi(a)}$ is called *completely positive* if for all $m \in \mathbb{N}$, the componentwise defined function

$$\varphi^{(m)}: \mathbf{M}_m(A) \to \mathbf{M}_m(\mathbf{C})$$

maps sums of hermitian squares to positive semidefinite matrices.

Remark 3.4. It is easily seen that a positive \mathbb{C} -linear functional $\varphi \colon A \to \mathbb{C}$ with $\varphi(a^*) = \overline{\varphi(a)}$ is always completely positive.

Example 3.5. The functionals from Examples 3.1 and 3.2 are positive, but not completely positive. Indeed, with $a = 1 + t^2$ and

$$M = \left(\begin{array}{cc} 1 & a \\ 0 & 0 \end{array}\right),$$

we find that

$$\varphi^{(2)}(M^*M) = \left(\begin{array}{cc} 1 & \varphi(a) \\ \varphi(a^*) & \varphi(a^*a) \end{array}\right)$$

is not positive semidefinite, since its determinant is negative in **R**.

Lemma 3.6. A \mathbb{C} -linear functional $\varphi \colon A \to \mathbf{C}$ with $\varphi(a^*) = \overline{\varphi(a)}$ is completely positive if and only if the \mathbf{C} -linear extension

$$\mathrm{id}\otimes\varphi\colon\mathbf{C}\otimes_{\mathbb{C}}A\to\mathbf{C}$$

is positive.

Proof. The condition that $id \otimes \varphi$ is positive is

$$0 \le (\mathrm{id} \otimes \varphi) \left(\left(\sum_{j=1}^m z_j \otimes a_j \right)^* \left(\sum_{j=1}^m z_j \otimes a_j \right) \right)$$
$$= (\mathrm{id} \otimes \varphi) \left(\sum_{j,k} \overline{z}_j z_k \otimes a_j^* a_k \right)$$
$$= \sum_{j,k} \overline{z}_j z_k \cdot \varphi(a_j^* a_k)$$

for all $m \in \mathbb{N}$, $z_j \in \mathbb{C}$, $a_j \in A$. But this just means that the matrix $(a_j^*a_k)_{j,k}$ is mapped to a positive semidefinite matrix under $\varphi^{(m)}$. Since every sum of hermitian squares in $M_m(A)$ is a finite sum of such rank-one squares, this proves the claim. \Box

Corollary 3.7. If $\varphi \colon A \to \mathbb{C}$ is completely positive, then it fulfills the Cauchy– Schwarz inequality

$$|\varphi(a^*b)|^2 \le \varphi(a^*a)\varphi(b^*b)$$

for all $a, b \in A$.

Proof. Either consider the positive and **C**-linear extension id $\otimes \varphi$ to $\mathbf{C} \otimes_{\mathbb{C}} A$ and use the standard proof for the inequality, or apply $\varphi^{(2)}$ to the sum of hermitian squares

$$\begin{pmatrix} a & b \\ 0 & 0 \end{pmatrix}^* \begin{pmatrix} a & b \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} a^*a & a^*b \\ b^*a & b^*b \end{pmatrix},$$

and use that the obtained matrix is positive semidefinite.

Remark 3.8. We see from the last proof that in fact only the 2-positivity of φ is needed for the Cauchy–Schwarz inequality.

Corollary 3.9. Let A be a \mathbb{C} -algebra with involution and \mathbb{R} a real closed field that contains \mathbb{R} . Let $\varphi \colon A \to \mathbb{C}$ be a completely positive \mathbb{C} -linear functional that satisfies $\varphi(a^*) = \overline{\varphi(a)}$ for all $a \in A$. The gauge $||a||_{\varphi} := \varphi(a^*a)^{1/2}$ satisfies

$$\|\lambda \cdot a\|_{\varphi} = |\lambda| \cdot \|a\|_{\varphi}$$

and

$$\|a+b\|_{\varphi} \leq \cdot \|a\|_{\varphi} + \|b\|_{\varphi}.$$

Proof. The first assertion is obvious. Let's compute

$$\begin{aligned} \|a+b\|_{\varphi}^{2} &= \varphi((a+b)^{*}(a+b)) \\ &= \varphi(a^{*}a) + \varphi(a^{*}b) + \varphi(b^{*}a) + \varphi(b^{*}b) \\ &\leq \|a\|_{\varphi}^{2} + \|b\|_{\varphi}^{2} + 2|\varphi(a^{*}b)| \\ &\leq \|a\|_{\varphi}^{2} + \|b\|_{\varphi}^{2} + 2\|a\|_{\varphi}\|b\|_{\varphi} \\ &= (\|a\|_{\varphi} + \|b\|_{\varphi})^{2}. \end{aligned}$$

This proves the claim.

It turns out that separation from the cone of sums of hermitian squares can often be done with a completely positive functional.

Definition 3.10. Let *A* be a \mathbb{C} -algebra with involution, not necessarily unital. Then *A* is called *real reduced* if $\sum_{i} a_i^* a_i = 0$ implies $a_i = 0$ for all *i* and $a_i \in A$.

Theorem 3.11. Let A be a \mathbb{C} -algebra with involution that is real reduced. Let $b \in A^h \setminus \Sigma^2 A$. Then there is real closed extension field **R** of \mathbb{R} and a completely positive \mathbb{C} -linear functional $\varphi \colon A \to \mathbb{C}$ with $\varphi(a^*) = \overline{\varphi(a)}$ such that

$$\varphi(b) < 0$$
 and $\varphi(a^*a) > 0$ for $a \in A \setminus \{0\}$.

 \Box

Proof. For any finite-dimensional subspace H of A, denote by

$$\Sigma^2 H = \left\{ \sum_i a_i^* a_i \mid a_i \in H \right\}$$

the set of sums of hermitian squares of elements from *H*. It is well known that $\Sigma^2 H$ is a closed convex cone in a finite-dimensional subspace of A^h . This follows from the fact that *A* is real reduced, using for example the approach from [Powers and Scheiderer 2001, Lemma 2.7]. It also follows that $\Sigma^2 H$ is salient, that is, it fulfills

$$\Sigma^2 H \cap -\Sigma^2 H = \{0\}.$$

So for each such H, there is an \mathbb{R} -linear functional $\varphi_H \colon A^h \to \mathbb{R}$ with

$$\varphi_H(b) < 0$$
 and $\varphi_H(a^*a) > 0$ for all $a \in H \setminus \{0\}$

Let \mathcal{G} be the set of all finite-dimensional subspaces H of A, equipped with an ultrafilter ω containing all sets $\{H \in \mathcal{G} \mid c \in H\}$ for $c \in A$. Define

$$\varphi \colon A^h \to \mathbb{R}^{\omega}, \quad \varphi(a) \coloneqq (\varphi_H(a))_{H \in \mathcal{G}}.$$

Then φ does the separation as desired. We consider the \mathbb{C} -linear extension of φ to A, and finally show that it is completely positive. The \mathbb{C} -linear extension of φ_H to A indeed maps a matrix $(a_i^*a_j)_{i,j} \in M_n(A)$ to a positive semidefinite hermitian matrix, at least if all $a_i \in H$, as is easily checked (compare to Remark 3.4). Since we can check positivity of the matrix $(\varphi(a_i^*a_j))_{i,j} \in M_n(\mathbb{R}^{\omega}[i])$ componentwisely in $M_n(\mathbb{R}[i])$, by the theorem of Łos, this finishes the proof.

Throughout, we will take the freedom to consider *-representations of *A* on vector spaces that carry a sesquilinear C-valued inner product, where $\mathbf{C} = \mathbf{R}[i]$ for some real closed field $\mathbf{R} \supset \mathbb{R}$. We call these representations *generalized representations*. For every completely positive functional $\varphi : A \rightarrow \mathbf{C}$, we can perform the usual GNS construction to construct such a representation (see the proof of Theorem 6.1 below for more technical details). The usual concepts of self-adjointness and positive semidefiniteness of operators on such a vector space can be defined without any problems. The first consequence is the following Positivstellensatz (compare to the standard Positivstellensatz from the introduction):

Theorem 3.12. Let A be a real reduced *-algebra, and $a \in A^h$. Then a is positive semidefinite in every generalized representation if and only if $a \in \Sigma^2 A$.

Proof. If $a \notin \Sigma^2 A$, then there exists a completely positive map $\varphi \colon A \to \mathbb{C}$ such that $\varphi(a) < 0$. Clearly, *a* will not be positive semidefinite in the generalized GNS representation associated with φ .

Examples for real reduced *-algebras are group algebras $\mathbb{C}[\Gamma]$. In Section 6, we will see that for particular groups, the study of generalized representations of $\mathbb{C}[\Gamma]$ can be reduced to the study of usual (finite-dimensional) unitary representations, using Tarski's transfer principle.

4. Sums of squares in the group algebra

Let Γ be a group and let $\mathbb{C}[\Gamma]$ denote the complex group algebra. A typical element in $\mathbb{C}[\Gamma]$ is denoted by $a = \sum_g a_g g$, where only finitely many of the $a_g \in \mathbb{C}$ are not zero. In $\mathbb{C}[\Gamma]$ we identify \mathbb{C} with $\mathbb{C} \cdot e$, where *e* denotes the neutral element of Γ . The group algebra comes equipped with an involution $(\sum_g a_g g)^* := \sum_g \bar{a}_g g^{-1}$ and a trace $\tau : \mathbb{C}[\Gamma] \to \mathbb{C}$ that is given by the formula $\tau (\sum_g a_g g) = a_e$. The faithfulness of the trace shows that $\mathbb{C}[\Gamma]$ is real reduced. Let $\Sigma^2 \mathbb{C}[\Gamma]$ denote the set of sums of hermitian squares in $\mathbb{C}[\Gamma]$. The following appears for example as Example 3 in [Cimprič 2009]:

Remark 4.1. $||a||_1^2 - a^*a \in \Sigma^2 \mathbb{C}[\Gamma]$ for all $a \in \mathbb{C}[\Gamma]$, where $||a||_1 = \sum_g |a_g|$.

Remark 4.2. From the identity

$$2\|a\|_{1} \cdot (\|a\|_{1} - a) = (\|a\|_{1} - a)^{*} (\|a\|_{1} - a) + (\|a\|_{1}^{2} - a^{*}a)$$

for $a \in \mathbb{C}[\Gamma]^h$, we see that 1 is an algebraic interior point of the cone $\Sigma^2 \mathbb{C}[\Gamma]$ in the real vector space $\mathbb{C}[\Gamma]^h$. That means $1 + \epsilon a \in \Sigma^2 \mathbb{C}[\Gamma]$ for all $a \in \mathbb{C}[\Gamma]^h$ and sufficiently small $\epsilon > 0$. In fact, the ϵ does only depend on $||a||_1$ here.

Remark 4.3. As explained in the introduction, for any element $a \in \mathbb{C}[\Gamma]^h$ that is positive semidefinite in each (usual) *-representation of $\mathbb{C}[\Gamma]$, one thus has $a + \epsilon \in \Sigma^2 \mathbb{C}[\Gamma]$, for all $\epsilon > 0$.

Remark 4.4. Since $\mathbb{C}[\Gamma]$ is real reduced and unital, the result of Theorem 3.12 holds here as well. So if *a* is positive semidefinite in each generalized representation, then $a \in \Sigma^2 \mathbb{C}[\Gamma]$.

We now consider the augmentation homomorphism $\varepsilon \colon \mathbb{C}[\Gamma] \to \mathbb{C}$, which is defined by $\varepsilon(\sum_{g} a_{g}g) = \sum_{g} a_{g}$. The *augmentation ideal* is

$$\omega(\Gamma) := \ker(\varepsilon) = \left\{ a \in \mathbb{C}[\Gamma] \mid \sum_{g} a_{g} = 0 \right\}.$$

We set c(g) := g - 1 and note that $\{c(g) \mid g \in \Gamma \setminus \{e\}\}$ is a basis of $\omega(\Gamma)$. The multiplication satisfies

$$c(g)c(h) = c(gh) - c(g) - c(h).$$

We denote by $\omega^2(\Gamma)$ the square of $\omega(\Gamma)$, that is, $\omega^2(\Gamma) = \operatorname{span}_{\mathbb{C}} \{ab \mid a, b \in \omega(\Gamma)\}$. Inside $\omega(\Gamma)$, we study the cone of sums of hermitian squares

$$\Sigma^2 \omega(\Gamma) := \left\{ \sum a_i^* a_i \mid a_i \in \omega(\Gamma) \right\}.$$

We are interested in interior points of this cone. Note that $\Sigma^2 \omega(\Gamma) \subset \omega^2(\Gamma)$ and $\omega(\Gamma)/\omega^2(\Gamma) = \mathbb{C} \otimes_{\mathbb{Z}} \Gamma_{ab}$, where $\Gamma_{ab} = \Gamma/[\Gamma, \Gamma]$ and $[\Gamma, \Gamma]$ denotes the subgroup of Γ generated by commutators. Hence, if Γ has nontorsion abelianization, then $\Sigma^2 \omega(\Gamma)$ is contained in a proper subspace of $\omega(\Gamma)$. However, we will show below that $\Sigma^2 \omega(\Gamma)$ always has an interior point in $\omega^2(\Gamma)^h$.

Lemma 4.5. For any group Γ , we have $\Sigma^2 \omega(\Gamma) = \Sigma^2 \mathbb{C}[\Gamma] \cap \omega(\Gamma)$.

Proof. The inclusion $\Sigma^2 \omega(\Gamma) \subset \Sigma^2 \mathbb{C}[\Gamma] \cap \omega(\Gamma)$ is obvious. If $\sum_i a_i^* a_i \in \omega(\Gamma)$ with $a_i \in \mathbb{C}[\Gamma]$, then $\sum_i |\varepsilon(a_i)|^2 = 0$, and hence $\varepsilon(a_i) = 0$ for all *i*. This proves the converse inclusion.

Remark 4.6. It turns out that we can always extend positive functionals φ on $\omega(\Gamma)$ to positive functionals $\overline{\varphi}$ on $\mathbb{C}[\Gamma]$, at least if we allow for an extension of the real closed field. Indeed, observe that

$$\left(\Sigma^2 \mathbb{C}[\Gamma] + \omega(\Gamma)^h\right) \cap -\left(\Sigma^2 \mathbb{C}[\Gamma] + \omega(\Gamma)^h\right) = \omega(\Gamma)^h,$$

which follows immediately from an application of the augmentation homomorphism ε . We can thus apply Theorem 2.3.

Lemma 4.7. Let **R** be a real closed extension field of \mathbb{R} , and $\varphi \colon \omega(\Gamma) \to \mathbf{C}$ a completely positive \mathbb{C} -linear functional with $\varphi(a^*) = \overline{\varphi(a)}$ for all $a \in \omega(\Gamma)$. Then for all $s, h \in \Gamma$,

$$\left|\varphi(c(s)^*c(h))\right| \le \frac{1}{\sqrt{2}} \cdot \left(\varphi(c(s)^*c(s)) + \varphi(c(h)^*c(h))\right),$$

$$\varphi(c(sh)^*c(sh)) \le 2 \cdot \left(\varphi(c(s)^*c(s)) + \varphi(c(h)^*c(h))\right).$$

Proof. The first inequality is an application of the Cauchy–Schwarz inequality (which is fulfilled by completely positive functionals) and the inequality $\lambda \mu \leq (\lambda + \mu)^2/2$. For the second inequality, first apply the triangle identity from Corollary 3.9 to the equation

$$c(sh) = c(s) + (c(s)c(h) + c(h)),$$

together with the well-known inequality $(a+b)^2 \le 2(a^2+b^2)$. Then use the easily verified identity $||c(s)c(h) + c(h)||^2 = ||c(h)||^2$.

Since $1 \notin \omega(\Gamma)$, we need to find a different candidate for an interior point of $\Sigma^2 \omega(\Gamma)$.

Definition 4.8. Let $S \subset \Gamma$ be a finite symmetric set, that is, $S^{-1} = S$. We define the Laplace operator on *S* to be

$$\Delta(S) := |S| - \sum_{s \in S} s.$$

Remark 4.9. Note that for every finite symmetric set $S \subset \Gamma$,

$$\Delta(S) = \frac{1}{2} \cdot \sum_{s \in S} c(s)^* c(s) \in \Sigma^2 \omega(\Gamma).$$

Proposition 4.10. Let Γ be a group generated by a finite symmetric set S. Then for any $b \in \omega^2(\Gamma)$, there exists a constant $C(b) \in \mathbb{R}$ such that for any real closed extension field \mathbf{R} of \mathbb{R} and any completely positive \mathbb{C} -linear functional $\varphi : \omega(\Gamma) \to \mathbf{C}$ with $\varphi(a^*) = \overline{\varphi(a)}$ for all $a \in \omega(\Gamma)$, one has

$$|\varphi(b)| \le C(b) \cdot \varphi(\Delta(S)).$$

Proof. Every element $b \in \omega^2(\Gamma)$ is a finite linear combination of $c(g)^*c(h)$ for $g, h \in \Gamma \setminus \{e\}$. This implies that $|\varphi(b)|$ is bounded by a constant times

$$\max\{\varphi(c(s)^*c(s)) \mid s \in S\},\$$

using Lemma 4.7 several times. However,

$$\max\left\{\varphi(c(s)^*c(s)) \mid s \in S\right\} \le 2 \cdot \varphi(\Delta(S))$$

follows from Remark 4.9. This proves the claim.

Theorem 4.11. Let Γ be a group with finite generating symmetric set S. Then for every $b \in \omega^2(\Gamma)^h$, there is a constant $C(b) \in \mathbb{R}$ such that

$$C(b) \cdot \Delta(S) \pm b \in \Sigma^2 \omega(\Gamma).$$

In particular, $\Delta(S)$ is an inner point of the cone $\Sigma^2 \omega(\Gamma)$ in $\omega^2(\Gamma)^h$. If

$$H_1(\Gamma, \mathbb{C}) = 0,$$

it is an inner point in $\omega(\Gamma)^h$.

Proof. In view of Proposition 4.10, we find that

$$C(b) \cdot \Delta(S) \pm b$$

is nonnegative under each completely positive real-closed valued \mathbb{C} -linear functional φ on $\omega(\Gamma)$. In view of Theorem 3.12, this means that $C(b) \cdot \Delta(S) \pm b$ is a sum of hermitian squares in $\omega(\Gamma)$. Note that we use that $\omega(\Gamma)$ is real reduced. Finally note that $H_1(\Gamma, \mathbb{C}) = 0$ just means that $\omega^2(\Gamma) = \omega(\Gamma)$.

5. Groups with Kazhdan's property (T)

We want to show that the constant C(b) from Theorem 4.11 can be chosen as a fixed multiple of $||b||_1$, in case the group Γ has Kazhdan's property.

Definition 5.1. Let Γ be a group and $\pi \colon \Gamma \to U(H)$ a unitary representation on a Hilbert space *H*.

- (1) A 1-cocycle with respect to the unitary representation π is a map $\delta \colon \Gamma \to H$ such that for all $g, h \in \Gamma$, we have $\delta(gh) = \pi(g)\delta(h) + \delta(g)$.
- (2) A 1-cocycle $\delta \colon \Gamma \to H$ is called inner if $\delta(g) = \pi(g)\xi \xi$ for some vector $\xi \in H$.

Definition 5.2. A group has Kazhdan's property (T) if every 1-cocycle with respect to every unitary representation is inner.

We will use several results on Kazhdan groups, which can be found, for example, in [Bekka et al. 2008]. It is well known that groups with Kazhdan's property (T) admit a finite generating set S, and that

$$\omega^2(\Gamma) = \omega(\Gamma)$$

holds. It is also known that for a fixed finite symmetric and generating set *S* in a Kazhdan group Γ , there is some $\epsilon > 0$ such that for any unitary representation $\pi \colon \Gamma \to U(H)$ without nonzero fixed vectors, one has

$$\langle \Delta(S)\xi,\xi\rangle \ge \epsilon \cdot \|\xi\|^2$$
 for all $\xi \in H$.

Such ϵ is called a *Kazhdan constant* for *S*.

Let's revisit the standard GNS representation in the context of $\omega(\Gamma)$. Let

$$\varphi \colon \omega(\Gamma) \to \mathbb{C}$$

be a positive linear functional with $\varphi(a^*) = \overline{\varphi(a)}$. We associate to φ a Hilbert space as follows. We define on $\omega(\Gamma)$ a positive semidefinite sesquilinear form

$$\langle a, b \rangle_{\varphi} := \varphi(b^*a)$$

and set $||a||_{\varphi} := \langle a, a \rangle_{\varphi}^{1/2}$. Let $N(\varphi) := \{a \in \omega(\Gamma) \mid ||a||_{\varphi} = 0\}$ and define $L^2(\omega(\Gamma), \varphi)$ to be the metric completion of $\omega(\Gamma)/N(\varphi)$ with respect to $|| \cdot ||_{\varphi}$. We denote the image of c(g) in $L^2(\omega(\Gamma), \varphi)$ by $\delta(g)$ and denote by $\delta(\Gamma)$ their complex linear span, which is dense by definition of $L^2(\omega(\Gamma), \varphi)$). It is standard that the left-multiplication of $\omega(\Gamma)$ on itself extends to a homomorphism $\pi^{\varphi} : \omega(\Gamma) \to D(\delta(\Gamma))$, where $D(\delta(\Gamma))$ denotes the algebra of densely defined linear operators mapping $\delta(\Gamma)$ into itself. Indeed, if $a \in N(\varphi)$ and $b \in \omega(\Gamma)$, then $ba \in N(\varphi)$ since

$$\varphi((ba)^*ba) = \varphi(a^*b^*ba) \le \|b\|_1^2 \cdot \varphi(a^*a),$$

by Remark 4.1. Note that

$$\pi^{\varphi}(c(g))\delta(h) = \delta(gh) - \delta(g) - \delta(h).$$

Now we define a unitary representation π_{φ} of Γ on $L^{2}(\omega(\Gamma), \varphi)$ by the rule

$$\pi_{\varphi}(g) := \pi^{\varphi}(c(g)) + \mathbf{1}_{L^{2}(\omega(\Gamma),\varphi)}.$$

Lemma 5.3. If $\omega^2(\Gamma) = \omega(\Gamma)$, then the representation π_{φ} has no fixed vectors.

Proof. Assume $\eta \in L^2(\omega(\Gamma), \varphi)$ is a fixed vector. By definition of π_{φ} , this means $\pi^{\varphi}(c(g))\eta = 0$ for all $g \in \Gamma$. Hence,

$$0 = \left\langle \pi^{\varphi}(c(g^{-1}))\eta, \delta(h) \right\rangle_{\varphi} = \langle \eta, \delta(gh) - \delta(g) - \delta(h) \rangle_{\varphi}.$$

Since c(g)c(h) = c(gh) - c(g) - c(h), the vectors $\delta(gh) - \delta(g) - \delta(h)$ span the image of $\omega^2(\Gamma)$ in $\delta(\Gamma)$, and hence $\delta(\Gamma)$, since $\omega^2(\Gamma) = \omega(\Gamma)$.

Note that the map $g \mapsto \delta(g)$ satisfies

$$\delta(gh) = \pi_{\varphi}(g)\delta(h) + \delta(g),$$

and hence defines a 1-cocycle with respect to the representation π_{φ} . If Γ is a Kazhdan group, then there exists $\Omega \in L^2(\omega(\Gamma), \varphi)$ such that

$$\delta(g) = \pi_{\varphi}(g)\Omega - \Omega.$$

Proposition 5.4. Let Γ be a Kazhdan group with finite symmetric generating set *S* and Kazhdan constant $\epsilon > 0$. Then for every nonzero $b \in \omega(\Gamma)^h$, every real closed extension field **R** of \mathbb{R} , and every positive nontrivial \mathbb{C} -linear functional $\varphi \colon \omega(\Gamma) \to \mathbf{C}$ with $\varphi(a^*) = \overline{\varphi(a)}$, one has

$$\epsilon \cdot \varphi(b) < 2 \|b\|_1 \cdot \varphi(\Delta(S)).$$

Proof. Let us first assume $\mathbf{R} = \mathbb{R}$ and $\mathbf{C} = \mathbb{C}$. We do the GNS construction as just described, and get some $\Omega \in L^2(\omega(\Gamma), \varphi)$ with $\delta(g) = \pi_{\varphi}(g)\Omega - \Omega$, for all $g \in \Gamma$. We set

$$\overline{\varphi} \colon \mathbb{C}[\Gamma] \to \mathbb{C}, \quad \overline{\varphi}(a) = \langle \pi_{\varphi}(a)\Omega, \Omega \rangle$$

and compute

$$\bar{\varphi}(c(h)^*c(g)) = \left\langle \pi_{\varphi}(c(g))\Omega, \, \pi_{\varphi}(c(h))\Omega \right\rangle_{\varphi} = \langle \delta(g), \, \delta(h) \rangle_{\varphi} = \varphi(c(h)^*c(g)).$$

This shows that $\overline{\varphi}$ and φ agree on $\omega^2(\Gamma)$, and hence on $\omega(\Gamma)$. If we now do the standard GNS construction with respect to $\overline{\varphi}$, we see that there is a natural Γ -equivariant identification of $L^2(\mathbb{C}[\Gamma], \overline{\varphi})$ and $L^2(\omega(\Gamma), \varphi)$. Since the representation π_{φ} has no fixed vectors, we get

$$\varphi(\Delta(S)) = \overline{\varphi}(\Delta(S)) = \langle \Delta(S)1, 1 \rangle_{\overline{\varphi}} \ge \epsilon \cdot \overline{\varphi}(1).$$

Since $\overline{\varphi}$ is positive and nontrivial, it follows from Remark 4.2 that $\overline{\varphi}(1) > 0$. So finally, again using Remark 4.2, we find

$$\epsilon \cdot \varphi(b) = \epsilon \cdot \overline{\varphi}(b) \le \epsilon \cdot \|b\|_1 \cdot \overline{\varphi}(1) < 2\|b\|_1 \cdot \varphi(\Delta(S)),$$

the desired result.

Now let **R** be arbitrary, and let $\varphi: \omega(\Gamma) \to \mathbf{C}$ be positive and nontrivial. From Theorem 4.11 it follows that $\varphi(\Delta(S)) > 0$. So we can assume without loss of generality that $\varphi(\Delta(S)) = 1$. Again from Theorem 4.11, we see that φ now only takes values in $\mathbb{O}[i]$, where \mathbb{O} is the convex hull of \mathbb{R} in **R**. It is well known that \mathbb{O} is a valuation ring in **R** with maximal ideal m, and that $\mathbb{O}/\mathbb{m} = \mathbb{R}$ (see for example [Prestel and Delzell 2001], especially the appendix on valued fields). The residue map $\mathbb{O} \to \mathbb{O}/\mathbb{m}$ maps nonnegative elements to nonnegative elements. So if we compose φ with the residue map on $\mathbb{O}[i]$, we get a positive linear functional to \mathbb{C} . Since we know that the desired strict inequality holds now, it was already valid for φ .

Theorem 5.5. Let Γ be a group with finite generating symmetric set S. Consider the statements:

- (1) Γ has Kazhdan's property (T).
- (2) $\Delta(S)$ is an algebraic interior point of $\Sigma^2 \omega(\Gamma)$ in the ℓ^1 -metric of $\omega(\Gamma)^h$. More precisely, there exists a constant $\epsilon > 0$ such that for every $b \in \omega(\Gamma)^h$ with $\|b\|_1 = 1$, we have

$$\Delta(S) + \epsilon \cdot b \in \Sigma^2 \omega(\Gamma).$$

The following implications hold: (1) implies (2), and (2) implies (1) under the additional assumption $H_2(\Gamma, \mathbb{C}) = 0$.

Proof. The implication $(1) \Rightarrow (2)$ is a direct consequence of Theorem 3.11 and Proposition 5.4. Let us now prove $(2) \Rightarrow (1)$ under the additional assumption $H_2(\Gamma, \mathbb{C}) = 0$. We first prove two lemmas.

Lemma 5.6. Let Γ be a group. There is an exact sequence as follows:

 $0 \to H_2(\Gamma, \mathbb{C}) \to \omega(\Gamma) \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma) \to \omega(\Gamma) \to H_1(\Gamma, \mathbb{C}) \to 0.$

Proof. It is well known that

$$\omega(\Gamma)/\omega^2(\Gamma) = \Gamma_{ab} \otimes_{\mathbb{Z}} \mathbb{C} = H_1(\Gamma, \mathbb{C}).$$

This shows exactness at $\omega(\Gamma)$, and it remains to show exactness at $\omega(\Gamma) \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma)$. Since $0 \to \omega(\Gamma) \to \mathbb{C}[\Gamma] \to \mathbb{C} \to 0$ is exact, we have

$$H_2(\Gamma, \mathbb{C}) = H_1(\Gamma, \omega(\Gamma)) = \operatorname{Tor}_1^{\mathbb{C}[\Gamma]}(\mathbb{C}, \omega(\Gamma)).$$

Again, we get an exact sequence

$$\operatorname{Tor}_{1}^{\mathbb{C}[\Gamma]}(\mathbb{C}[\Gamma], \omega(\Gamma)) \to \operatorname{Tor}_{1}^{\mathbb{C}[\Gamma]}(\mathbb{C}, \omega(\Gamma)) \to \omega(\Gamma) \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma) \to \mathbb{C}[\Gamma] \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma).$$

This finishes the proof, since $\operatorname{Tor}_{1}^{\mathbb{C}[\Gamma]}(\mathbb{C}[\Gamma], \omega(\Gamma)) = 0.$

If $\Delta(S)$ is an algebraic interior point of $\Sigma^2 \omega(\Gamma)$ in $\omega(\Gamma)^h$, then $\omega(\Gamma) = \omega^2(\Gamma)$, that is, $H_1(\Gamma, \mathbb{C}) = 0$. Hence $H_2(\Gamma, \mathbb{C}) = 0$ ensures that the natural map

 $\omega(\Gamma) \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma) \to \omega(\Gamma)$

is an isomorphism. This is what we are going to use.

Lemma 5.7. Let $\pi : \Gamma \to U(H)$ be a unitary representation and let $\delta : \Gamma \to H$ be a 1-cocycle with respect to H. Then

$$\varphi(c(h)^* \otimes c(g)) := \langle \delta(g), \delta(h) \rangle$$

yields a well-defined positive linear functional on $\omega(\Gamma) = \omega(\Gamma) \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma)$.

Proof. It is clear that $(c(h)^*, c(g)) \mapsto \langle \delta(g), \delta(h) \rangle$ defines a bilinear map on $\omega(\Gamma)$, that is, a linear map $\varphi' \colon \omega(\Gamma) \otimes_{\mathbb{C}} \omega(\Gamma) \to \mathbb{C}$. We show that this map passes to $\omega(\Gamma) \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma)$. Let $g, h, k \in \Gamma$; then

$$\begin{split} \varphi'(c(h)^*k \otimes c(g)) &= \varphi'(c(h^{-1})k \otimes c(g)) \\ &= \varphi'\big((c(h^{-1}k) - c(k)) \otimes c(g)\big) \\ &= \varphi'(c(k^{-1}h)^* \otimes c(g)) - \varphi'(c(k^{-1})^* \otimes c(g)) \\ &= \langle \delta(g), \delta(k^{-1}h) \rangle - \langle \delta(g), \delta(k^{-1}) \rangle \\ &= \langle \delta(g), \pi(k^{-1})\delta(h) \rangle \\ &= \langle \pi(k)\delta(g), \delta(h) \rangle \\ &= \langle \delta(kg), \delta(h) \rangle - \langle \delta(k), \delta(h) \rangle \\ &= \varphi'(c(h)^* \otimes c(kg)) - \varphi'(c(h)^* \otimes c(k)) \\ &= \varphi'(c(h^*) \otimes kc(g)). \end{split}$$

We can now understand φ' as a linear map on $\omega(\Gamma)$ via the above isomorphism to $\omega(\Gamma) \otimes_{\mathbb{C}[\Gamma]} \omega(\Gamma)$. Since a^*a corresponds to $a^* \otimes a$, one easily checks that φ is positive on $\omega(\Gamma)$.

We continue with the proof of Theorem 5.5. Condition (2) in Theorem 5.5 and Lemma 5.7 imply that any 1-cocycle with respect to any unitary representation is bounded. This is well known to imply Kazhdan's property (T) for Γ .

Remark 5.8. It is not clear whether the condition $H_2(\Gamma, \mathbb{C}) = 0$ is necessary.

There is an analogue of the implication $(1) \Rightarrow (2)$ in Theorem 5.5 in

$$\ell^1 \Gamma := \left\{ \sum_{g \in \Gamma} a_g g \mid \sum_{g \in \Gamma} |a_g| < \infty \right\}.$$

We set $\omega^1 \Gamma := \overline{\omega(\Gamma)}^{\|\cdot\|_1}$ and define

$$\Sigma^{2,1}\omega(\Gamma) := \left\{ \sum_{i=1}^{\infty} a_i^* a_i \mid a_i \in \omega[\Gamma], \sum_{i=1}^{\infty} \|a_i\|_1^2 < \infty \right\}$$

and

$$\Sigma^{2,1}\ell^{1}(\Gamma) := \left\{ \sum_{i=1}^{\infty} a_{i}^{*}a_{i} \mid a_{i} \in \ell^{1}[\Gamma], \sum_{i=1}^{\infty} \|a_{i}\|_{1}^{2} < \infty \right\}.$$

We note that $||a||_1 - a \in \Sigma^{2,1} \ell^1(\Gamma)$ for every hermitian element $a \in \ell^1 \Gamma$. Indeed,

$$\|a\|_{1} - a = \sum_{g \in G} 2|a_{g}| - a_{g}g - \bar{a}_{g}g^{-1} = \sum_{g \in G} \left(|a_{g}|^{1/2} - \frac{a_{g}}{|a_{g}|^{1/2}}g\right)^{*} \left(|a_{g}|^{1/2} - \frac{a_{g}}{|a_{g}|^{1/2}}g\right).$$

Hence, for $\varphi \colon \ell^1 \Gamma \to \mathbb{C}$, \mathbb{C} -linear and positive on $\Sigma^{2,1}\ell^1\Gamma$, $|\varphi(a)| \le ||a||_1$ for all $a \in \ell^1\Gamma$. A priori, there is no reason to assume that $\Sigma^{2,1}\ell^1\Gamma$ or $\Sigma^{2,1}\omega^1\Gamma$ are closed or have nontrivial interior. Nevertheless, our result shows:

Corollary 5.9. Let Γ be a Kazhdan group with finite generating symmetric set S and Kazhdan constant ϵ . Then for every $b \in \omega^1(\Gamma)^h$ with $||b||_1 = 1$, we have

$$\Delta(S) + \varepsilon \cdot b \in \Sigma^{2,1} \omega(\Gamma).$$

6. Group algebras of free groups

In this section, let $\Gamma = F_n$ be the free group on *n* generators g_1, \ldots, g_n or $\Gamma = F_\infty$. Schmüdgen (private communication, 2011) has proven that an element from the group algebra $\mathbb{C}[\Gamma]$ that is nonnegative under each finite-dimensional *-representation is a sum of squares. We demonstrate how his result can be reproved with our real closed separation approach. The main idea of our proof is the same as in Schmüdgen's work. However, instead of a partial GNS construction, we use a full GNS construction, but over a general real closed field. We then reduce to the standard real numbers by Tarski's transfer principle.

Theorem 6.1 (Schmüdgen). Let $\Gamma = F_n$ be the free group on *n* generators. If $b \in \mathbb{C}[\Gamma]^h$ is mapped to a positive semidefinite matrix under each finite-dimensional *-representation of $\mathbb{C}[\Gamma]$, then $b \in \Sigma^2 \mathbb{C}[\Gamma]$.

Proof. Assume that $b \notin \Sigma^2 \mathbb{C}[\Gamma]$. By Theorem 3.11, there is a real closed extension field **R** of \mathbb{R} and a completely positive \mathbb{C} -linear functional $\varphi \colon \mathbb{C}[\Gamma] \to \mathbf{C}$ with $\varphi(a^*) = \overline{\varphi(a)}$, such that $\varphi(b) < 0$. By Lemma 3.6, the canonical **C**-linear extension

of φ to $A = \mathbb{C} \otimes_{\mathbb{C}} \mathbb{C}[\Gamma]$ is still positive, and we denote it again by φ . We apply the usual GNS construction to A. We note that

$$N = \{a \in A \mid \varphi(a^*a) = 0\}$$

is a *-subspace of the C-vector space A, which follows from the Cauchy–Schwarz inequality, as shown in Corollary 3.9. We denote the quotient space A/N by H, and note that

$$\langle a+N, c+N \rangle_{\varphi} := \varphi(c^*a)$$

is a well-defined and positive definite C-valued sesquilinear form on H. We also note that left-multiplication from $\mathbb{C}[\Gamma]$ on A is well-defined on H, as explained in Section 5. So we have a \mathbb{C} -linear *-representation

$$\pi: \mathbb{C}[\Gamma] \to \mathscr{L}(H)$$

with $\langle \pi(b)\xi, \xi \rangle_{\varphi} < 0$, where $\xi = 1 + N$.

Now let H' be a finite-dimensional *-subspace of H, containing the residue classes of all words in the generators g_i of length at most d, where d is the maximal word length in b. We can choose an orthonormal basis $v_1, \ldots v_m$ of H', using the usual Gram–Schmidt procedure over **C**. So there is an orthogonal projection map $p: H \rightarrow H'$, defined as

$$p: h \mapsto \sum_{i=1}^m \langle h, v_i \rangle v_i.$$

Define

$$M_i := p \circ \pi(g_i) \in \mathcal{L}(H').$$

It is easy to see that all M_i are contractions; thus the linear operators $\sqrt{1 - M_i^* M_i}$ and $\sqrt{1 - M_i M_i^*}$ exist on H'. Using Choi's matrix trick [1980, Theorem 7], we define

$$U_i := \begin{pmatrix} M_i & \sqrt{I - M_i M_i^*} \\ \sqrt{I - M_i^* M_i} & -M_i^* \end{pmatrix} \in \mathcal{L}(H' \oplus H')$$

The U_i are checked to be unitary operators, and thus yield a \mathbb{C} -linear *-representation $\widetilde{\pi}$ of $\mathbb{C}[\Gamma]$ on $H' \oplus H'$. Since the residue classes of all words occurring in *b* belong to H', and by the definition of the U_i , we find

$$\langle \widetilde{\pi}(b)\xi',\xi' \rangle_{H'\oplus H'} = \langle \pi(b)\xi,\xi \rangle_{\varphi} < 0,$$

where $\xi' = (\xi, 0)$. Now finally, since $H' \oplus H'$ is finite-dimensional, the existence of such a representation over **C** implies the existence over \mathbb{C} , by Tarski's transfer principle. This finishes the proof.

Remark 6.2. The proof becomes even simpler when considering the *-algebra of polynomials in noncommuting variables $\mathbb{C}\langle y_1, \ldots, y_n, z_1, \ldots, z_n \rangle$ with $y_i^* = z_i$, or $\mathbb{C}\langle z_1, \ldots, z_n \rangle$ with $z_i^* = z_i$, instead of the group algebra of a free group. The reason

is that one is not forced to make the matrices M_i unitary (only hermitian in the second case). So Theorem 6.1 also holds for these polynomial algebras. This was first proven by Helton [2002].

It is an interesting problem to study the class of groups Γ for which positivity of $a \in \mathbb{C}[\Gamma]$ in every finite-dimensional unitary representation implies that $a \in \Sigma^2 \mathbb{C}[\Gamma]$. It is clear that in order for an analogous argument to work, Γ has to be residually finite-dimensional in a very strong sense. Residual finite-dimensionality of a group means that every unitary representation on a Hilbert space can be approximated in the Fell topology by finite-dimensional representations; see [Brown and Ozawa 2008] for details. If — more generally — every generalized unitary representation of Γ on a Hilbert space can be approximated on finitely many vectors by generalized finite-dimensional unitary representations, then everything works. With additional work, this can be carried out for virtually free groups (Schmüdgen, private communication).

Deep results of Scheiderer [2006] imply that the conclusion holds for \mathbb{Z}^2 . By a classical result [Rudin 1963], however, the group \mathbb{Z}^3 does not satisfy the desired conclusion, and the same holds for every group containing \mathbb{Z}^3 . This is also implied by seminal work of Scheiderer [2000, Theorem 6.2], who showed that the existence of positive elements that are not sums of squares under general assumptions in dimension ≥ 3 .

This shows that the theory of generalized unitary representations is fundamentally different and new pathologies occur.

An intriguing and possibly manageable case is that of surface groups. Lubotzky and Shalom [2004] showed that surface groups are residually finite-dimensional. It is quite possible that their methods extend and lead to a resolution of the case of surface groups.

Conjecture 6.3. Let Γ be a surface group. Every element $a \in \mathbb{C}[\Gamma]^h$ that is positive semidefinite in every finite-dimensional unitary representation lies in $\Sigma^2 \mathbb{C}[\Gamma]$.

Similar questions can be studied if one allows the unitary representations to be infinite-dimensional. Again, the only known obstruction is $\mathbb{Z}^3 \subset \Gamma$.

References

[[]Barvinok 2002] A. Barvinok, *A course in convexity*, Graduate Studies in Mathematics **54**, American Mathematical Society, Providence, RI, 2002. MR 2003j:52001 Zbl 1014.52001

[[]Bekka et al. 2008] B. Bekka, P. de la Harpe, and A. Valette, *Kazhdan's property (T)*, New Mathematical Monographs **11**, Cambridge University Press, Cambridge, 2008. MR 2009i:22001 Zbl 1146. 22009

[[]Brown and Ozawa 2008] N. P. Brown and N. Ozawa, *C*-algebras and finite-dimensional approximations*, Graduate Studies in Mathematics **88**, American Mathematical Society, Providence, RI, 2008. MR 2009h:46101 Zbl 1160.46001

- [Choi 1980] M. D. Choi, "The full *C**-algebra of the free group on two generators", *Pacific J. Math.* **87**:1 (1980), 41–48. MR 82b:46069 Zbl 0463.46047
- [Cimprič 2009] J. Cimprič, "A representation theorem for Archimedean quadratic modules on *rings", *Canad. Math. Bull.* **52**:1 (2009), 39–52. MR 2010c:46154 Zbl 1173.16024
- [Cimprič et al. 2011] J. Cimprič, M. Marshall, and T. Netzer, "Closures of quadratic modules", *Israel J. Math.* **183** (2011), 445–474. MR 2012e:13047 Zbl 1235.13021
- [Helton 2002] J. W. Helton, "'Positive' noncommutative polynomials are sums of squares", *Ann. of Math.* (2) **156**:2 (2002), 675–694. MR 2003k:12002 Zbl 1033.12001
- [Lubotzky and Shalom 2004] A. Lubotzky and Y. Shalom, "Finite representations in the unitary dual and Ramanujan groups", pp. 173–189 in *Discrete geometric analysis* (Sendai, 2002), edited by M. Kotani et al., Contemp. Math. **347**, Amer. Math. Soc., Providence, RI, 2004. MR 2005e:22011 Zbl 1080.22006
- [Powers and Scheiderer 2001] V. Powers and C. Scheiderer, "The moment problem for non-compact semialgebraic sets", *Adv. Geom.* 1:1 (2001), 71–88. MR 2002c:14086 Zbl 0984.44012
- [Prestel and Delzell 2001] A. Prestel and C. N. Delzell, *Positive polynomials: from Hilbert's 17th problem to real algebra*, Springer, Berlin, 2001. MR 2002k:13044 Zbl 0987.13016
- [Rudin 1963] W. Rudin, "The extension problem for positive-definite functions", *Illinois J. Math.* **7** (1963), 532–539. MR 27 #1779 Zbl 0114.31003
- [Scheiderer 2000] C. Scheiderer, "Sums of squares of regular functions on real algebraic varieties", *Trans. Amer. Math. Soc.* **352**:3 (2000), 1039–1069. MR 2000j:14090 Zbl 0941.14024
- [Scheiderer 2006] C. Scheiderer, "Sums of squares on real algebraic surfaces", *Manuscripta Math.* **119**:4 (2006), 395–410. MR 2006m:14079 Zbl 1120.14047
- [Schmüdgen 1991] K. Schmüdgen, "The *K*-moment problem for compact semi-algebraic sets", *Math. Ann.* **289**:2 (1991), 203–206. MR 92b:44011 Zbl 0744.44008
- [Schmüdgen 2009] K. Schmüdgen, "Noncommutative real algebraic geometry: some basic concepts and first ideas", pp. 325–350 in *Emerging applications of algebraic geometry*, edited by M. Putinar and S. Sullivant, IMA Vol. Math. Appl. 149, Springer, New York, 2009. MR 2010m:13036 Zbl 1158.13313

Received January 24, 2012. Revised August 8, 2012.

TIM NETZER MATHEMATISCHES INSTITUT UNIVERSITY OF LEIPZIG PF 10 09 20 D-04009 LEIPZIG GERMANY netzer@math.uni-leipzig.de

ANDREAS THOM MATHEMATISCHES INSTITUT UNIVERSITY OF LEIPZIG PF 10 09 20 D-04009 LEIPZIG GERMANY thom@math.uni-leipzig.de

UNIQUENESS THEOREM FOR ORDINARY DIFFERENTIAL EQUATIONS WITH HÖLDER CONTINUITY

YIFEI PAN, MEI WANG AND YU YAN

We study ordinary differential equations of the type $u^{(n)}(t) = f(u(t))$, with initial conditions $u(0) = u'(0) = \cdots = u^{(m-1)}(0) = 0$ and $u^{(m)}(0) \neq 0$, where $m \geq n$; no additional assumption is made on f. We establish some uniqueness results and show that f is always Hölder continuous.

1. Introduction

The question of finding criteria for the uniqueness of solutions has been a constant theme in the study of ordinary differential equations for a very long time, and a wealth of results have been established. The one most quoted in textbooks is perhaps the Lipschitz uniqueness theorem, which states that in the equation $y^{(n)}(x) = f(x, y, y', \dots, y^{(n-1)})$, if the function $f(x, z_1, z_2, \dots, z_n)$ is Lipschitz continuous with respect to z_1, z_2, \dots, z_n , then the initial value problem has a unique local solution. Generally speaking, to ensure the uniqueness of solutions to an ODE, we need to assume some condition on the function f besides continuity, the Lipschitz condition being one example. Most of the research in this topic has been devoted to finding the appropriate condition, and there are many nice results, such as the classical theorems by Peano, Osgood, Montel and Tonelli, and Nagumo. An extensive and systematic treatment of the available results is provided in [Agarwal and Lakshmikantham 1993].

In this paper, we approach the uniqueness problem from a different perspective and relate it to the unique continuation problem. We study autonomous ODEs of the type $u^{(n)}(t) = f(u(t))$, where $u \in C^{\infty}([0, 1])$ and no additional assumption is made on the function f.

If we assume the initial conditions $u(0) = u'(0) = \cdots = u^{(n-1)}(0) = 0$, the solution is not unique. The following is a trivial example.

Example 1. $u(t) = t^3$ satisfies $u''(t) = 6u^{1/3}$ and u(0) = u'(0) = 0. Another solution to this initial value problem is $u \equiv 0$.

MSC2010: 34A12.

The research of Yu Yan was partially supported by the Scholar in Residence program at Indiana University–Purdue University Fort Wayne.

Keywords: uniqueness of solutions, ordinary differential equations.

It is no surprise that uniqueness fails in this example, because the function $f(u) = 6u^{1/3}$ has fairly strong singularity at 0. From another perspective, this example shows that if a solution and its derivatives up to order n - 1 all vanish at 0, it is not guaranteed to be the zero function. On the other hand, even if all its derivatives vanish at 0, the solution still may not be identically 0.

Example 2. The function

$$u(t) = \begin{cases} e^{-1/t} & 0 < t \le 1, \\ 0 & t = 0 \end{cases}$$

is in $C^{\infty}([0, 1])$, and

(1)
$$u^{(k)}(0) = 0$$
 for all $k \in \mathbb{N}$.

Let

$$f(s) = \begin{cases} s(\ln s)^2 & s > 0, \\ 0 & s = 0. \end{cases}$$

Then u(t) satisfies the equation

$$u' = f(u).$$

However, this equation has another solution, $u \equiv 0$, which also satisfies (1).

This function u(t) is also a classical example in the study of the unique continuation problem, which asks when we can conclude that a function is locally identically zero if its derivatives all vanish at a point. Here is one result in this line:

Theorem 1.1 [Pan and Wang 2008]. *Let* $g(x) \in C^{\infty}([a, b]), 0 \in [a, b]$, *and*

(2)
$$|g^{(n)}(x)| \le C \sum_{k=0}^{n-1} \frac{|g^{(k)}(x)|}{|x|^{n-k}}, \quad x \in [a, b]$$

for some constant *C* and some $n \ge 1$. Then

$$g^{(k)}(0) = 0 \quad for \ all \quad k \ge 0$$

implies

$$g \equiv 0$$
 on $[a, b]$.

The order of singularity of |x| at 0 in (2), that is, n - k, is sharp, as one example in [Pan and Wang 2008] shows. This theorem is crucial to the proof of our main theorem below.

The previous two examples suggest that to guarantee uniqueness near 0, the solution needs to vanish to sufficiently high order, but not to the infinite order. So we assume that it satisfies the initial conditions $u(0) = u'(0) = \cdots = u^{(m-1)}(0) = 0$ and $u^{(m)}(0) = a \neq 0$, where $m \ge n$. That is, the order of the lowest nonvanishing derivative of u at 0 is no less than the order of the equation. From the equation

it is not difficult to see that f is differentiable away from 0; however, it is not differentiable at 0, as shown by Example 1 with m = 3.

Due to the lack of information about the regularity of f, the available uniqueness theory no longer applies to this type of equation. We will show that because u has sufficiently high vanishing order at 0, such solutions are unique near 0. Specifically, we have the following result.

Theorem 1.2. Let $u(t) \in C^{\infty}([0, 1))$ be a solution of the differential equation

(3)
$$u^{(n)}(t) = f(u(t)),$$

where $n \ge 1$ and f is a function. Assume that u satisfies

(4)
$$u(0) = u'(0) = \dots = u^{(m-1)}(0) = 0$$
 and $u^{(m)}(0) = a \neq 0$,

with $m \ge n$. Then such a solution u(t) is unique for t near 0.

The proof of Theorem 1.2 is carried out in two steps. First, we show the following result concerning the derivatives of u at 0.

Lemma 1.3. Let u(t) be a solution that satisfies Equations (3) and (4). The derivative of u at 0 of any order equal to or higher than m, that is, $u^{(k)}(0)$ for any $k \ge m$, depends only on m, n, and the behavior of the function f near 0.

In the second step, suppose there are two solutions u and v, both satisfying (3) and (4); then by Lemma 1.3, the function u(t) - v(t) and all its derivatives vanish at 0. Making use of Theorem 1.1, we can show that $u - v \equiv 0$.

Typically, for an *n*-th order ODE, we need only *n* initial conditions. Theorem 1.1 shows that in some sense, the lack of information about f can be compensated by assuming additional derivative information at the initial point.

Interestingly, it turns out that the solution is unique as long as the vanishing order is no less than the order of the equation, but the actual vanishing order and the value of the lowest nonzero derivative are not essential.

Theorem 1.4. Suppose u_1 and u_2 are two solutions of (3) that satisfy

(5)
$$u_1(0) = u'_1(0) = \dots = u_1^{(m-1)}(0) = 0, \quad u_1^{(m)}(0) = a \neq 0$$

and

(6)
$$u_2(0) = u'_2(0) = \dots = u_2^{(l-1)}(0) = 0, \quad u_2^{(l)}(0) = b \neq 0,$$

where $m, l \ge n$. Then m = l, a = b, and $u_1 \equiv u_2$ for small t.

The proof of Lemma 1.3 will be given in Section 2, and the proofs of Theorems 1.2 and 1.4 will be given in Section 3.

Naturally, we would like to ask if the result still holds if one of the solutions in Theorem 1.4 has vanishing order lower than n, the order of the equation.

Conjecture. Suppose (3) has a solution u(t) that satisfies (4) with $m \ge n$. Then it cannot possess another solution v(t) that satisfies initial conditions

 $v(0) = v'(0) = \dots = v^{(l-1)}(0) = 0$ and $v^{(l)}(0) = b \neq 0$,

where l < n.

We can show that this conjecture is true if m = n + 1 or l and m are relatively prime. However, there are some difficulties in the general case and we have not been able to prove the full conjecture.

Although in Theorems 1.2 and 1.4 we do not need to make any assumptions about the function f, we can actually obtain interesting information about it. Suppose there is a function u(t) that satisfies Condition (4); then as shown in Section 2, locally t can be expressed as a function of u, and therefore we can express $u^{(n)}(t)$ locally as a function f of u, so $u^{(n)}(t) = f(u)$. The next theorem shows that the function f is Hölder continuous in an interval $[0, \delta]$ for small $\delta > 0$.

Theorem 1.5. Suppose a function u(t) satisfies Condition (4); then Equation (3) holds for some function f, where $m \ge n$ and there is a constant $\delta > 0$ such that f is uniformly Hölder continuous in the interval $[0, \delta]$.

This theorem is proved after Theorem 1.4 in Section 3.

A summary of Theorems 1.2 and 1.5 is that any smooth function of finite order vanishing at 0 is a unique local solution of a differential equation in the form of Equation (3), where f is differentiable in the interior and uniformly Hölder continuous up to the boundary.

In Theorems 1.2 and 1.4, the high-order vanishing condition (4) allows us to obtain uniqueness results without any extra assumption on f. This phenomenon is only found in autonomous equations like (3). For general ODEs of the form

$$\frac{d^n u}{dt^n} = f\left(t, u, \frac{du}{dt}, \dots, \frac{d^{n-1}u}{dt^{n-1}}\right),$$

such results cannot be expected because there is more than one expression for f. For example:

Example 3. Let $u(t) = t^4$. It satisfies the initial conditions

$$u(0) = u'(0) = u''(0) = u^{(3)}(0) = 0$$
 and $u^{(4)}(0) = 24$.

Its derivatives are

$$u'(t) = 4t^3$$
 and $u''(t) = 12t^2$

We can express u'' as a function f of u and u' in different ways, such as

$$u''(t) = \frac{192u^2}{(u')^2}, \qquad u''(t) = \frac{3(u')^2}{4u}, \qquad u''(t) = 12(\frac{1}{4}uu')^{2/7},$$

or

$$u''(t) = 12u^{1/4} \left(\frac{1}{4}u'\right)^{1/3}$$

In the first two equations f is not continuous at the origin, while in the last two equations f is Hölder continuous at the origin. This simple example shows that the function f can be expressed in various ways and that in order to study the uniqueness, we need to impose very specific assumptions on f.

Our work was motivated by that of Li and Nirenberg [2006], who studied a similar second-order PDE: $\Delta u = f(u)$, where $u = u(t, x) \in C^{\infty}(\mathbb{R}^{k+1})$ has a nonvanishing partial derivative at 0 that can be expressed in the form $u(t, x) = at^m + O(t^{m+1})$, $a \neq 0, t \in \mathbb{R}$, and $x \in \mathbb{R}^k$. They showed that if two solutions u and v satisfy $u \ge v$, then $u \equiv v$. Theorem 1.2 can be viewed as an improvement of their result in the one-dimensional case to arbitrary order and without the comparison condition $u \ge v$.

2. The proof of Lemma 1.3

Without loss of generality, we can assume that a > 0.

First, we show that a, the *m*-th derivative of u at 0, only depends on m, n, and the function f.

Define

(7)
$$\tilde{x} = \left(\frac{u}{a}\right)^{1/m}.$$

Then

$$\tilde{x} = (t^m + O(t^{m+1}))^{1/m} = t(1 + O(t)).$$

This implies that

(8)
$$\frac{\dot{x}}{t} \to 1 \quad \text{as} \quad t \to 0$$

We can also write $u = a\tilde{x}^m$. Taking the derivative with respect to t, we get

$$\frac{du}{dt} = am\tilde{x}^{m-1}\frac{d\tilde{x}}{dt},$$
$$amt^{m-1} + O(t^m) = am\tilde{x}^{m-1}\frac{d\tilde{x}}{dt},$$
$$\frac{t^{m-1}}{\tilde{x}^{m-1}} + O\left(\frac{t^m}{\tilde{x}^{m-1}}\right) = \frac{d\tilde{x}}{dt}.$$

In the second equation above and in the analysis that follows, we formally differentiate the Taylor expansion with the big-*O* notation. A detailed discussion of this differentiation is provided in the Appendix. In light of (8), it follows that

(9)
$$\frac{d\tilde{x}}{dt}|_{t=0} = 1.$$

By the inverse function theorem, t can be expressed as a function of \tilde{x} :

$$t = \tilde{x} + O(\tilde{x}^2).$$

Then

$$t^{m-n} = (\tilde{x} + O(\tilde{x}^2))^{m-n} = \tilde{x}^{m-n}(1 + O(\tilde{x})).$$

Similarly, $t^{m-n+1} = \tilde{x}^{m-n+1}(1 + O(\tilde{x}))$. Thus

(10)
$$f(u) = u^{(n)}$$

= $am(m-1) \dots (m-n+1)t^{m-n} + O(t^{m-n+1})$
= $am(m-1) \dots (m-n+1)\tilde{x}^{m-n} + O(\tilde{x}^{m-n+1})$
= $am(m-1) \dots (m-n+1) \left(\frac{u}{a}\right)^{(m-n)/m} + O(u^{(m-n+1)/m})$ by (7)
= $a^{n/m}m(m-1) \dots (m-n+1)u^{(m-n)/m} + O(u^{(m-n+1)/m}).$

Therefore,

(11)
$$a^{n/m} = \lim_{u \to 0} \frac{f(u)}{m(m-1)\dots(m-n+1)u^{(m-n)/m}}.$$

This shows that a is completely determined by m, n, and the behavior of the function f near 0.

Next, we show that the (m + 1)-th derivative of u at 0 also only depends on m, n, and f.

Write *u* as

(12)
$$u(t) = at^{m} + a_{m+1}t^{m+1} + O(t^{m+2}).$$

We will show that a_{m+1} only depends on m, n, and the behavior of f at 0.

Express *t* as

(13)
$$t = \tilde{x} + b_2 \tilde{x}^2 + O(\tilde{x}^3).$$

We would like to obtain an expression for b_2 in terms of the derivatives of u at 0. To do this, we take the derivative with respect to t on both sides of

$$\frac{d\tilde{x}}{dt} \cdot \frac{dt}{d\tilde{x}} = 1.$$

By the product rule and chain rule, we have

(14)
$$\frac{d}{dt}\left(\frac{d\tilde{x}}{dt}\right) \cdot \frac{dt}{d\tilde{x}} + \frac{d\tilde{x}}{dt} \cdot \frac{d}{dt}\left(\frac{dt}{d\tilde{x}}\right) = 0,$$
$$\frac{d^2\tilde{x}}{dt^2} \cdot \frac{dt}{d\tilde{x}} + \frac{d\tilde{x}}{dt} \cdot \left(\frac{d^2t}{d\tilde{x}^2} \cdot \frac{d\tilde{x}}{dt}\right) = 0,$$
$$\frac{d^2\tilde{x}}{dt^2} \cdot \frac{dt}{d\tilde{x}} + \left(\frac{d\tilde{x}}{dt}\right)^2 \cdot \frac{d^2t}{d\tilde{x}^2} = 0.$$

We would like to evaluate (14) at t = 0.

From

$$\begin{split} \tilde{x} &= \left(\frac{u}{a}\right)^{1/m} = \left(t^m + \frac{a_{m+1}}{a}t^{m+1} + O(t^{m+2})\right)^{1/m} \\ &= t\left(1 + \frac{a_{m+1}}{a}t + O(t^2)\right)^{1/m} \\ &= t\left(1 + \frac{1}{m}\left(\frac{a_{m+1}}{a}t + O(t^2)\right) + \frac{1}{2} \cdot \frac{1}{m}\left(\frac{1}{m} - 1\right)\left(\frac{a_{m+1}}{a}t + O(t^2)\right)^2 + O(t^3)\right) \\ &= t + \frac{a_{m+1}}{ma}t^2 + O(t^3), \end{split}$$

we know that

(15)
$$\frac{d^2 \tilde{x}}{dt^2}|_{t=0} = \frac{2a_{m+1}}{ma}$$

From (13) we know that

$$\frac{dt}{d\tilde{x}}|_{t=0} = 1 \quad \text{and} \quad \frac{d^2t}{d\tilde{x}^2}|_{t=0} = 2b_2.$$

Thus if we evaluate (14) at t = 0, we get

$$\frac{2a_{m+1}}{ma} \cdot 1 + 1 \cdot 2b_2 = 0,$$

and therefore

(16) $b_2 = -\frac{a_{m+1}}{ma}.$

Now, from (13) we have

$$t^{m-n} = (\tilde{x} + b_2 \tilde{x}^2 + O(\tilde{x}^3))^{m-n}$$

= $\tilde{x}^{m-n} [1 + b_2 \tilde{x} + O(\tilde{x}^2)]^{m-n}$
= $\tilde{x}^{m-n} [1 + (m-n)(b_2 \tilde{x} + O(\tilde{x}^2)) + O(\tilde{x}^2)]$
= $\tilde{x}^{m-n} + (m-n)b_2 \tilde{x}^{m-n+1} + O(\tilde{x}^{m-n+2}).$

Similarly,

$$t^{m-n+1} = \tilde{x}^{m-n+1} + (m-n+1)b_2\tilde{x}^{m-n+2} + O(\tilde{x}^{m-n+3}),$$

$$t^{m-n+2} = O(\tilde{x}^{m-n+2}).$$

Then from (12) and the above expressions for the powers of t, we have

$$\begin{split} u^{(n)} &= am(m-1)\dots(m-n+1)t^{m-n} \\ &+ a_{m+1}(m+1)m\dots(m-n+2)t^{m-n+1} + O(t^{m-n+2}) \\ &= am(m-1)\dots(m-n+1)\big(\tilde{x}^{m-n} + (m-n)b_2\tilde{x}^{m-n+1} + O(\tilde{x}^{m-n+2})\big) \\ &+ a_{m+1}(m+1)m\dots(m-n+2) \\ &\times \big(\tilde{x}^{m-n+1} + (m-n+1)b_2\tilde{x}^{m-n+2} + O(\tilde{x}^{m-n+3})\big) \\ &+ O(\tilde{x}^{m-n+2}) \\ &= am(m-1)\dots(m-n+1)\tilde{x}^{m-n} \\ &+ \big(am(m-1)\dots(m-n)b_2 + a_{m+1}(m+1)m\dots(m-n+2)\big)\tilde{x}^{m-n+1} \\ &+ O(\tilde{x}^{m-n+2}). \end{split}$$

Thus, in view of $f(u) = u^{(n)}$ and (7), we have

This means that

$$a^{\frac{n-1}{m}}m(m-1)\dots(m-n)b_2 + a^{\frac{n-m-1}{m}}a_{m+1}(m+1)m\dots(m-n+2)$$
$$= \lim_{u \to 0} \frac{f(u) - a^{n/m}m(m-1)\dots(m-n+1)u^{\frac{m-n}{m}}}{u^{\frac{m-n+1}{m}}}.$$

By (16), this can be written as

$$a^{\frac{n-1}{m}}m(m-1)\dots(m-n)\left(-\frac{a_{m+1}}{ma}\right) + a^{\frac{n-m-1}{m}}a_{m+1}(m+1)m\dots(m-n+2)$$
$$= \lim_{u \to 0} \frac{f(u) - a^{\frac{n}{m}}m(m-1)\dots(m-n+1)u^{\frac{m-n}{m}}}{u^{\frac{m-n+1}{m}}}$$

After collecting similar terms, we get

$$a_{m+1}a^{\frac{n-m-1}{m}}((m+1)m\dots(m-n+2) - (m-1)(m-2)\dots(m-n))$$

=
$$\lim_{u \to 0} \frac{f(u) - a^{\frac{n}{m}}m(m-1)\dots(m-n+1)u^{\frac{m-n}{m}}}{u^{\frac{m-n+1}{m}}}.$$

Consequently,

(17)
$$a_{m+1} = \frac{a^{\frac{m-n+1}{m}} \left(\lim_{u \to 0} \frac{f(u) - a^{\frac{n}{m}} m(m-1) \dots (m-n+1) u^{\frac{m-n}{m}}}{u^{\frac{m-n+1}{m}}} \right)}{(m+1)m \dots (m-n+2) - (m-1)(m-2) \dots (m-n)}.$$

Since we have proved that *a* only depends on *m*, *n*, and *f*, Equation (17) shows that a_{m+1} is also completely determined by *m*, *n*, and the behavior of *f* near 0. By (16), this also shows that b_2 depends only on *m*, *n*, and *f*.

Now we will use mathematical induction to show that all the derivatives of u at 0 of order higher than m are completely determined by m, n, and f.

Express u and t as

(18)
$$u(t) = at^m + a_{m+1}t^{m+1} + \dots + a_{m+k}t^{m+k} + a_{m+k+1}t^{m+k+1} + O(t^{m+k+2})$$

and

(19)
$$t = \tilde{x} + b_2 \tilde{x}^2 + \dots + b_{k+1} \tilde{x}^{k+1} + b_{k+2} \tilde{x}^{k+2} + O(\tilde{x}^{k+3}).$$

Suppose that for $k \ge 1$, $a, a_{m+1}, \ldots, a_{m+k}, b_2, \ldots, b_{k+1}$ are all determined only by m, n, and f; we will show that a_{m+k+1} and b_{k+2} also are determined only by m, n, and f.

We start by obtaining an expression for b_{k+2} in terms of a_{m+1}, \ldots, a_{m+k} and a_{m+k+1} .

Taking the derivative with respect to t on both sides of (14), we obtain

$$0 = \frac{d^3\tilde{x}}{dt^3} \cdot \frac{dt}{d\tilde{x}} + \frac{d^2\tilde{x}}{dt^2} \cdot \left(\frac{d^2t}{d\tilde{x}^2} \cdot \frac{d\tilde{x}}{dt}\right) + 2 \cdot \frac{d\tilde{x}}{dt} \cdot \frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^2t}{d\tilde{x}^2} + \left(\frac{d\tilde{x}}{dt}\right)^2 \cdot \left(\frac{d^3t}{d\tilde{x}^3} \cdot \frac{d\tilde{x}}{dt}\right)$$

and

$$0 = \frac{d^3 \tilde{x}}{dt^3} \cdot \frac{dt}{d\tilde{x}} + 3 \cdot \frac{d\tilde{x}}{dt} \cdot \frac{d^2 \tilde{x}}{dt^2} \cdot \frac{d^2 t}{d\tilde{x}^2} + \left(\frac{d\tilde{x}}{dt}\right)^3 \cdot \frac{d^3 t}{d\tilde{x}^3}.$$

Taking the derivative of both sides of these equations, we get

$$(20) \quad 0 = \left(\frac{d^4\tilde{x}}{dt^4} \cdot \frac{dt}{d\tilde{x}} + \frac{d^3\tilde{x}}{dt^3} \cdot \frac{d^2t}{d\tilde{x}^2} \cdot \frac{d\tilde{x}}{dt}\right) \\ + 3\left(\frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^2t}{d\tilde{x}^2} + \frac{d\tilde{x}}{dt} \cdot \frac{d^3\tilde{x}}{dt^3} \cdot \frac{d^2t}{d\tilde{x}^2} + \frac{d\tilde{x}}{dt} \cdot \frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^3t}{d\tilde{x}^3} \cdot \frac{d\tilde{x}}{dt}\right) \\ + \left(3\left(\frac{d\tilde{x}}{dt}\right)^2 \cdot \frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^3t}{d\tilde{x}^3} + \left(\frac{d\tilde{x}}{dt}\right)^3 \cdot \frac{d^4t}{d\tilde{x}^4} \cdot \frac{d\tilde{x}}{dt}\right), \\ 0 = \frac{d^4\tilde{x}}{dt^4} \cdot \frac{dt}{d\tilde{x}} + 4 \cdot \frac{d^3\tilde{x}}{dt^3} \cdot \frac{d\tilde{x}}{dt} \cdot \frac{d^2t}{d\tilde{x}^2} + 3 \cdot \frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^2t}{d\tilde{x}^2} \\ + 6\left(\frac{d\tilde{x}}{dt}\right)^2 \cdot \frac{d^2\tilde{x}}{dt^2} \cdot \frac{d^3t}{d\tilde{x}^3} + \left(\frac{d\tilde{x}}{dt}\right)^4 \cdot \frac{d^4t}{d\tilde{x}^4}.$$

If we take the derivative k times with respect to t and collect the similar terms after each differentiation as shown above, eventually we arrive at an expression of the form

(21)
$$0 = \frac{d^{k+2}\tilde{x}}{dt^{k+2}} \cdot \frac{dt}{d\tilde{x}} + \left(\text{terms involving } \frac{d^{k+1}\tilde{x}}{dt^{k+1}}, \frac{d^k\tilde{x}}{dt^k}, \dots, \frac{d\tilde{x}}{dt}, \frac{dt}{d\tilde{x}}, \frac{d^2t}{d\tilde{x}^2}, \dots, \frac{d^{k+1}t}{d\tilde{x}^{k+1}}\right) + \left(\frac{d\tilde{x}}{dt}\right)^{k+2} \cdot \frac{d^{k+2}t}{d\tilde{x}^{k+2}}.$$

From (19) we know that

(22)

$$\frac{dt}{d\tilde{x}}\Big|_{t=0} = 1, \\
\frac{d^{2}t}{d\tilde{x}^{2}}\Big|_{t=0} = 2b_{2}, \\
\vdots \\
\frac{d^{k+1}t}{d\tilde{x}^{k+1}}\Big|_{t=0} = (k+1)! b_{k+1}, \\
\frac{d^{k+2}t}{d\tilde{x}^{k+2}}\Big|_{t=0} = (k+2)! b_{k+2}.$$

Then we look at
$$\frac{d\tilde{x}}{dt}\Big|_{t=0}, \dots, \frac{d^{k+1}\tilde{x}}{dt^{k+1}}\Big|_{t=0}$$
, and $\frac{d^{k+2}\tilde{x}}{dt^{k+2}}\Big|_{t=0}$. By (7) and (18),

$$\begin{split} \tilde{x} &= \left(\frac{at^m + a_{m+1}t^{m+1} + \dots + a_{m+k}t^{m+k} + a_{m+k+1}t^{m+k+1} + O(t^{m+k+2})}{a}\right)^{1/m} \\ &= t \left(1 + \frac{a_{m+1}}{a}t + \frac{a_{m+2}}{a}t^2 + \dots + \frac{a_{m+k}}{a}t^k + \frac{a_{m+k+1}}{a}t^{k+1} + O(t^{k+2})\right)^{1/m} \\ &= t \left\{1 + \frac{1}{m} \left[\frac{a_{m+1}}{a}t + \frac{a_{m+2}}{a}t^2 + \dots + \frac{a_{m+k}}{a}t^k + \frac{a_{m+k+1}}{a}t^{k+1} + O(t^{k+2})\right] \\ &+ \frac{1}{2} \cdot \frac{1}{m} \left(\frac{1}{m} - 1\right) \\ &\qquad \left[\frac{a_{m+1}}{a}t + \frac{a_{m+2}}{a}t^2 + \dots + \frac{a_{m+k}}{a}t^k + \frac{a_{m+k+1}}{a}t^{k+1} + O(t^{k+2})\right]^2 \\ &+ \dots \\ &+ \frac{1}{(k+1)!} \cdot \frac{1}{m} \cdot \left(\frac{1}{m} - 1\right) \cdots \left(\frac{1}{m} - k\right) \\ &\qquad \left[\frac{a_{m+1}}{a}t + \frac{a_{m+2}}{a}t^2 + \dots + \frac{a_{m+k}}{a}t^k + \frac{a_{m+k+1}}{a}t^{k+1} + O(t^{k+2})\right]^{k+1} \\ &+ O(t^{k+2}) \bigg\} \end{split}$$

After collecting similar terms, we can write

(23)
$$\tilde{x} = t + \lambda_2 t^2 + \lambda_3 t^3 + \dots + \lambda_{k+1} t^{k+1} + \left(\frac{a_{m+k+1}}{ma} + \lambda_{k+2}\right) t^{k+2} + O(t^{k+3}),$$

where

÷

- λ_2 is a constant involving *m*, *a*, and a_{m+1} ,
- λ_3 is a constant involving m, a, a_{m+1} and a_{m+2} ,
- λ_{k+1} is a constant involving $m, a, a_{m+1}, \ldots, a_{m+k-1}, a_{m+k}$, and
- λ_{k+2} is a constant involving $m, a, a_{m+1}, \ldots, a_{m+k-1}, a_{m+k}$.

By the inductive hypothesis, $\lambda_2, \lambda_3, \ldots, \lambda_{k+1}, \lambda_{k+2}$ are all constants that only depend on *m*, *n*, and the function *f*.

From (23), we obtain

$$\frac{d\tilde{x}}{dt}\Big|_{t=0} = 1,$$
$$\frac{d^2\tilde{x}}{dt^2}\Big|_{t=0} = 2\lambda_2,$$

÷

(24)

$$\frac{d^{k+1}\tilde{x}}{dt^{k+1}}\Big|_{t=0} = (k+1)!\,\lambda_{k+1},\\ \frac{d^{k+2}\tilde{x}}{dt^{k+2}}\Big|_{t=0} = (k+2)!\,\left(\frac{a_{m+k+1}}{ma} + \lambda_{k+2}\right).$$

Now we evaluate (21) at t = 0 and make use of (22) and (24):

$$0 = (k+2)! \left(\frac{a_{m+k+1}}{ma} + \lambda_{k+2}\right) \cdot 1$$

+ (terms involving $b_2, \dots, b_{k+1}, \lambda_2, \dots, \lambda_{k+1}$) + 1 \cdot (k+2)! b_{k+2} .

Thus we obtain

(25)
$$b_{k+2} = -\frac{a_{m+k+1}}{ma} + Q,$$

where Q is a constant depending on $b_2, \ldots, b_{k+1}, \lambda_2, \ldots, \lambda_{k+1}, \lambda_{k+2}$, and hence Q is completely determined by m, n, and f.

Next we will analyze a_{m+k+1} . From (19) we have

$$\begin{split} t^{m-n} &= \tilde{x}^{m-n} \Big(1 + b_2 \tilde{x} + \dots + b_{k+1} \tilde{x}^k + b_{k+2} \tilde{x}^{k+1} + O(\tilde{x}^{k+2}) \Big)^{m-n} \\ &= \tilde{x}^{m-n} \bigg\{ 1 + (m-n) \big(b_2 \tilde{x} + \dots + b_{k+1} \tilde{x}^k + b_{k+2} \tilde{x}^{k+1} + O(\tilde{x}^{k+2}) \big) \\ &+ \frac{(m-n)(m-n-1)}{2} \big(b_2 \tilde{x} + \dots + b_{k+1} \tilde{x}^k + b_{k+2} \tilde{x}^{k+1} + O(\tilde{x}^{k+2}) \big)^2 + \dots \\ &+ \frac{(m-n)(m-n-1)\dots(m-n-k)}{(k+1)!} \\ &\quad (b_2 \tilde{x} + \dots + b_{k+1} \tilde{x}^k + b_{k+2} \tilde{x}^{k+1} + O(\tilde{x}^{k+2}) \big)^{k+1} \\ &+ O(\tilde{x}^{k+2}) \bigg\}. \end{split}$$

After collecting similar terms, we can express t^{m-n} as

$$t^{m-n} = \tilde{x}^{m-n} \{ 1 + c_{1,m-n} \tilde{x} + c_{2,m-n} \tilde{x}^2 + \dots + c_{k,m-n} \tilde{x}^k + ((m-n)b_{k+2} + c_{k+1,m-n}) \tilde{x}^{k+1} + O(\tilde{x}^{k+2}) \},\$$

where $c_{1,m-n}$ is a constant depending on m and b_2 ; $c_{2,m-n}$ is a constant depending on m, b_2 , and b_3 ; ...; $c_{k,m-n}$ is a constant depending on m, b_2 , ..., b_{k+1} ; $c_{k+1,m-n}$ is a constant depending on m, b_2 , ..., b_{k+1} .

By the inductive hypothesis, $c_{1,m-n}, c_{2,m-n}, \ldots, c_{k,m-n}$ and $c_{k+1,m-n}$ are all determined only by m, n, and f. Thus we have

(26)
$$t^{m-n} = \tilde{x}^{m-n} + c_{1,m-n}\tilde{x}^{m-n+1} + c_{2,m-n}\tilde{x}^{m-n+2} + \dots + c_{k,m-n}\tilde{x}^{m-n+k} + ((m-n)b_{k+2} + c_{k+1,m-n})\tilde{x}^{m-n+k+1} + O(\tilde{x}^{m-n+k+2}),$$

where $c_{1,m-n}, c_{2,m-n}, \ldots, c_{k,m-n}$ and $c_{k+1,m-n}$ are constants depending on m, n, and f.

By the same type of analysis we obtain similar expressions for the other powers of *t*:

(27)
$$t^{m-n+1} = \tilde{x}^{m-n+1} + c_{1,m-n+1}\tilde{x}^{m-n+2} + c_{2,m-n+1}\tilde{x}^{m-n+3} + \cdots + c_{k,m-n+1}\tilde{x}^{m-n+k+1} + ((m-n+1)b_{k+2} + c_{k+1,m-n+1})\tilde{x}^{m-n+k+2} + O(\tilde{x}^{m-n+k+3}),$$

where $c_{1,m-n+1}, c_{2,m-n+1}, \ldots, c_{k,m-n+1}$, and $c_{k+1,m-n+1}$ are constants depending on *m*, *n*, and *f*.

(28)
$$t^{m-n+2} = \tilde{x}^{m-n+2} + c_{1,m-n+2}\tilde{x}^{m-n+3} + c_{2,m-n+2}\tilde{x}^{m-n+4} + \cdots + c_{k,m-n+2}\tilde{x}^{m-n+k+2} + ((m-n+2)b_{k+2} + c_{k+1,m-n+2})\tilde{x}^{m-n+k+3} + O(\tilde{x}^{m-n+k+4}),$$

where $c_{1,m-n+2}, c_{2,m-n+2}, \ldots, c_{k,m-n+2}$ and $c_{k+1,m-n+2}$ are constants depending on *m*, *n*, and *f*. Proceeding inductively,

(29)
$$t^{m-n+k} = \tilde{x}^{m-n+k} + c_{1,m-n+k} \tilde{x}^{m-n+k+1} + c_{2,m-n+k} \tilde{x}^{m-n+k+2} + \cdots + c_{k,m-n+k} \tilde{x}^{m-n+2k} + ((m-n+k)b_{k+2} + c_{k+1,m-n+k}) \tilde{x}^{m-n+2k+1} + O(\tilde{x}^{m-n+2k+2}),$$

where $c_{1,m-n+k}, c_{2,m-n+k}, \ldots, c_{k,m-n+k}$ and $c_{k+1,m-n+k}$ are constants depending on *m*, *n*, and *f*.

(30)
$$t^{m-n+k+1} = \tilde{x}^{m-n+k+1} + c_{1,m-n+k+1}\tilde{x}^{m-n+k+2} + c_{2,m-n+k+1}\tilde{x}^{m-n+k+3} + \dots + c_{k,m-n+k+1}\tilde{x}^{m-n+2k+1} + ((m-n+k+1)b_{k+2} + c_{k+1,m-n+k+1})\tilde{x}^{m-n+2k+2} + O(\tilde{x}^{m-n+2k+3}),$$

where $c_{1,m-n+k+1}, \ldots, c_{k,m-n+k+1}$ and $c_{k+1,m-n+k+1}$ are constants depending on m, n, and f.

(31)
$$t^{m-n+k+2} = O(\tilde{x}^{m-n+k+2}).$$

From (18) we obtain

$$u^{(n)} = am(m-1)\dots(m-n+1)t^{m-n} + a_{m+1}(m+1)m\dots(m-n+2)t^{m-n+1} +\dots + a_{m+k}(m+k)(m+k-1)\dots(m-n+k+1)t^{m-n+k} + a_{m+k+1}(m+k+1)(m+k)\dots(m-n+k+2)t^{m-n+k+1} + O(t^{m-n+k+2}).$$

Then by (26) to (31), we can write

$$\begin{split} u^{(m)} &= am(m-1)\dots(m-n+1)\Big\{\tilde{x}^{m-n}+c_{1,m-n}\tilde{x}^{m-n+1}+c_{2,m-n}\tilde{x}^{m-n+2}+\dots\\ &+c_{k,m-n}\tilde{x}^{m-n+k}+((m-n)b_{k+2}+c_{k+1,m-n})\tilde{x}^{m-n+k+1}+O(\tilde{x}^{m-n+k+2})\Big\}\\ &+a_{m+1}(m+1)m\dots(m-n+2)\Big\{\tilde{x}^{m-n+1}+c_{1,m-n+1}\tilde{x}^{m-n+2}\\ &+c_{2,m-n+1}\tilde{x}^{m-n+3}+\dots+c_{k,m-n+1}\tilde{x}^{m-n+k+1}\\ &+((m-n+1)b_{k+2}+c_{k+1,m-n+1})\tilde{x}^{m-n+k+2}+O(\tilde{x}^{m-n+k+3})\Big\}+\dots\\ &+a_{m+k}(m+k)(m+k-1)\dots(m+k-n+1)\Big\{\tilde{x}^{m-n+k}\\ &+c_{1,m-n+k}\tilde{x}^{m-n+k+1}+c_{2,m-n+k}\tilde{x}^{m-n+k+2}+\dots\\ &+c_{k,m-n+k}\tilde{x}^{m-n+2k}+((m-n+k)b_{k+2}+c_{k+1,m-n+k})\tilde{x}^{m-n+2k+1}\\ &+O(\tilde{x}^{m-n+2k+2})\Big\}\\ &+a_{m+k+1}(m+k+1)(m+k)\dots(m+k-n+2)\Big\{\tilde{x}^{m-n+k+1}\\ &+c_{1,m-n+k+1}\tilde{x}^{m-n+2k+1}\\ &+((m-n+k+1)b_{k+2}+c_{k+1,m-n+k+1})\tilde{x}^{m-n+2k+2}+O(\tilde{x}^{m-n+2k+3})\Big\}\\ &+O(\tilde{x}^{m-n+k+2})\\ &=am(m-1)\dots(m-n+1)\tilde{x}^{m-n}+C(m,a,a_{m+1},c_{1,m-n})\tilde{x}^{m-n+1}\\ &+C(m,a,a_{m+1},a_{m+2},c_{1,m-n+1},c_{2,m-n})\tilde{x}^{m-n+2}+\dots\\ &+C(m,a,a_{m+1},\dots,a_{m+k},c_{k,1,m-n},c_{k,m-n+1},\dots,c_{1,m-n+k})\Big)\tilde{x}^{m-n+k+1}\\ &+O(\tilde{x}^{m-n+k+2}). \end{split}$$

Here $C(m, a, a_{m+1}, c_{1,m-n})$ is a constant depending on m, a, a_{m+1} , and $c_{1,m-n}$; we denote it as p_{m-n+1} to simplify notations. Since a, a_{m+1} , and $c_{1,m-n}$ only depend on m, n, and f, we know that p_{m-n+1} only depends on m, n, and f.

Similarly, the other constants $C(m, a, a_{m+1}, a_{m+2}, c_{1,m-n+1}, c_{2,m-n}), \ldots$, and $C(m, a, a_{m+1}, \ldots, a_{m+k}, c_{k+1,m-n}, c_{k,m-n+1}, \ldots, c_{1,m-n+k})$ all depend on m, n, and f only, and can be denoted simply as $p_{m-n+2}, \ldots, p_{m-n+k}$, and $p_{m-n+k+1}$. Thus we can rewrite the above equation as

(32)
$$u^{(n)} = am(m-1)\dots(m-n+1)\tilde{x}^{m-n} + p_{m-n+1}\tilde{x}^{m-n+1} + p_{m-n+2}\tilde{x}^{m-n+2} + \dots + p_{m-n+k}\tilde{x}^{m-n+k} + (am(m-1)\dots(m-n)b_{k+2} + a_{m+k+1}(m+k+1)(m+k)\dots(m+k-n+2) + p_{m-n+k+1})\tilde{x}^{m-n+k+1} + O(\tilde{x}^{m-n+k+2}).$$

Now because of $u^{(n)} = f(u)$ and definition (7), we have

$$f(u) = am(m-1)\dots(m-n+1)\left(\frac{u}{a}\right)^{(m-n/m)} + p_{m-n+1}\left(\frac{u}{a}\right)^{(m-n+1/m)} + p_{m-n+2}\left(\frac{u}{a}\right)^{(m-n+2)/m} + \dots + p_{m-n+k}\left(\frac{u}{a}\right)^{(m-n+k)/m} + (am(m-1)\dots(m-n)b_{k+2} + a_{m+k+1}(m+k+1)(m+k)\dots(m+k-n+2) + p_{m-n+k+1}\right)\left(\frac{u}{a}\right)^{(m-n+k+1)/m} + O(u^{(m-n+k+2)/m}).$$

Due to (25), we can rewrite this equation as

$$(33) \quad f(u) = am(m-1)\dots(m-n+1)\left(\frac{u}{a}\right)^{(m-n)/m} + p_{m-n+1}\left(\frac{u}{a}\right)^{(m-n+1)/m} + p_{m-n+2}\left(\frac{u}{a}\right)^{(m-n+2)/m} + \dots + p_{m-n+k}\left(\frac{u}{a}\right)^{(m-n+k)/m} + \left\{\left((m+k+1)(m+k)\dots(m+k-n+2)\right) - (m-1)(m-2)\dots(m-n)a_{m+k+1} + am(m-1)\dots(m-n)Q + p_{m-n+k+1}\right\}\left(\frac{u}{a}\right)^{(m-n+k+1)/m} + O(u^{(m-n+k+2)/m}).$$

From (33) we get

$$\left((m+k+1)(m+k)\dots(m+k-n+2) - (m-1)(m-2)\dots(m-n) \right) a_{m+k+1} + am(m-1)\dots(m-n)Q + p_{m-n+k+1} = \lim_{u \to 0} \frac{f(u) - am\dots(m-n+1) \left(\frac{u}{a}\right)^{\frac{m-n}{m}} - p_{m-n+1} \left(\frac{u}{a}\right)^{\frac{m-n+k}{m}} - \dots - p_{m-n+k} \left(\frac{u}{a}\right)^{\frac{m-n+k}{m}} }{\left(\frac{u}{a}\right)^{\frac{m-n+k+1}{m}}}$$

Note that $(m+k+1)(m+k)...(m+k-n+2)-(m-1)(m-2)...(m-n) \neq 0$; then since the constants $Q, a, p_{m-n+1}, ..., p_{m-n+k+1}$ all depend only on m, n, and f, we know that a_{m+k+1} only depends on m, n, and f. Consequently, b_{k+2} also only depends on m, n, and f because of (25).

Therefore, by mathematical induction, all derivatives of f at 0 are determined completely by m, n, and f. This completes the proof of Lemma 1.3.

3. The proofs of Theorems 1.2, 1.4 and 1.5

Proof of Theorem 1.2. Suppose there are two solutions u(t) and v(t), both satisfying Equations (3) and (4). By Lemma 1.3, at t = 0, u and v have the same derivative of any order. Let w = u - v; then

$$w^{(k)}(0) = 0$$
 for any integer $k \ge 0$.

In order to apply Theorem 1.1, we need to show that w satisfies Condition (2).

$$w^{(n)}(t) = u^{(n)}(t) - v^{(n)}(t) = f(u(t)) - f(v(t)).$$

Without loss of generality, we assume a > 0. By Equation (10), we can write

(34)
$$f(u) = \left[a^{n/m}m(m-1)\dots(m-n+1)u^{(m-n)/m} + \alpha(u)\right]$$

and

$$f(v) = \left[a^{n/m}m(m-1)\dots(m-n+1)v^{(m-n)/m} + \alpha(v)\right],$$

where α is a function with the order

$$\alpha(s) = O(s^{(m-n+1)/m}).$$

So we can write

(35)
$$w^{(n)}(t) = a^{n/m}m(m-1)\dots(m-n+1)(u^{(m-n)/m}-v^{(m-n)/m}) + (\alpha(u)-\alpha(v)).$$

If m = n, then

$$a^{n/m}m(m-1)\dots(m-n+1)(u^{(m-n)/m}-v^{(m-n)/m})=0.$$

If m > n, by the mean value theorem,

(36)
$$|u^{(m-n)/m} - v^{(m-n)/m}| \le \frac{m-n}{m} \zeta^{-n/m} |u-v|.$$

where $\zeta(t)$ is between u(t) and v(t). Since $u(t) = at^m + O(t^{m+1})$ and $v(t) = at^m + O(t^{m+1})$, we know that $\zeta(t) = at^m + O(t^{m+1})$, which implies

$$\zeta^{n/m} = a^{n/m} t^n (1 + O(t)) \ge C t^n$$

for some constant C > 0 when t is sufficiently small. Thus

$$\zeta^{-n/m}|u-v| \le C^{-1} \frac{|u-v|}{t^n},$$

and by (36), we know that

(37)
$$a^{n/m}m(m-1)\dots(m-n+1)\left|u^{(m-n)/m}-v^{(m-n)/m}\right| \le C\frac{|u-v|}{t^n}$$

for another constant C > 0.

Next, we estimate $|\alpha(u) - \alpha(v)|$.

From (3), we know that f is differentiable with respect to t, since $u^{(n)}(t)$ is differentiable with respect to t. Condition (4) shows that $(du/dt)(t) \neq 0$ when t > 0 is sufficiently small. Then by the inverse function theorem, t is differentiable with respect to u. Thus, when u is small and positive, f is differentiable with respect to u and

$$\frac{df}{du} = \frac{df}{dt} \cdot \frac{dt}{du}$$

Then by (34), since f is differentiable on a small interval $(0, \delta)$, α is also differentiable on a small interval $(0, \delta)$. By the mean value theorem,

$$\alpha(u) - \alpha(v) = \alpha'(\eta)(u - v),$$

where $\eta(t)$ is between u(t) and v(t). Because $u(t) = at^m + O(t^{m+1})$ and $v(t) = at^m + O(t^{m+1})$, we know that $\eta(t) = at^m + O(t^{m+1}) \ge Ct^m$ when t is small, and thus

$$\eta^{-(n-1)/m} = O(t^{-(n-1)}).$$

From $\alpha(s) = O(s^{(m-n+1)/m})$ we get $\alpha'(s) = O(s^{-(n-1)/m})$. Therefore

$$\alpha'(\eta) = O(\eta^{-(n-1)/m}) = O(t^{-(n-1)}).$$

Thus for some C > 0,

(38)
$$|\alpha(u) - \alpha(v)| \le Ct^{-(n-1)}|u - v|$$
$$\le Ct^{-n}|u - v| \quad \text{since} \quad 0 < t < 1.$$

Combining Equations (35), (37), and (38), we conclude

$$|w^{(n)}(t)| \le C \frac{|u(t) - v(t)|}{t^n} = C \frac{|w(t)|}{t^n}.$$

Finally, extend the domain of w(t) to [-1, 1] by defining w(t) = w(-t) when $-1 \le t < 0$. Then $w \in C^{\infty}([-1, 1])$ and it satisfies Condition (2). By Theorem 1.1, $w \equiv 0$, which means $u \equiv v$.

This completes the proof of Theorem 1.2.

Proof of Theorem 1.4. Without loss of generality we assume a > 0. We apply the same analysis as in the proof of Lemma 1.3 to u_1 and u_2 , respectively. Similar to (11), we have

$$a^{n/m} = \lim_{u_1 \to 0} \frac{f(u_1)}{m(m-1)\dots(m-n+1)u_1^{(m-n)/m}}$$
$$= \lim_{s \to 0} \frac{f(s)}{m(m-1)\dots(m-n+1)s^{(m-n)/m}}$$

and

$$b^{n/l} = \lim_{u_2 \to 0} \frac{f(u_2)}{l(l-1)\dots(l-n+1)u_2^{(l-n)/l}} = \lim_{s \to 0} \frac{f(s)}{l(l-1)\dots(l-n+1)s^{(l-n)/l}}.$$

Suppose $m \neq l$; without loss of generality we assume m < l. Dividing the two equations, we get

$$a^{n/m}b^{-n/l} = \frac{l(l-1)\dots(l-n+1)}{m(m-1)\dots(m-n+1)} \lim_{s \to 0} s^{(l-n)/l-(m-n)/m}$$
$$= \frac{l(l-1)\dots(l-n+1)}{m(m-1)\dots(m-n+1)} \lim_{s \to 0} s^{n/m-n/l}.$$

Since m < l, $\lim_{s\to 0} s^{n/m-n/l} = 0$. However, $a^{n/m}b^{-n/l} \neq 0$. This is a contradiction.

Therefore m = l, and consequently a = b. Then by Theorem 1.2, we know that $u_1 \equiv u_2$ for small *t*.

Proof of Theorem 1.5. The proof of Lemma 1.3 shows that near 0, t is a function of u, and therefore $u^{(n)}(t)$ can be expressed as a function f of u. Thus (3) holds when t > 0 is small. From Condition (4), we define f(0) = 0.

By the first two equations in (10) and the discussions in the Appendix, we know that there is a function h that is C^1 on the closed interval $[0, \epsilon]$ for some $\epsilon > 0$, such that

$$f(u) = am(m-1)\dots(m-n+1)\tilde{x}^{m-n} + h(\tilde{x})\tilde{x}^{m-n}$$

By definition (7), we have

(39)
$$f(u) = am(m-1)\dots(m-n+1)\left(\frac{u}{a}\right)^{(m-n)/m} + h\left(\left(\frac{u}{a}\right)^{1/m}\right)\left(\frac{u}{a}\right)^{(m-n)/m}$$
Since $0 \le (m-n)/m < 1$ and 0 < 1/m < 1, it is well known that $u^{(m-n)/m}$ and $u^{1/m}$ are Hölder continuous on the closed interval [0, 1] with Hölder coefficients (m-n)/m and 1/m, respectively. This implies that the first term in (39) is Hölder continuous on [0, 1].

Since *h* is C^1 on $[0, \epsilon]$, it is also Hölder continuous on $[0, \epsilon]$. Then since the composition of two Hölder continuous functions is Hölder continuous, we know that $h((u/a)^{1/m})$ is Hölder continuous with respect to *u* on a closed interval $[0, \delta]$ with $\delta > 0$. Next, because the product of two Hölder continuous functions is also Hölder continuous, we know that $h((u/a)^{1/m}) \cdot (u/a)^{(m-n)/m}$ is Hölder continuous. Thus the second term in (39) is Hölder continuous on $[0, \delta]$.

Therefore, f is Hölder continuous on $[0, \delta]$ and the theorem is proved.

Appendix: Differentiation of the Taylor expansion

We will discuss the regularity of the remainder term in the Taylor expansion of a function that is used in the proof of Theorem 1.5 and the differentiation of the Taylor expansion that is frequently used in the proof of Lemma 1.3.

In general, consider a function $g(x) \in C^{k+1}([a, b])$; by the Taylor theorem, we can write

(40)
$$g(x) = g(a) + g'(a)(x - a)$$

 $+ \frac{g''(a)}{2!}(x - a)^2 + \dots + \frac{g^{(k)}(a)}{k!}(x - a)^k + h(x)(x - a)^k,$

where $\lim_{x \to a} h(x) = 0$. An explicit expression for h(x) is

(41)
$$h(x) = \frac{g^{(k+1)}(\xi)}{(k+1)!}(x-a).$$

where $a < \xi < x$.

From (40) we know that h(x) is C^1 on (a, b]. Next we show that it is actually C^1 up to the boundary, on [a, b].

Taking the derivative on both sides of (40), we get

(42)
$$g'(x) = g'(a) + g''(a)(x-a) + \cdots + \frac{g^{(k)}(a)}{(k-1)!}(x-a)^{k-1} + h'(x)(x-a)^k + kh(x)(x-a)^{k-1}.$$

Define h(a) = 0, so h is continuous on [a, b]. Write

$$P(x) = g(a) + g'(a)(x-a) + \frac{g''(a)}{2!}(x-a)^2 + \dots + \frac{g^{(k)}(a)}{k!}(x-a)^k;$$

then $h(x) = \frac{g(x) - P(x)}{(x - a)^k}$. By the definition of limits,

(43)
$$h'(a) = \lim_{x \to a} \frac{h(x) - h(a)}{x - a} = \lim_{x \to a} \frac{g(x) - P(x)}{(x - a)^{k+1}}$$

$$\vdots$$

$$= \frac{g^{(k+1)}(a) - P^{(k+1)}(a)}{(k+1)!} \quad \text{applying l'Hospital's rule } k + 1 \text{ times}$$

$$= \frac{g^{(k+1)}(a)}{(k+1)!},$$

where we have used the fact that $P^{(k+1)}(a) = 0$.

When x > a,

$$h'(x) = \frac{d}{dx} \left(\frac{g(x) - P(x)}{(x - a)^k} \right)$$

= $\frac{(g'(x) - P'(x))(x - a)^k - (g(x) - P(x))k(x - a)^{k-1}}{(x - a)^{2k}}$
= $\frac{g'(x) - P'(x)}{(x - a)^k} - \frac{k(g(x) - P(x))}{(x - a)^{k+1}}.$

By repeatedly applying l'Hospital's rule, we know that

$$\lim_{x \to a} \frac{g'(x) - P'(x)}{(x-a)^k} = \frac{g^{(k+1)}(a) - P^{(k+1)}(a)}{k!} = \frac{g^{(k+1)}(a)}{k!}$$

and

$$\lim_{x \to a} \frac{k(g(x) - P(x))}{(x - a)^{k+1}} = \frac{k\left(g^{(k+1)}(a) - P^{(k+1)}(a)\right)}{(k+1)!} = \frac{kg^{(k+1)}(a)}{(k+1)!}.$$

Therefore

(44)
$$\lim_{x \to a} h'(x) = \frac{g^{(k+1)}(a)}{k!} - \frac{kg^{(k+1)}(a)}{(k+1)!} = \frac{g^{(k+1)}(a)}{(k+1)!}.$$

Equations (43) and (44) show that h(x) is C^1 on the closed interval [a, b].

Furthermore, we know that for any $x \in [a, b]$, $|h'(x)| \le C$ for some constant C_1 , and thus

$$|h'(x)(x-a)^k| \le C_1|x-a|^k.$$

Since $g(x) \in C^{k+1}([a, b])$, from (41) we know that $|h(x)| \le C_2|x-a|$ for some constant C_2 , and thus

$$|kh(x)(x-a)^{k-1}| \le kC_2|x-a|^k.$$

Therefore, (42) can be written as

(45)
$$g'(x) = g'(a) + g''(a)(x-a) + \dots + \frac{g^{(k)}(a)}{(k-1)!}(x-a)^{k-1} + O(x-a)^k.$$

Since the first, second, ..., and (k-1)-th derivatives of g'(x) at *a* are g''(a), $g^{(3)}(a)$, ..., and $g^{(k)}(a)$, respectively, Equation (45) is the Taylor expansion of g'(x) at *a* to order k-1.

Usually we write (40) as

(46)
$$g(x) = g(a) + g'(a)(x - a)$$

 $+ \frac{g''(a)}{2!}(x - a)^2 + \dots + \frac{g^{(k)}(a)}{k!}(x - a)^k + O((x - a)^{k+1}).$

This shows that we can formally differentiate (46) to get (45).

References

- [Agarwal and Lakshmikantham 1993] R. P. Agarwal and V. Lakshmikantham, *Uniqueness and nonuniqueness criteria for ordinary differential equations*, Series in Real Analysis **6**, World Scientific Publishing, River Edge, NJ, 1993. MR 96e:34002 Zbl 0785.34003
- [Li and Nirenberg 2006] Y. Y. Li and L. Nirenberg, "A geometric problem and the Hopf lemma, II", *Chinese Ann. Math. Ser. B* **27**:2 (2006), 193–218. MR 2007i:53006 Zbl 1149.53302
- [Pan and Wang 2008] Y. Pan and M. Wang, "When is a function not flat?", *J. Math. Anal. Appl.* **340**:1 (2008), 536–542. MR 2008m:26048 Zbl 1131.26016

Received January 25, 2012. Revised February 20, 2012.

YIFEI PAN DEPARTMENT OF MATHEMATICAL SCIENCES INDIANA UNIVERSITY-PURDUE UNIVERSITY FORT WAYNE FORT WAYNE, IN 46805 UNITED STATES and SCHOOL OF MATHEMATICS AND INFORMATICS JIANGXI NORMAL UNIVERSITY NANCHANG CHINA pan@ipfw.edu

MEI WANG DEPARTMENT OF STATISTICS UNIVERSITY OF CHICAGO CHICAGO, IL 60637 UNITED STATES meiwang@galton.uchicago.edu

YU YAN DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE HUNTINGTON UNIVERSITY HUNTINGTON, IN 46750 UNITED STATES yyan@huntington.edu

AN ANALOGUE TO THE WITT IDENTITY

G. A. T. F. DA COSTA AND G. A. ZIMMERMANN

We solve combinatorial and algebraic problems associated with a multivariate identity first considered by Sherman, which he called an analog to the Witt identity. We extend previous results obtained for the univariate case.

1. Introduction

S. Sherman [1962] considered the formal identity in the indeterminates z_1, \ldots, z_n ,

(1-1)
$$\prod_{m_1,\dots,m_R\geq 0} (1+z_1^{m_1}\cdots z_R^{m_R})^{N_+} (1-z_1^{m_1}\cdots z_R^{m_R})^{N_-} = \prod_{j=1}^K (1+z_j)^2,$$

where N_+ and N_- are the number of distinct classes of equivalence of nonperiodic closed paths with positive and negative signs, respectively, which traverse without backtracking m_i times edge i, i = 1, ..., R, of a graph G_R with R > 1 edges forming loops counterclockwise oriented and hooked to a single vertex, $\sum m_i \ge 1$.

Sherman [1962] refers to (1-1) as *an analog to the Witt identity*. The reason will become clear soon. The *Sherman identity*, as we call it, is a special nontrivial case of another identity called the *Feynman identity*, first conjectured by Richard Feynman. This identity relates the Euler polynomial of a graph to a formal product over the classes of equivalence of closed nonperiodic paths with no backtracking in the graph, and it is an important ingredient in a combinatorial formulation of the Ising model in two dimensions, much studied in physics. The Feynman identity was proved for planar and toroidal graphs by Sherman [1960], and in great generality by M. Loebl [2004] and D. Cimasoni [2010].

Sherman compared (1-1) with the multivariate Witt identity [Witt 1937]:

P

(1-2)
$$\prod_{m_1,\dots,m_R\geq 0} (1-z_1^{m_1}\cdots z_R^{m_R})^{\mathcal{M}(m_1,\dots,m_R)} = 1-\sum_{i=1}^{K} z_i,$$

(1-3)
$$\mathcal{M}(m_1, \dots, m_R) = \sum_{g \mid m_1, \dots, m_R} \frac{\mu(g)}{g} \frac{(N/g)!}{(N/g)(m_1/g)! \cdots (m_R/g)!}$$

MSC2010: primary 05C30; secondary 05C25, 05C38.

Keywords: Sherman identity, paths counting, (generalized) Witt formula, free Lie algebras.

where $N = m_1 + \cdots + m_R > 0$, μ is the Möbius function defined by the rules

- (a) $\mu(+1) = +1$,
- (b) $\mu(g) = 0$, for $g = p_1^{e_1} \cdots p_q^{e_q}$, p_1, \dots, p_q primes, and any $e_i > 1$,
- (c) $\mu(p_1 \cdots p_q) = (-1)^q$.

The summation runs over all the common divisors of m_1, \ldots, m_R .

Originally, the Witt identity appeared associated with Lie algebras. In this context the formula gives the dimensions of the homogeneous subspaces of a finitely generated free Lie algebra *L*. If $L(m_1, \ldots, m_R)$ is the subspace of *L* generated by all homogeneous elements of multidegree (m_1, \ldots, m_R) , then dim $L = \mathcal{M}$. However, formula (1-3) has many applications in combinatorics as well [Moree 2005]. Especially relevant is that \mathcal{M} can be interpreted as the number of equivalence classes of closed nonperiodic paths which traverse counterclockwise the edges of G_R , the same graph associated to the Sherman identity (1-1). This property is stated in [Sherman 1962] without a proof, but this combinatorial interpretation of the Witt formula can be reinterpreted as a coloring problem of a necklace with *N* beads with colors chosen out of a set of *R* colors such that the colored beads form a nonperiodic configuration. In other words, $\mathcal{M}(m_1, \ldots, m_R)$ is the number of nonperiodic colored necklaces composed of m_i occurrences of the color $i, i = 1, \ldots, R$.

Sherman [1962] called attention to this association of identities (1-1) and (1-2) to paths in the same graph, which motivated him to consider the problem of finding a relation between (1-1) and Lie algebras. To interpret (1-1) in algebraic terms means to relate the exponents N_{\pm} to some Lie algebraic data.

An investigation of Sherman's problem was initiated in [da Costa 1997; da Costa and Variane 2005] and a solution obtained for the univariate case of identity (1-1). In the present paper we solve the problem in the multivariate formal case, which requires important improvements. The counting method developed in [da Costa 1997; da Costa and Variane 2005] is based on a sign formula for a path given in terms of data encoded in the word representation for the path. It played a crucial role in getting formulas for N_{\pm} in the univariate case. However, the counting method based on this sign formula is complicated. In the present paper we make improvements in the counting method in order to apply it to the multivariate case without depending too much on the sign formula. The formula is used here only to prove a simple lemma.

S-J. Kang and M-H. Kim [1999] derived dimension formulas for the homogeneous spaces of general free graded Lie algebras. We use some of their results to solve Sherman's problem. At the same time our results give a combinatorial realization for some of theirs in terms of paths in a graph.

The paper is organized as follows. In Section 2, we recall the word representation of a path and some basic definitions. We prove a basic lemma about the distribution

of signs in the set of words of a given length. In Section 3, we compute formulas for the number of equivalence classes of closed nonperiodic paths of given length. The first of these generalizes Witt's formula in the sense that it counts paths that traverse the edges of the graph without backtracking. The other formulas give the exponents in Sherman's identity (1-1). We also interpret these formulas in terms of a coloring problem. Sherman's problem, that is, to give an algebraic meaning to the exponents in (1-1) is solved in Section 4.

2. Preliminaries

A path in G_R is an ordered sequence of the edges which does not necessarily respect their orientation. A path is closed and subjected to the constraint that it never goes immediately backwards over a previous edge.

Given $G_r \subseteq G_R$, denote by i_1, \ldots, i_r an enumeration of the edges of G_r in increasing order. A closed path of length $N \ge r$ in G_r is best represented by a word of the form

(2-1)
$$D_{j_1}^{e_{j_1}} D_{j_2}^{e_{j_2}} \cdots D_{j_l}^{e_{j_l}}$$

where $l = r, r + 1, ..., N, j_k \in \{i_1, ..., i_r\}, j_k \neq j_{k+1}, j_l \neq j_1$, and

$$\sum_{k=1}^l |e_{j_k}| = N.$$

All edges of G_r are traversed by a path such that each i_k appears at least once in the sequence $S_l = (j_1, j_2, ..., j_l)$. The order in which the symbols $D_j^{e_j}$ appear in the word indicates the edges traversed by p and in which order. If the sign of e_j is positive, the path traverses edge j exactly e_j times following the edge's orientation; if negative, it traverses the edge $|e_j|$ times in the opposite direction.

A word is called *periodic* if it equals

$$(D_{j_1}^{e_{j_1}}D_{j_2}^{e_{j_2}}\cdots D_{j_{\alpha}}^{e_{j_{\alpha}}})^g$$

for some g > 1. The number g is called the *period* of the word if the word in parentheses is nonperiodic.

Permuting the symbols $D_j^{e_j}$ in (2-1) cyclically, one gets *l* words that represent the same closed path. (For example, $D_1^{-2}D_2^{+1}D_1^{+1}D_2^{+3}$ is a cyclic permutation of $D_2^{+1}D_1^{+1}D_2^{-3}D_1^{-2}$.) Words obtained from one another by a cyclic permutation are taken to be equivalent for this reason. Although the word (2-1) and its inverse

$$D_{j_l}^{-e_{j_l}}\cdots D_{j_1}^{-e_{j_1}}$$

also represent the same path, they are not taken as equivalent here. This is the reason for the exponent 2 on the right side of (1-1), also present in [Sherman 1962].

In Section 3 we consider signed paths. The sign of a path is given by the formula

(2-2)
$$\operatorname{sign}(p) = (-1)^{1+n(p)},$$

where n(p) is the number of integral revolutions of the tangent vector of p. From this definition it follows that if $p = (h)^g$ is a periodic path with odd period g, then sign(p) = sign(h). If g is even, sign(p) = -1. The sign of a path can be computed from its word representation (2-1) using the formula [da Costa and Variane 2005]

$$(2-3) \qquad \qquad (-1)^{N+l+T+s+1}.$$

where *T* is the number of subsequences in the decomposition of S_l into subsequences (see [da Costa and Variane 2005] for the definition and an example of a decomposition) and *s* is the number of negative exponents in (2-1). It follows from the previous sign formulas that periodic words with even period have negative sign.

The following lemma is important in the proof of several results in Section 3. It was assumed in [da Costa 1997; da Costa and Variane 2005] without a proof.

Lemma 2.1. Given $G_r \subseteq G_R$, consider all paths that traverse each edge of G_r at least once (no backtracking allowed) and the set of all representative words (periodic or not, cyclic permutations and inversions included) of fixed length $N \ge r > 1$. Then half of the words have positive sign and the other half have negative sign.

Proof. It suffices to consider the subset of words associated to a fixed sequence $S_l = (j_1, j_2, ..., j_l)$. For this sequence the numbers N, l, and T are fixed. The words with these numbers have signs which depend only on $s \in \{0, 1, 2, ..., l\}$. For N + l + T even, the sign of a word is $(-1)^{s+1}$. If l = 2k, then, for each odd value of s, there are

$$\left(\begin{array}{c}2k\\s\end{array}\right)$$

words with positive sign. Summing over the odd values of *s*, we get the total number of 2^{2k-1} words with positive sign. Summing over the even values of *s*, we get the same number of words with negative sign. If l = 2k + 1, a similar counting gives 2^{2k} words with positive (negative) signs. The case N + l + T odd is analogous. \Box

3. Counting paths in G_r

Fix a subgraph $G_r \subseteq G_R$. Given distinct edges i_1, \ldots, i_r in G_r and positive integers m_{i_1}, \ldots, m_{i_r} , with $m_{i_1} + \cdots + m_{i_r} = N > r$, let $\theta_{\pm}(m_{i_1}, \ldots, m_{i_r})$ be the number of equivalence classes of closed nonperiodic paths of length N with \pm signs that traverse each edge i_j exactly m_{i_j} times, for $j = 1, \ldots, r$, with no backtracking, and traverse the edges in $G_R \setminus G_r$ zero times. In this section we derive formulas for $\theta := \theta_+ + \theta_-$ and θ_{\pm} . Notice that θ_{\pm} is just another name for the exponents N_{\pm} in (1-1) showing only the nonzero entries in N_{\pm} .

Firstly, we compute θ . In the case r = 1, a path with $m_i > 1$ is periodic. The nonperiodic ones are two, the path with length N = 1 and its inversion so that $\theta(m_i) = 0$ if $m_i > 1$ and $\theta(m_i) = 2$ if $m_i = 1$. In other cases, θ is given as follows.

Theorem 3.1. For r = 2, define

(3-1)
$$\mathscr{F}\left(\frac{m_{i_1}}{g}, \frac{m_{i_2}}{g}\right) = \sum_{a=1}^{M/g} \frac{2^{2a}}{a} \binom{m_{i_1}/g - 1}{a - 1} \binom{m_{i_2}/g - 1}{a - 1}$$

where $M = \min\{m_{i_1}, m_{i_2}\}$ and, if $r \ge 3$,

(3-2)
$$\mathscr{F}\left(\frac{m_{i_1}}{g}, \dots, \frac{m_{i_r}}{g}\right) = \sum_{a=r}^{N/g} \frac{2^a}{a} \sum_{\{S_a\}} \prod_{c=1}^r \binom{m_{i_c}/g - 1}{t_{i_c} - 1}$$

where $\{S_a\}$ is the set of sequences (j_1, \ldots, j_a) such that $j_k \in \{i_1, \ldots, i_r\}$ and $j_k \neq j_{k+1}, j_a \neq j_1$. Number t_{i_c} counts how many times edge i_c occurs in a sequence S_a . Use is made of the convention that the combination symbol in (3-2) is zero whenever $t_{i_c} > m_{i_c}/g$. Then

(3-3)
$$\theta(m_{i_1},\ldots,m_{i_r}) = \sum_{g|m_{i_1},\ldots,m_{i_r}} \frac{\mu(g)}{g} \mathscr{F}\left(\frac{m_{i_1}}{g},\ldots,\frac{m_{i_r}}{g}\right).$$

The summation is over all the common divisors g of m_{i_1}, \ldots, m_{i_r} , and $\mu(g)$ is the Möbius function.

Proof. The number $\mathcal{K}(l, m_{i_1}, \ldots, m_{i_r})$ of words that have the same values of m_{i_1}, \ldots, m_{i_r} and $l \in \{r, r+1, \ldots, N\}$ is given by

$$\mathscr{K}(l, m_{i_1}, \ldots, m_{i_r}) = 2^l \sum_{\{S_l\}} \prod_{c=1}^r \binom{m_{i_c} - 1}{n_{i_c} - 1}.$$

Let's explain this formula a bit. The number n_{i_c} counts the number of occurrences of edge i_c in a sequence $S_l = (j_1, \ldots, j_l)$. The combination symbol counts the number of unrestricted partitions of m_{i_c} into n_{i_c} nonzero positive parts [Andrews 1976]; thus the product times 2^l (there are 2^l ways of assigning \pm signs to the exponents in (2-1)) gives the total number of words representing paths traversing each edge i_j of $G_r \subseteq G_R$ exactly m_{i_j} times in all possible ways. Then we sum over all sequences S_l with the convention that a combination symbol equals zero when m < n.

In the set of $\mathcal{K}(l, m_{i_1}, \ldots, m_{i_r})$ words, there is the subset of nonperiodic words plus their cyclic permutations and inversions, and the subset of periodic words, if any, whose periods are the common divisors of $l, m_{i_1}, \ldots, m_{i_r}$ plus their cyclic permutations and inversions. Denote by $\overline{\mathcal{K}(l, m_{i_1}, \ldots, m_{i_r})}$ the number of elements

in the former set. The words with period g are of the form

$$(D_{k_1}^{e_{k_1}}D_{k_2}^{e_{k_2}}\cdots D_{k_{\alpha}}^{e_{k_{\alpha}}})^g$$

where $\alpha = l/g$ and $D_{k_1}^{e_{k_1}} D_{k_2}^{e_{k_2}} \cdots D_{k_{\alpha}}^{e_{k_{\alpha}}}$ is nonperiodic so that the number of periodic words with period g plus their cyclic permutations and inversions is given by $\frac{\mathcal{K}(l/g, m_{i_1}/g, \dots, m_{i_r}/g)}{\mathcal{K}(l/g, m_{i_1}/g, \dots, m_{i_r}/g)}$. Therefore,

$$\mathscr{K}(l, m_{i_1}, \ldots, m_{i_r}) = \sum_{g \mid l, k, m_{i_1}, \ldots, m_{i_r}} \mathscr{K}\left(\frac{l}{g}, \frac{m_{i_1}}{g}, \ldots, \frac{m_{i_r}}{g}\right).$$

The summation is over all the common divisors g of $l, m_{i_1}, \ldots, m_{i_r}$.

Applying the Möbius inversion formula [Apostol 1976], it follows that

(3-4)
$$\overline{\mathscr{K}(l,m_{i_1},\ldots,m_{i_r})} = \sum_{g|(l,m_{i_1},\ldots,m_{i_r})} \mu(g) \mathscr{K}\left(\frac{l}{g},\frac{m_{i_1}}{g},\ldots,\frac{m_{i_r}}{g}\right).$$

where μ is the Möbius function. To eliminate cyclic permutations divide (3-4) by *l*. Summing over all possible values of *l* one gets

(3-5)
$$\theta(m_{i_1}, \dots, m_{i_r}) = \sum_{l=r}^{N} \frac{\overline{\mathcal{K}(l, m_{i_1}, \dots, m_{i_r})}}{l}$$

Upon substitution of (3-4) into (3-5) one gets, for the case $r \ge 3$,

$$\theta(m_{i_1},\ldots,m_{i_r}) = \sum_{l=r}^N \frac{1}{l} \sum_{g \mid (l,m_{i_1},\ldots,m_{i_r})} \mu(g) 2^{l/g} \sum_{\{S_{l/g}\}} \prod_{c=1}^r \binom{m_{i_c}/g - 1}{n_{i_c}/g - 1}.$$

Proceed now as follows. For a given common divisor g of m_{i_1}, \ldots, m_{i_r} , sum over all values of l which are multiples of g. Then sum over all possible divisors of m_{i_1}, \ldots, m_{i_r} . Write l = ag, and n = tg. If $r \ge 3$, one has $r/g \le a \le N/g$, but, unless g = 1, it is not admissible to have a < r, because all r edges of the graph should be traversed. For this reason, $r \le a \le N/g$. Result (3-2) follows. If r = 2, l is even and, for each l, only sequences of the form $(i_1, i_2, \ldots, i_1, i_2)$ with $n_{i_1} = n_{i_2} = l/2$ are possible. Put l = 2a, $a = 1, 2, \ldots, M = \min\{m_1, m_2\}$ to get (3-1).

Example 1. From (3-1), we have

$$\begin{aligned} & \mathcal{F}(1, 1) = \mathcal{F}(1, 2) = \mathcal{F}(2, 1) = \mathcal{F}(1, 3) = \mathcal{F}(3, 1) = 4 \\ & \mathcal{F}(2, 2) = 12, \\ & \mathcal{F}(1, 4) = \mathcal{F}(4, 1) = \mathcal{F}(1, 5) = \mathcal{F}(5, 1) = 4, \\ & \mathcal{F}(2, 3) = \mathcal{F}(3, 2) = 20, \\ & \mathcal{F}(2, 4) = \mathcal{F}(4, 2) = 28, \\ & \mathcal{F}(3, 3) = 172/3. \end{aligned}$$

From (3-3),

$$\begin{aligned} \theta(1,1) &= \theta(1,2) = \theta(2,1) = \theta(1,3) \\ &= \theta(3,1) = \theta(1,4) = \theta(4,1) = \theta(1,5) = \theta(5,1) = 4, \\ \theta(2,2) &= 10, \quad \theta(2,3) = \theta(3,2) = 20, \quad \theta(3,3) = 56. \end{aligned}$$

Example 2. From (3-2),

$$\begin{aligned} &\mathcal{F}(1,1,1) = 16, \\ &\mathcal{F}(1,1,2) = \mathcal{F}(1,2,1) = \mathcal{F}(2,1,1) = 32, \\ &\mathcal{F}(1,2,2) = \mathcal{F}(2,1,2) = \mathcal{F}(2,2,1) = 112, \\ &\mathcal{F}(1,1,3) = \mathcal{F}(1,3,1) = \mathcal{F}(3,1,1) = 48, \\ &\mathcal{F}(1,1,4) = \mathcal{F}(1,4,1) = \mathcal{F}(4,1,1) = 64, \\ &\mathcal{F}(1,2,3) = \mathcal{F}(3,1,2) = \mathcal{F}(2,3,1) = \mathcal{F}(3,2,1) = \mathcal{F}(1,3,2) = \mathcal{F}(2,1,3) = 256, \\ &\mathcal{F}(2,2,2) = 1056. \end{aligned}$$

From (3-3),

 $\begin{aligned} \theta(1, 1, 1) &= 16, \\ \theta(1, 1, 2) &= \theta(2, 1, 1) = \theta(1, 2, 1) = 32, \\ \theta(1, 2, 2) &= \theta(2, 1, 2) = \theta(2, 2, 1) = 112, \\ \theta(1, 1, 3) &= \theta(3, 1, 1) = \theta(1, 3, 1) = 48, \\ \theta(1, 1, 4) &= \theta(4, 1, 1) = \theta(1, 4, 1) = 64, \\ \theta(1, 2, 3) &= \theta(3, 1, 2) = \theta(2, 3, 1) = \theta(3, 2, 1) = \theta(1, 3, 2) = \theta(2, 1, 3) = 256, \\ \theta(2, 2, 2) &= 1048. \end{aligned}$

Remark. (a) Notice that θ , and likewise the Witt formula, is given in terms of the Möbius function. However, formula (3-3) counts closed nonperiodic paths traversing the edges of G_R in all directions (without backtracking) and, in that sense, generalizes the Witt formula. Also, our formula has an algebraic meaning of a dimension. See Section 4.

(b) If m_{i_1}, \ldots, m_{i_r} are coprime, $\mathcal{F} = \theta$. Otherwise, \mathcal{F} can be rational. For instance, $\mathcal{F}(3,3) = 172/3$. But $\mathcal{F}' := N\mathcal{F}$, $N = m_{i_1} + \cdots + m_{i_r}$, is always a positive integer which counts the number of words of length N. For example, in the case N = 4, $m_1 = m_2 = 2$, $\mathcal{F}' = 48$. The words are

$$D_1^{\pm 2} D_2^{\pm 2}, \quad D_1^{-1} D_2^{+1} D_1^{+1} D_2^{+1}, \quad D_1^{+1} D_2^{-1} D_1^{+1} D_2^{+1}, \quad D_1^{-1} D_2^{-1} D_1^{+1} D_2^{+1}, \\ D_1^{-1} D_2^{+1} D_1^{+1} D_2^{-1}, \quad D_1^{-1} D_2^{-1} D_1^{-1} D_2^{+1}, \quad D_1^{-1} D_2^{-1} D_1^{+1} D_2^{-1},$$

plus four cyclic permutations for each of them, and the four periodic words $(D_1^{\pm 1}D_2^{\pm 1})^2$ plus two cyclic permutations for each.

In terms of \mathcal{F}' ,

$$\theta(m_{i_1},\ldots,m_{i_r})=\frac{1}{N}\sum_{g\mid m_{i_1},\ldots,m_{i_r}}\mu(g)\mathcal{F}'\left(\frac{m_{i_1}}{g},\ldots,\frac{m_{i_r}}{g}\right).$$

Although the Möbius function is negative for some divisors g, the right hand side is nevertheless always a positive number because $\mathcal{F}'(m_{i_1}/g, \ldots, m_{i_r}/g)$ counts words in a subset of the words counted by $\mathcal{F}'(m_{i_1}, \ldots, m_{i_r})$.

(c) Given a circular necklace with N beads, consider the problem of counting inequivalent nonperiodic colorings of these beads with 2r colors $\{c_i, \bar{c}_i\}$, i = 1, ..., r, with m_i occurrences of the index i, $N = \sum m_i$, with the restriction that no two colors c_i and \bar{c}_i (same index) occur adjacent in a coloring. Now, consider an oriented graph with r loops hooked to a single vertex. Each loop edge corresponds to a color c_i . A nonperiodic closed nonbacktracking path of length N in the graph corresponds to a coloring, and a color \bar{c}_i corresponds to an edge being traversed in the opposite orientation. The presence of a single vertex in the graph reflects the fact that adjacent to a bead with, say, color c_i , any other with distinct index may follow. The number of inequivalent colorings is given by θ .

As a basic test of our counting ideas, we prove Sherman's statement [1962] relating the Witt formula to paths in G_R .

Proposition 3.2. Relative to graph G_R , formula (1-2) gives the number \mathcal{M} of equivalence classes of closed nonperiodic paths of length N > 0 that traverse each edge i counterclockwise $m_i \ge 0$ times (i = 1, 2, ..., R), where $m_1 + \cdots + m_R = N$.

Proof. Denote by $m_{i_1}, \ldots, m_{i_r}, r \leq R$, the nonzero entries in $\mathcal{M}(m_1, \ldots, m_R)$, which we call $\mathcal{M}_r(m_{i_1}, \ldots, m_{i_r})$. Words representing counterclockwise paths have positive exponents so that the factors 2^{2a} and 2^a in formulas (3-1) and (3-2) are not needed. Hence

(3-6)
$$\mathcal{M}_r(m_{i_1},\ldots,m_{i_r}) = \sum_{g|m_{i_1},\ldots,m_{i_r}} \frac{\mu(g)}{g} \mathcal{F}_c\left(\frac{m_{i_1}}{g},\ldots,\frac{m_{i_r}}{g}\right)$$

where

(3-7)
$$\mathscr{F}_{c}\left(\frac{m_{i_{1}}}{g}, \frac{m_{i_{2}}}{g}\right) = \sum_{a=1}^{M/g} \frac{1}{a} \binom{m_{i_{1}}/g - 1}{a - 1} \binom{m_{i_{2}}/g - 1}{a - 1}$$
 if $r = 2$,

with $M = \min\{m_{i_1}, m_{i_2}\}$, and

(3-8)
$$\mathscr{F}_{c}\left(\frac{m_{i_{1}}}{g},\ldots,\frac{m_{i_{r}}}{g}\right) = \sum_{a=r}^{N/g} \frac{1}{a} \sum_{\{S_{a}\}} \prod_{c=1}^{r} \binom{m_{i_{c}}/g-1}{t_{i_{c}}-1}$$
 if $r \ge 3$.

In the case r = 2 suppose $m_{i_1} \le m_{i_2}$. Using formula (A-3) from the Appendix (with l = 2), it follows that

$$\sum_{a=1}^{m_{i_1}/g} \frac{1}{a} \binom{m_{i_1}/g - 1}{a - 1} \binom{m_{i_2}/g - 1}{a - 1} = \frac{g}{m_{i_2}} \binom{m_{i_1}/g + m_{i_2}/g - 1}{m_{i_1}/g}$$
$$= \frac{(N/g)!}{(N/g)(m_{i_1}/g)!(m_{i_2}/g)!}.$$

Similarly if $m_{i_2} \le m_{i_1}$. In the case $r \ge 3$ define

(3-9)
$$I = \sum_{\substack{m_i > 0 \\ m_{i_1} + \dots + m_{i_r} = N}} \mathcal{F}_c\left(\frac{m_{i_1}}{g}, \dots, \frac{m_{i_r}}{g}\right).$$

Upon substituting (3-8) into (3-9) and exchanging the summation symbols, we get

$$I = \sum_{a=r}^{N/g} \frac{1}{a} \sum_{\{S_a\}} \sum_{\substack{m_i > 0 \\ m_{i_1} + \dots + m_{i_r} = N}} \prod_{c=1}^r \left(\frac{m_{i_c}}{t_{i_c}} - 1 \right).$$

Applying Lemma A.2,

$$I = \sum_{a=r}^{N/g} \frac{1}{a} \sum_{\{S_a\}} \binom{N/g - 1}{a - 1} = \sum_{a=r}^{N/g} \frac{1}{a} \binom{N/g - 1}{a - 1} r w_r(a)$$

where

$$rw(a) = \sum_{j=1}^{r} (-1)^{r+j} \binom{r}{j} (j-1)^a + (-1)^{a+r}$$

is the number of sequences in $\{S_a\}$ [da Costa and Variane 2005]. Using that

$$\sum_{a=r}^{N/g} \frac{1}{a} \binom{N/g-1}{a-1} (j-1)^a = \frac{g}{N} (j^{N/g} - 1)$$

and

$$\sum_{a=r}^{N/g} \frac{1}{a} \binom{N/g-1}{a-1} (-1)^{a+r} = (-1)^{r+1} \frac{g}{N},$$

we get

(3-10)
$$I = \frac{g}{N} \sum_{j=1}^{r} (-1)^{r+j} {r \choose j} j^{N/g}.$$

The Stirling numbers S(N/g, r) of the second kind are given by [Chen and Koh 1992]

(3-11)
$$S\left(\frac{N}{g}, r\right) = \frac{1}{r!} \sum_{k=0}^{r} (-1)^k \binom{r}{k} (r-k)^{N/g} = \frac{1}{r!} \sum_{j=0}^{r} (-1)^{r+j} \binom{r}{j} j^{N/g}$$

so that

(3-12)
$$I = r! \frac{g}{N} S\left(\frac{N}{g}, r\right).$$

Stirling numbers have the property that

(3-13)
$$\sum_{\substack{m_i > 0 \\ m_{i_1} + \dots + m_{i_r} = N}} \frac{(N/g)!}{(m_{i_1}/g)! \cdots (m_{i_r}/g)!} = r! S\left(\frac{N}{g}, r\right).$$

Comparing relations (3-12), (3-13), and (3-9),

(3-14)
$$\mathscr{F}_c\left(\frac{m_{i_1}}{g},\ldots,\frac{m_{i_r}}{g}\right) = \frac{g}{N} \frac{(N/g)!}{(m_1/g)!\cdots(m_r/g)!}$$

Upon substitution of (3-14) into (3-6), the result follows.

In the following we compute formulas for θ_+ and θ_- .

Theorem 3.3. Suppose any of the following conditions is satisfied:

(a)
$$N = m_{i_1} + \dots + m_{i_r} < 2r$$
.

- (b) m_{i_1}, \ldots, m_{i_r} are coprime.
- (c) m_{i_1}, \ldots, m_{i_r} are neither all odd nor all even.
- (d) m_{i_1}, \ldots, m_{i_r} are all odd.

Then

(3-15)
$$\theta_{-}(m_{i_1},\ldots,m_{i_r}) = \theta_{+}(m_{i_1},\ldots,m_{i_r}).$$

Proof. The proof is similar to that of [da Costa 1997, Theorem 1] and uses Lemma 2.1. \Box

Theorem 3.4. The number $\theta_+(m_{i_1}, \ldots, m_{i_r})$ is given by

(3-16)
$$\theta_{+}(m_{i_{1}},\ldots,m_{i_{r}}) = \sum_{\text{odd } g \mid m_{i_{1}},\ldots,m_{i_{r}}} \frac{\mu(g)}{g} \mathscr{G}\left(\frac{m_{i_{1}}}{g},\ldots,\frac{m_{i_{r}}}{g}\right)$$

where the summation is over all the common odd divisors of m_{i_1}, \ldots, m_{i_r} , and $\mathcal{G} = \mathcal{F}/2$ with \mathcal{F} as in (3-1) and (3-2). If m_{i_1}, \ldots, m_{i_r} are all even numbers, then

(3-17)
$$\theta_{-}(m_{i_1},\ldots,m_{i_r}) = \theta_{+}(m_{i_1},\ldots,m_{i_r}) - \theta_{+}\left(\frac{m_{i_1}}{2},\ldots,\frac{m_{i_r}}{2}\right).$$

Proof. First, suppose that all common divisors of m_{i_1}, \ldots, m_{i_r} are odd numbers. In this case,

$$\theta(m_{i_1},\ldots,m_{i_r})=\sum_{\text{odd }g\mid m_{i_1},\ldots,m_{i_r}}\frac{\mu(g)}{g}\mathscr{F}\left(\frac{m_{i_1}}{g},\ldots,\frac{m_{i_r}}{g}\right).$$

Since $\theta = \theta_+ + \theta_-$ and $\theta_+ = \theta_-$ (Theorem 3.3) it follows that $\theta = 2\theta_+$, hence

(3-18)
$$\theta_{+} = \frac{1}{2} \sum_{\text{odd } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F}\left(\frac{m_{i_1}}{g}, \dots, \frac{m_{i_r}}{g}\right).$$

If the numbers m_{i_1}, \ldots, m_{i_r} are all even, θ_+ is again given by (3-18) because, in this case, the m_i 's have common divisors which are even numbers, but since periodic words with even period have negative sign, only the odd divisors are relevant to getting θ_+ . The reason one should have the factor $\frac{1}{2}$ is that, by Lemma 2.1, when one considers the set of all possible words representing paths of a given length which traverse the edges of $G_r m_{i_1}, \ldots, m_{i_r}$ times, half of them have positive sign and the other half have negative sign. To account for the positive half, one needs the factor $\frac{1}{2}$. Let's now compute θ_- in the even case. Write

$$\begin{aligned} \theta &= \sum_{\text{odd } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F} + \sum_{\text{even } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F} \\ &= \frac{1}{2} \sum_{\text{odd } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F} + \frac{1}{2} \sum_{\text{odd } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F} + \sum_{\text{even } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F} \\ &= 2\theta_+ + \sum_{\text{even } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F}. \end{aligned}$$

Using that $\theta = \theta_+ + \theta_-$, we obtain

$$\theta_{-} = \theta_{+} + \sum_{\text{even } g \mid m_{i_1}, \dots, m_{i_r}} \frac{\mu(g)}{g} \mathcal{F}.$$

Now the relevant even divisors are $\{2n\}$ where *n* are the odd common divisors of $\{m_i\}$. For the other possible divisors, if any, use that $\mu(2^j n) = 0$, $j \ge 2$. Using the equality $\mu(2n) = -\mu(n)$, we see that the summation over the even divisors is equal to $-\theta_+(m_{i_1}/2, \ldots, m_{i_r}/2)$, proving the result.

Remark. Like θ , the numbers θ_{\pm} can be interpreted as the number of inequivalent nonperiodic colorings of a circular necklace with *N* beads. However, now these colorings are classified as positive or negative according to formula (2-3). It is positive (negative) if the number N + l + T + s is odd (even). In this case, *s* is the number \bar{c} of colors present in a coloring. Interpret *T* in terms of the color indices.

Definition. Let s_1, \ldots, s_r be arbitrary positive integers. Let the number \mathcal{P} be defined as follows. If s_1, \ldots, s_r are all even numbers,

(3-19)
$$\mathscr{P}(s_1,\ldots,s_r) = \sum_{\text{even } g \mid s_1,\ldots,s_r} \frac{\mu(g)}{g} \mathscr{G}\left(\frac{s_1}{g},\ldots,\frac{s_r}{g}\right),$$

otherwise, $\mathcal{P}(s_1, \ldots, s_r) = 0$. Also, define

(3-20)
$$\mathcal{H}(s_1, \dots, s_r)$$
 if s_1, \dots, s_r not all even,

$$= \begin{cases} \mathcal{G}(s_1, \dots, s_r) & \text{if } s_1, \dots, s_r \text{ not all even,} \\ \mathcal{G}(s_1, \dots, s_r) - \sum_{k \mid s_1, \dots, s_r} \frac{1}{k} \mathcal{P}\left(\frac{s_1}{k}, \dots, \frac{s_r}{k}\right) & \text{otherwise.} \end{cases}$$

Lemma 3.5.
$$\mathcal{P} = \sum_{g \mid s_1, \dots, s_r} \frac{\mu(g)}{g} (\mathcal{G} - \mathcal{H}).$$

Proof. From the definition, $\mathcal{G} = \mathcal{H}$ if s_1, \ldots, s_r are not all even. Otherwise,

(3-21)
$$\mathscr{G} - \mathscr{H} = \sum_{g|s_1,\dots,s_r} \frac{1}{g} \mathscr{P}\left(\frac{s_1}{g},\dots,\frac{s_r}{g}\right).$$

Now apply Lemma A.1 to get the result.

Theorem 3.6.
$$\theta_+(m_{i_1},\ldots,m_{i_r}) = \sum_{g \mid m_{i_1},\ldots,m_{i_r}} \frac{\mu(g)}{g} \mathcal{H}\left(\frac{m_{i_1}}{g},\ldots,\frac{m_{i_r}}{g}\right).$$

Proof. When m_{i_1}, \ldots, m_{i_r} are not all even, their odd divisors are the only possible common divisors. In this case, $\mathcal{P} = 0$ and

(3-22)
$$\theta_{+} = \sum_{\text{odd } g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} \mathcal{H},$$

with $\mathcal{H} = \mathcal{G}$. If m_{i_1}, \ldots, m_{i_r} are all even, the sum over odd divisors of m_{i_1}, \ldots, m_{i_r} can be expressed as

$$\begin{aligned} \theta_{+} &= \sum_{\text{odd } g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} \mathcal{G} \\ &= \sum_{g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} \mathcal{G} - \sum_{\text{even } g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} \mathcal{G} = \sum_{g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} \mathcal{G} - \mathcal{P} \\ &= \sum_{g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} \mathcal{G} - \sum_{g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} (\mathcal{G} - \mathcal{H}) = \sum_{g \mid m_{i_{1}}, \dots, m_{i_{r}}} \frac{\mu(g)}{g} \mathcal{H}. \end{aligned}$$

Example 3.

$$\begin{aligned} \theta_{\pm}(1,1) &= \theta_{\pm}(1,2) = \theta_{\pm}(2,1) = \theta_{\pm}(1,3) = \theta_{\pm}(3,1) = \theta_{\pm}(1,4) = \theta_{\pm}(4,1) \\ &= \theta_{\pm}(1,5) = \theta_{\pm}(5,1) = 2, \\ \theta_{+}(2,2) &= 6, \qquad \theta_{-}(2,2) = 4, \qquad \theta_{\pm}(2,3) = \theta_{\pm}(3,2) = 10, \qquad \theta_{+}(2,4) = 14, \\ \theta_{-}(2,4) &= 12, \qquad \theta_{+}(4,2) = 14, \qquad \theta_{-}(4,2) = 12, \qquad \theta_{\pm}(3,3) = 28. \end{aligned}$$

Example 4.

 $\begin{aligned} \theta_{\pm}(1,1,1) &= 8, \\ \theta_{\pm}(1,1,2) &= \theta_{\pm}(2,1,1) = \theta_{\pm}(1,2,1) = 16, \\ \theta_{\pm}(1,2,2) &= \theta_{\pm}(2,1,2) = \theta_{\pm}(2,2,1) = 56, \\ \theta_{\pm}(1,1,3) &= \theta_{\pm}(3,1,1) = \theta_{\pm}(1,3,1) = 24, \\ \theta_{\pm}(1,1,4) &= \theta_{\pm}(4,1,1) = \theta_{\pm}(1,4,1) = 32, \\ \theta_{\pm}(1,2,3) &= \theta_{\pm}(3,1,2) = \theta_{\pm}(2,3,1) = \theta_{\pm}(3,2,1) = \theta_{\pm}(1,3,2) = \theta_{\pm}(2,1,3) = 128, \\ \theta_{+}(2,2,2) &= 524, \\ \theta_{-}(2,2,2) &= 516. \end{aligned}$

4. Sherman identity and Lie algebras

In this section we relate our previous results with Lie algebras and solve Sherman's problem. The solution is provided by the following proposition.

Proposition 4.1 [Kang and Kim 1999]. Let $V = \bigoplus_{(k_1,...,k_r) \in \mathbb{Z}'_{>0}} V_{(k_1,...,k_r)}$ be a $\mathbb{Z}'_{>0}$ -graded vector space over \mathbb{C} with dim $V_{(k_1,...,k_r)} = d(k_1,...,k_r) < \infty$, for all $(k_1,...,k_r) \in \mathbb{Z}'_{>0}$, and let

$$L = \bigoplus_{(k_1, \dots, k_r) \in \mathbb{Z}_{>0}^r} L_{(k_1, \dots, k_r)}$$

be the free Lie algebra generated by V. Then the dimensions of the subspaces $L_{(k_1,...,k_r)}$ are given by

(4-1)
$$\dim L_{(k_1,...,k_r)} = \sum_{g \mid (k_1,...,k_r)} \frac{\mu(g)}{g} \mathcal{W}\left(\frac{k_1}{g},\ldots,\frac{k_r}{g}\right)$$

where summation is over all common divisors g of k_1, \ldots, k_r and W is given by

(4-2)
$$\mathscr{W}(k_1,\ldots,k_r) = \sum_{s\in T(k_1,\ldots,k_r)} \frac{(|s|-1)!}{s!} \prod_{i_1,\ldots,i_r=1}^{\infty} d(i_1,\ldots,i_r)^{s_{i_1,\ldots,i_r}}.$$

The exponents $s_{i_1,...,i_r}$ are the components of $s \in T$,

(4-3)
$$T(k_1, \dots, k_r) = \{s = (s_{i_1, \dots, i_r}) | s_{i_1, \dots, i_r} \in \mathscr{Z}_{\geq 0}, \\ \sum_{i_1, \dots, i_r=1}^{\infty} s_{i_1, \dots, i_r}(i_1, \dots, i_r) = (k_1, \dots, k_r)\},$$

and

(4-4)
$$|s| = \sum_{i_1,\dots,i_r=1}^{\infty} s_{i_1,\dots,i_r}, \quad s! = \prod_{i_1,\dots,i_r=1}^{\infty} s_{i_1,\dots,i_r}!.$$

Moreover, the numbers dim $L_{(k_1,...,k_r)}$ *satisfy*

(4-5)
$$\prod_{k_1,\dots,k_r=1}^{\infty} (1 - z_1^{k_1} \cdots z_r^{k_r})^{\dim L_{(k_1,\dots,k_r)}} = 1 - f(z_1,\dots,z_r)$$

where

(4-6)
$$f(z_1, \ldots, z_r) := \sum_{k_1, \ldots, k_r=1}^{\infty} d(k_1, \ldots, k_r) z_1^{k_1} \cdots z_r^{k_r}.$$

This function is associated with the generating function of the W's,

(4-7)
$$g(z_1, \ldots, z_r) := \sum_{k_1, \ldots, k_r=1}^{\infty} \mathcal{W}(k_1, \ldots, k_r) z_1^{k_1} \cdots z_r^{k_r},$$

by the relation

(4-8)
$$e^{-g} = 1 - f.$$

Identity (4-5) is a consequence of the famous *Poincaré–Birkhoff–Witt theorem* for the free Lie algebra. Computation of the formal logarithm of the left-hand side of (4-5) and its expansion gives that the infinite product equals the exponential in (4-8). Raise both members of (4-5) to the power -1, compute the formal logarithm of both members, and expand them. Identification of the coefficients of the same order, definition (4-2), and application of the Möbius inversion give (4-1). See [Kang and Kim 1999] for details. In [da Costa and Variane 2005], (4-1) is called the *generalized Witt formula*, W is called the *Witt partition function*, and (4-5) the *generalized Witt identity*.

Formulas (3-3) and (3-20) have exactly the form of (4-1) with corresponding *Witt partition functions* given by \mathcal{F} and \mathcal{H} , respectively, so we interpret θ and θ_+ as giving the dimensions of the homogeneous spaces of graded Lie algebras. In each case, the algebra is generated by a graded vector space whose dimensions can be computed recursively from (4-2) as a function of the Witt partition function.

However, a general formula can be obtained from (4-8) using (4-6) as the formal Taylor expansion of $1 - e^{-g}$. This gives

(4-9)
$$d(k_1, \dots, k_r) = \frac{1}{k_1! \cdots k_r!} \frac{\partial^{|k|}}{\partial z_1^{k_1} \cdots \partial z_r^{k_r}} (1 - e^{-g})|_{z_1 = \dots = z_r = 0}$$

with

(4-10)
$$g(z_1, \ldots, z_r) := \sum_{k_1, \ldots, k_r=1}^{\infty} \mathcal{W}(k_1, \ldots, k_r) z^{k_1} \cdots z^{k_r}$$

and $\mathcal{W} = \mathcal{F}$, \mathcal{H} given by (3-1), (3-2), and (3-20). Furthermore, dim $L_{(k_1,...,k_r)} = \theta$, θ_+ given by (3-3), 3.6 satisfy the *generalized Witt identity* (4-5) with the corresponding dimensions given by (4-9). In fact, an explicit formula for (4-9) can be derived:

Theorem 4.2. The numbers $d(k_1, \ldots, k_r)$ are given by the formula

(4-11)
$$d(k_1,\ldots,k_r) = \sum_{\lambda=1}^{|k|} (-1)^{\lambda+1} \sum_{p(\lambda,k)} \prod_{i=1}^q \frac{[{}^{\mathfrak{W}}(l_{i1},\ldots,l_{ir})]^{a_i}}{a_i!},$$

where $|k| = k_1 + \dots + k_r$, $q = -1 + \prod_{i=1}^r (k_i + 1)$, $p_{\lambda,k}$ is the set of all $a_i \in \{0, 1, 2, \dots\}$ such that $\sum_{i=1}^q a_i = \lambda$, $\sum_{i=1}^q a_i l_{ij} = k_j$, and the vectors $l_i = (l_{i1}, \dots, l_{ir})$, l_{ij} satisfying $0 \le l_{ij} \le k_j$, $\forall j = 1, \dots, r$, $\forall i = 1, \dots, q$ and $\sum_{j=1}^r l_{ij} > 0$. Set $\mathcal{W}(l_i) = 0$ if $l_{ij} = 0$ for some j; otherwise, \mathcal{W} is the Witt partition function.

Proof. A generalization of Faà di Bruno's relation [Constantine and Savits 1996; Savits 2006] gives a formula for the |k|-th derivative of the exponential of a function $g(z_1, \ldots, z_r)$. From this formula and (4-9), (4-11) follows.

Example 5. We compute d(2, 2) explicitly. In this case, $k_1 = k_2 = 2$, |k| = 4, q = 8. The possible vectors $l \le (2, 2)$ are $l_1 = (0, 1)$, $l_2 = (1, 0)$, $l_3 = (1, 1)$, $l_4 = (0, 2)$, $l_5 = (2, 0)$, $l_6 = (2, 1)$, $l_7 = (1, 2)$, and $l_8 = (2, 2)$. Next we give the values of $a_1, \ldots, a_8 \ge 0$ satisfying

$$\sum_{i=1}^{8} a_i = \lambda, \quad \sum_{i=1}^{8} a_i l_i = (2, 2).$$

Define the vector $a = (a_1, \ldots, a_8)$. The possible *a*'s for each λ are

for $\lambda = 1$, a = (0, ..., 0, 1); for $\lambda = 2$, a = (0, 1, 0, 0, 0, 0, 1, 0), (0, 0, 2, 0, 0, 0, 0, 0), (0, 0, 0, 1, 1, 0, 0, 0), (1, 0, 0, 0, 0, 1, 0, 0); for $\lambda = 3$, a = (0, 2, 0, 1, 0, 0, 0, 0), (2, 0, 0, 0, 1, 0, 0, 0), (1, 1, 1, 0, 0, 0, 0, 0, 0); for $\lambda = 4$, a = (2, 2, 0, 0, 0, 0, 0, 0). We get

$$d(2,2) = \mathcal{W}(2,2) - \frac{1}{2}\mathcal{W}(1,1)^2.$$

The dimensions up to d(3, 3) are

$$\begin{split} N &= 2, \quad d(1,1) = \mathcal{W}(1,1), \\ N &= 3, \quad d(1,2) = \mathcal{W}(1,2), \quad d(2,1) = \mathcal{W}(2,1), \\ N &= 4, \quad d(1,3) = \mathcal{W}(1,3), \quad d(3,1) = \mathcal{W}(3,1), \\ d(2,2) &= \mathcal{W}(2,2) - \frac{1}{2}\mathcal{W}(1,1)^2, \\ N &= 5, \quad d(1,4) = \mathcal{W}(1,4), \quad d(4,1) = \mathcal{W}(4,1), \\ d(2,3) &= \mathcal{W}(2,3) - \mathcal{W}(1,1)\mathcal{W}(1,2), \\ d(3,2) &= \mathcal{W}(3,2) - \mathcal{W}(1,1)\mathcal{W}(2,1), \\ N &= 6, \quad d(1,5) = \mathcal{W}(1,5), \quad d(5,1) = \mathcal{W}(5,1), \\ d(2,4) &= \mathcal{W}(2,4) - \mathcal{W}(1,1)\mathcal{W}(1,3) - \frac{1}{2}\mathcal{W}(1,2)^2, \\ d(4,2) &= \mathcal{W}(4,2) - \mathcal{W}(1,1)\mathcal{W}(3,1) - \frac{1}{2}\mathcal{W}(2,1)^2 \\ d(3,3) &= \mathcal{W}(3,3) - \mathcal{W}(1,1)\mathcal{W}(2,2) - \mathcal{W}(1,2)\mathcal{W}(2,1) + \frac{1}{6}\mathcal{W}(1,1)^3. \end{split}$$

For r = 3, the dimensions up to d(2, 2, 2) are

$$\begin{split} N &= 3, \quad d(1,1,1) = \mathcal{W}(1,1,1), \\ N &= 4, \quad d(1,1,2) = \mathcal{W}(1,1,2), \quad d(1,2,1) = \mathcal{W}(1,2,1), \quad d(2,1,1) = \mathcal{W}(2,1,1), \\ N &= 5, \quad d(1,2,2) = \mathcal{W}(1,2,2), \quad d(2,1,2) = \mathcal{W}(2,1,2), \quad d(2,2,1) = \mathcal{W}(2,2,1), \\ \quad d(1,1,3) = \mathcal{W}(1,1,3), \quad d(1,3,1) = \mathcal{W}(1,3,1), \quad d(3,1,1) = \mathcal{W}(3,1,1), \\ N &= 6, \quad d(1,1,4) = \mathcal{W}(1,1,4), \quad d(1,4,1) = \mathcal{W}(1,4,1), \quad d(4,1,1) = \mathcal{W}(4,1,1), \\ \quad d(1,2,3) = \mathcal{W}(1,2,3), \quad d(3,1,2) = \mathcal{W}(3,1,2), \quad d(2,3,1) = \mathcal{W}(2,3,1), \\ \quad d(3,2,1) = \mathcal{W}(3,2,1), \quad d(1,3,2) = \mathcal{W}(1,3,2), \quad d(2,1,3) = \mathcal{W}(2,1,3), \\ \quad d(2,2,2) = \mathcal{W}(2,2,2) - \frac{1}{2}\mathcal{W}^2(1,1,1). \end{split}$$

Example 6. Relative to θ with $\mathcal{W} = \mathcal{G}$ and applying data from previous examples, for the case r = 2, we find the dimensions

$$d(1, 1) = d(1, 2) = d(2, 1) = d(1, 3) = d(3, 1) = d(1, 4) = d(4, 1) = d(2, 3)$$
$$= d(3, 2) = d(1, 5) = d(5, 1) = d(2, 2) = d(2, 4) = d(4, 2) = d(3, 3) = 4.$$

In the case r = 3, the dimensions are

$$d(1,1,1) = 8,$$

$$d(1,1,2) = d(2,1,1) = d(1,2,1) = 16, \quad d(1,2,2) = d(2,1,2) = d(2,2,1) = 56,$$

$$d(1,1,3) = d(3,1,1) = d(1,3,1) = 24, \quad d(1,1,4) = d(4,1,1) = d(1,4,1) = 32,$$

$$d(1,2,3) = d(3,1,2) = d(2,3,1) = d(3,2,1) = d(1,3,2) = d(2,1,3) = 128,$$

$$d(2,2,2) = 496.$$

Example 7. Relative to θ_+ with $\mathcal{W} = \mathcal{H}$, we find for the case r = 2

$$d(1, 1) = d(1, 2) = d(2, 1) = d(1, 3) = d(3, 1) = d(1, 4) = d(4, 1) = d(2, 3)$$

= d(3, 2) = d(1, 5) = d(5, 1) = 2,
d(2, 2) = 5, d(2, 4) = d(4, 2) = 9, d(3, 3) = 28,
and for r = 3
d(1, 1, 1) = 8,
d(1, 1, 2) = d(2, 1, 1) = d(1, 2, 1) = 16, d(1, 2, 2) = d(2, 1, 2) = d(2, 2, 1) = 56,
d(1, 1, 3) = d(3, 1, 1) = d(1, 3, 1) = 24, d(1, 1, 4) = d(4, 1, 1) = d(1, 4, 1) = 32,
d(1, 2, 3) = d(3, 1, 2) = d(2, 3, 1) = d(3, 2, 1) = d(1, 3, 2) = d(2, 1, 3) = 128,
d(2, 2, 2) = 504.

Remark. In spite of the negative terms in the formulas for the dimensions, they give positive results. To understand why, consider, for example, the case

$$d(2,2) = \mathcal{W}(2,2) - \frac{1}{2}\mathcal{W}(1,1)^2$$

with $W(a, b) = \mathcal{F}' = (a + b)\mathcal{F}$. So d(2, 2) is four times the result in example 6. In the set of words counted by $\mathcal{F}'(2, 2) = 48$ there is a subset whose elements are words that are obtained by gluing together the words in the set counted by W(1, 1) = 8. The gluing produces an overcounting which is corrected by the one half factor. So d(2, 2) is positive. The same argument can be used to get positivity for the other formulas.

Theorem 4.3. For each $G_r \subseteq G_R$, we have

(4-12)
$$\prod_{m_{i_1},\dots,m_{i_r}=1}^{\infty} (1+z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_+} = e^{-g(z_{i_1}^2,\dots,z_{i_r}^2)+g(z_{i_1},\dots,z_{i_r})},$$

and

(4-13)
$$\prod_{m_{i_1},\dots,m_{i_r}=1}^{\infty} (1-z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_-} = e^{+g(z_{i_1}^2,\dots,z_{i_r}^2)-g(z_{i_1},\dots,z_{i_r})}.$$

Proof. To prove (4-12), multiply and divide its left-hand side by

$$\prod_{m_{i_1},\ldots,m_{i_r}=1}^{\infty} (1-z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_+}$$

and use (4-8). To get (4-13), write

$$\prod_{m_{i_1},\dots,m_{i_r}=1}^{\infty} (1-z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_-} = \prod_{N=r}^{\infty} \prod_{\substack{m_i>0\\m_{i_1}+\cdots+m_{i_r}=N}} (1-z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_-}.$$

Decompose the product over *N* into three products, namely, one over all N < 2r, one over all even $N \ge 2r$, and another over all odd N > 2r. Then apply Theorems 3.3 and 3.4 and formula (3-17).

$$\prod_{m_{i_1},\dots,m_{i_r}=1}^{\infty} (1+z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_+} (1-z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_-} = 1$$

Proof. Multiply (4-12) and (4-13).

The left side of (1-1) equals

$$\prod_{j=1}^{R} (1+z_j)^2 \prod_{r=2}^{R} \prod_{G_r} \prod_{m_{i_1},\dots,m_{i_r}>0} (1+z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_+} (1-z_{i_1}^{m_{i_1}}\cdots z_{i_r}^{m_{i_r}})^{\theta_-}.$$

The Sherman identity now follows from Theorem 4.4.

Appendix

Lemma A.1. If

(A-1)
$$g(n_1,\ldots,n_k) = \sum_{d\mid n_1,\ldots,n_k} \frac{\mu(d)}{d} f\left(\frac{n_1}{d},\ldots,\frac{n_k}{d}\right),$$

then

(A-2)
$$f(n_1, ..., n_k) = \sum_{d \mid n_1, ..., n_k} \frac{1}{d} g\left(\frac{n_1}{d}, ..., \frac{n_1}{d}\right).$$

Proof. Set $G(n_1, ..., n_k) := (n_1 + \dots + n_k)g(n_1, \dots, n_k)$ and

$$F\left(\frac{n_1}{d},\ldots,\frac{n_k}{d}\right) := \left(\frac{n_1}{d}+\cdots+\frac{n_k}{d}\right) f\left(\frac{n_1}{d},\ldots,\frac{n_k}{d}\right).$$

Then (A-1) can be expressed in the form

$$G(n_1,\ldots,n_k)=\sum_{d\mid n_1,\ldots,n_k}\mu(d)F\bigg(\frac{n_1}{d},\ldots,\frac{n_1}{d}\bigg).$$

Möbius inversion gives

$$F(n_1,\ldots,n_k) = \sum_{d\mid n_1,\ldots,n_k} G\left(\frac{n_1}{d},\ldots,\frac{n_k}{d}\right).$$

Therefore,

$$(n_1 + \dots + n_k) f(n_1, \dots, n_k) = \sum_{d \mid n_1, \dots, n_k} \left(\frac{n_1}{d} + \dots + \frac{n_k}{d} \right) g\left(\frac{n_1}{d}, \dots, \frac{n_k}{d} \right). \quad \Box$$

The converse is also true.

Lemma A.2. Let $N \ge \alpha = n_1 + \cdots + n_l, n_1, \ldots, n_l, n_i > 0$, be a partition of α . Then

(A-3)
$$\sum_{\sum_{i=1}^{l} k_i = N} \prod_{i=1}^{l} \binom{k_i - 1}{n_i - 1} = \binom{N - 1}{\alpha - 1}$$

with the convention that a bracket in the left side is zero whenever $k_i < n_i$.

Proof. Using

$$\frac{q^{\alpha}}{(1-q)^{\alpha}} = \sum_{N=\alpha}^{\infty} \binom{N-1}{\alpha-1} q^N,$$

it follows that

$$\frac{q^{\alpha}}{(1-q)^{\alpha}} = \prod_{i=1}^{l} \frac{q^{n_i}}{(1-q)^{n_i}} = \prod_{i=1}^{l} \sum_{k_i=n_i}^{\infty} \binom{k_i-1}{n_i-1} q^{k_i} = \sum_{N=\alpha}^{\infty} \sum_{\substack{k_i \ge n_i \\ \sum_{i=1}^{l} k_i=N}} \prod_{i=1}^{l} \binom{k_i-1}{n_i-1} q^N.$$

Comparison with the previous expression and the convention gives the result. \Box

Acknowledgements

We thank Professor Peter Moree (Max Planck Institute for Mathematics, Bonn) and Professor Thomas Ward (University of East Anglia, UK) for email correspondence regarding the positivity of the Möbius inversion formula.

References

[Apostol 1976] T. M. Apostol, *Introduction to analytic number theory*, Springer, New York, 1976. MR 55 #7892 Zbl 0335.10001

[Chen and Koh 1992] C. C. Chen and K. M. Koh, *Principles and techniques in combinatorics*, World Scientific, River Edge, NJ, 1992. MR 93h:05001 Zbl 0786.05002

[[]Andrews 1976] G. E. Andrews, *The theory of partitions*, Encyclopedia of Mathematics and its Applications **2**, Addison-Wesley, Reading, MA, 1976. MR 58 #27738 Zbl 0371.10001

- [Cimasoni 2010] D. Cimasoni, "A generalized Kac–Ward formula", J. Stat. Mech. 2010 (2010), Article ID #P07023.
- [Constantine and Savits 1996] G. M. Constantine and T. H. Savits, "A multivariate Faà di Bruno formula with applications", *Trans. Amer. Math. Soc.* **348**:2 (1996), 503–520. MR 96g:05008 Zbl 0846.05003
- [da Costa 1997] G. A. T. F. da Costa, "Feynman identity: a special case, I", *J. Math. Phys.* **38**:2 (1997), 1014–1034. MR 98b:82012 Zbl 0869.05039
- [da Costa and Variane 2005] G. A. T. F. da Costa and J. Variane, Jr., "Feynman identity: a special case revisited", *Lett. Math. Phys.* **73**:3 (2005), 221–235. MR 2007a:82012 Zbl 1101.82007
- [Kang and Kim 1999] S.-J. Kang and M.-H. Kim, "Dimension formula for graded Lie algebras and its applications", *Trans. Amer. Math. Soc.* **351**:11 (1999), 4281–4336. MR 2000b:17009 Zbl 0926.17002
- [Loebl 2004] M. Loebl, "A discrete non-Pfaffian approach to the Ising problem", pp. 145–154 in *Graphs, morphisms and statistical physics* (Piscataway, NJ, 2001), edited by J. Nešetřil and P. Winkler, DIMACS Ser. Discrete Math. Theoret. Comput. Sci. 63, Amer. Math. Soc., Providence, RI, 2004. MR 2005k:05066 Zbl 1069.82003
- [Moree 2005] P. Moree, "The formal series Witt transform", *Discrete Math.* **295**:1-3 (2005), 143–160. MR 2006b:05015 Zbl 1064.05025
- [Savits 2006] T. H. Savits, "Some statistical applications of Faa di Bruno", *J. Multivariate Anal.* **97**:10 (2006), 2131–2140. MR 2008g:62144 Zbl 1101.62041
- [Sherman 1960] S. Sherman, "Combinatorial aspects of the Ising model for ferromagnetism, I: A conjecture of Feynman on paths and graphs", *J. Math. Phys.* **1** (1960), 202–217. MR 22 #10273 Zbl 0123.45501
- [Sherman 1962] S. Sherman, "Combinatorial aspects of the Ising model for ferromagnetism, II: An analogue to the Witt identity", *Bull. Amer. Math. Soc.* **68** (1962), 225–229. MR 25 #2003
- [Witt 1937] E. Witt, "Treue Darstellung Liescher Ring", J. Reine Angew. Math. 177 (1937), 152–160. JFM 63.0089.02

Received May 5, 2012.

G. A. T. F. DA COSTA DEPARTAMENTO DE MATEMÁTICA Universidade Federal de Santa Catarina 88040-900 Florianópolis SC Brazil

gatcosta@mtm.ufsc.br

G. A. ZIMMERMANN DEPARTAMENTO DE MATEMÁTICA UNIVERSIDADE FEDERAL DE SANTA CATARINA 88040-900 FLORIANÓPOLIS SC BRAZIL

graciele@ifsc.edu.br

ON THE CLASSIFICATION OF STABLE SOLUTIONS TO BIHARMONIC PROBLEMS IN LARGE DIMENSIONS

JUNCHENG WEI, XINGWANG XU AND WEN YANG

We give a new bound on the exponent for nonexistence of stable solutions to the biharmonic problem $\Delta^2 u = u^p$ in \mathbb{R}^n , where u > 0, p > 1, and $n \ge 20$.

1. Introduction

Of concern is the biharmonic equation

(1-1)
$$\Delta^2 u = u^p, \quad u > 0 \quad \text{in } \mathbb{R}^n$$

where $n \ge 5$ and p > 1. Set

(1-2)
$$\Lambda_u(\varphi) := \int_{\mathbb{R}^n} |\Delta \varphi|^2 dx - p \int_{\mathbb{R}^n} u^{p-1} \varphi^2 dx \quad \text{for all } \varphi \in H^2(\mathbb{R}^n).$$

The Morse index $\operatorname{ind}(u)$ of a classical solution to (1-1) is defined as the maximal dimension of all subspaces of $H^2(\mathbb{R}^n)$ such that $\Lambda_u(\varphi) < 0$ in $H^2(\mathbb{R}^n) \setminus \{0\}$. We say u is a stable solution to (1-1) if $\Lambda_u(\varphi) \ge 0$ for any test function $\varphi \in H^2(\mathbb{R}^n)$; that is, if the Morse index is zero.

In the first part of the paper, we obtain the following classification result on stable solutions of (1-1).

Theorem 1.1. Let $n \ge 20$ and 1 . Then (1-1) has no stable solutions.

Here p^* stands for the smallest real root greater than $\frac{n-4}{n-8}$ of the algebraic equation

$$512(2-n)x^{6} + 4(n^{3} - 60n^{2} + 670n - 1344)x^{5} - 2(13n^{3} - 424n^{2} + 3064n - 5408)x^{4} + 2(27n^{3} - 572n^{2} + 3264n - 5440)x^{3} - (49n^{3} - 772n^{2} + 3776n - 5888)x^{2} + 4(5n^{3} - 66n^{2} + 288n - 416)x - 3(n^{3} - 12n^{2} + 48n - 64) = 0.$$

MSC2010: primary 35B20; secondary 35J60.

Keywords: stable solutions, biharmonic superlinear equations.

The first author was supported from an earmarked grant ("On Elliptic Equations with Negative Exponents") from RGC of Hong Kong.

Some remarks are in order. Let us recall that for the second-order problem

(1-3)
$$\Delta u + u^p = 0 \quad u > 0 \text{ in } \mathbb{R}^n, \ p > 1,$$

Farina gave a complete classification of all finite Morse index solutions. The main result of [Farina 2007] is that no stable solution exists to (1-3) if either $n \le 10$, p > 1 or $n \ge 11$, $p < p_{JL}$. Here p_{JL} denotes the Joseph–Lundgren exponent [Gui et al. 1992]. On the other hand, a stable radial solution exists for $p \ge p_{JL}$.

For the fourth-order case, the nonexistence of positive solutions to (1-1) is shown if $p < \frac{n+4}{n-4}$, and all entire solutions are classified if $p = \frac{n+4}{n-4}$. See [Lin 1998; Wei and Xu 1999]. When $p > \frac{n+4}{n-4}$, radially symmetric solutions to (1-1) are completely classified in [Ferrero et al. 2009; Gazzola and Grunau 2006; Guo and Wei 2010]. The radial solutions are shown to be stable if and only if $p \ge p'_{JL}$ and $n \ge 13$, where p'_{JL} stands for the corresponding Joseph–Lundgren exponent (see [Ferrero et al. 2009; Gazzola and Grunau 2006]). In the general nonradial case, Wei and Ye [Wei and Ye 2010] showed the nonexistence of stable or finite Morse index solutions when either $n \le 8$, p > 1 or $n \ge 9$, $p \le \frac{n}{n-8}$. In dimensions $n \ge 9$, a perturbation argument is used to show the nonexistence of stable solutions for $p < \frac{n}{n-8} + \varepsilon_n$ for some $\varepsilon_n > 0$. However, no explicit value of ε_n was given. The proof of Wei and Ye [2010] follows an earlier idea of Cowan, Esposito and Ghoussoub [2010] in which a similar problem in a bounded domain was studied. Theorem 1.1 gives an explicit value on ε_n for $n \ge 20$.

In the second-order case, the proof of Farina uses basically the Moser iterations: namely multiply (1-3) by the power of u, like u^q , q > 1. Moser iteration works because of the following simple identity

$$\int_{\mathbb{R}^n} u^q (-\Delta u) = \frac{4q}{(q+1)^2} \int_{\mathbb{R}^n} |\nabla u|^{\frac{q+1}{2}} |^2, \, \forall u \in C_0^1(\mathbb{R}^n).$$

In the fourth-order case, such equality does not hold, and in fact we have

$$\int_{\mathbb{R}^n} u^q (\Delta^2 u) = \frac{4q}{(q+1)^2} \int_{\mathbb{R}^n} |\Delta u^{\frac{q+1}{2}}|^2 - q(q-1)^2 \int_{\mathbb{R}^n} u^{q-3} |\nabla u|^4, \forall u \in C_0^2(\mathbb{R}^n).$$

The additional term $\int_{\mathbb{R}^n} u^{q-3} |\nabla u|^4$ makes the Moser iteration argument difficult to use. Wei and Ye [2010] used instead the new test function $-\Delta u$ and showed that $\int_{\mathbb{R}^2} |\Delta u|^2$ is bounded. Thus the exponent $\frac{n}{n-8}$ is obtained. In this paper, we use the Moser iteration for the fourth-order problem and give a control on the term $\int_{\mathbb{R}^n} u^{q-3} |\nabla u|^4$ (Lemma 2.3). As a result, we obtain a better exponent $\frac{n}{n-8} + \varepsilon_n$ where ε_n is explicitly given. As far as we know, this seems to be the first result for Moser iteration for a fourth-order problem.

In the second part of this paper, we show that the same idea can be used to establish the regularity of extremal solutions to

(1-4)
$$\begin{cases} \Delta^2 u = \lambda (u+1)^p, \ \lambda > 0 \ \text{in } \Omega, \\ u > 0 \qquad \qquad \text{in } \Omega, \\ u = \Delta u = 0 \qquad \qquad \text{on } \partial \Omega. \end{cases}$$

where Ω is a smooth and bounded convex domain in \mathbb{R}^n .

For problem (1-4), it is known [Berchio and Gazzola 2005] that for $p > \frac{n+4}{n-4}$ there exists a critical value $\lambda^* > 0$ depending on p > 1 and Ω such that

- If $\lambda \in (0, \lambda^*)$, (1-4) has a minimal and classical solution which is stable;
- If λ = λ*, a unique weak solution, called the extremal solution u* exists for (1-4);
- No weak solution of (1-4) exists whenever $\lambda > \lambda^*$.

The regularity of the extremal solution of problem (1-4) at $\lambda = \lambda_*$ has been studied in [Cowan et al. 2010; Wei and Ye 2010], where it was shown that the extremal solution is bounded provided $n \le 8$ or $p < \frac{n}{n-8} + \varepsilon_n$, $n \ge 9$ (ε_n very small). Here, we also give a explicit bound for the exponent p in large dimensions and our second result is the following.

Theorem 1.2. The extremal solution u^* of (1-4) when $\lambda = \lambda^*$ is bounded provided that $n \ge 20$ and $1 , where <math>p^*$ is defined as above.

As $n \to +\infty$, the value ε_n is asymptotically $8\sqrt{8/3}/(n-8)^{3/2}$ and thus the upper bound for *p* has the expansion

(1-5)
$$1 + \frac{8}{n-8} + \frac{8\sqrt{8/3}}{(n-8)^{3/2}} + O\left(\frac{1}{(n-8)^2}\right).$$

On the other hand, for radial solutions, the Joseph–Lundgren exponent [Gui et al. 1992] has the following asymptotic expansion

(1-6)
$$1 + \frac{8}{n-8} + \frac{16}{(n-8)^{3/2}} + O\left(\frac{1}{(n-8)^2}\right).$$

In this paper, we have only considered fourth-order problems with power-like nonlinearity. Other kinds of nonlinearity, such as exponential and negative powers, also appear in many applications; see [Cowan et al. 2010]. However, our technique here yields no improvements of results of that reference in the case of exponential and negative nonlinearities.

This paper is organized as follows. We prove Theorem 1.1 and Theorem 1.2 respectively in Section 2 and Section 3. Some technical inequalities are given in the Appendix.

2. Proof of Theorem 1.1

Lemma 2.1. For any $\varphi \in C_0^4(\mathbb{R}^n)$ with $\varphi \ge 0$, any $\gamma > 1$ and $\varepsilon > 0$ an arbitrary small number, we have

(2-1)
$$\int_{\mathbb{R}^{n}} (\Delta(u^{\gamma}\varphi^{\gamma}))^{2} \leq \int_{\mathbb{R}^{n}} ((\Delta u^{\gamma}\varphi^{\gamma})^{2} + \varepsilon |\nabla u|^{4}\varphi^{2\gamma}u^{2\gamma-4} + Cu^{2\gamma} \|\nabla^{4}(\varphi^{2\gamma})\|),$$

$$(2-2) \quad \int_{\mathbb{R}^{n}} (\Delta(u^{\gamma}\varphi^{\gamma}))^{2} \ge \int_{\mathbb{R}^{n}} ((\Delta u^{\gamma}\varphi^{\gamma})^{2} - \varepsilon |\nabla u|^{4} \varphi^{2\gamma} u^{2\gamma-4} - C u^{2\gamma} ||\nabla^{4}(\varphi^{2\gamma})||),$$

$$(2-3) \quad \int ((u^{\gamma})_{i\,i})^{2} \varphi^{2\gamma} < \int ((u^{\gamma}\varphi^{\gamma})_{i\,i})^{2} + \varepsilon \int |\nabla u|^{4} u^{2\gamma-4} \varphi^{2\gamma}$$

$$(2-3) \quad \int_{\mathbb{R}^n} ((u^{\gamma})_{ij})^2 \varphi^{2\gamma} \leq \int_{\mathbb{R}^n} ((u^{\gamma} \varphi^{\gamma})_{ij})^2 + \varepsilon \int_{\mathbb{R}^n} |\nabla u|^4 u^{2\gamma - 4} \varphi^{2\gamma} + C \int_{\mathbb{R}^n} u^{2\gamma} \|\nabla^4(\varphi^{2\gamma})\|_{\mathcal{H}^{1,2}}$$

where *C* is a positive number that only depends on γ and ε , and $\|\nabla^4(\varphi^{2\gamma})\|$ is defined by

$$\|\nabla^4(\varphi^{2\gamma})\|^2 = \varphi^{-2\gamma} |\nabla\varphi^{\gamma}|^4 + |\varphi^{\gamma}(\Delta^2 \varphi^{\gamma})| + |\nabla^2 \varphi^{\gamma}|^2.$$

In the following, unless said otherwise, the constant *C* always denotes a positive number which may change term by term but only depends on γ , ε .

Proof. Since φ is compactly supported, we can use integration by parts without considering the boundary terms. First, by direct calculation, we get

$$(2-4) \quad (\Delta(u^{\gamma}\varphi^{\gamma}))^{2} = [(\Delta u^{\gamma})\varphi^{\gamma}]^{2} + 4\nabla u^{\gamma}\nabla\varphi^{\gamma}\Delta\varphi^{\gamma}u^{\gamma} + 4\nabla u^{\gamma}\nabla\varphi^{\gamma}\Delta u^{\gamma}\varphi^{\gamma} + 4(\nabla u^{\gamma}\nabla\varphi^{\gamma})^{2} + 2\Delta u^{\gamma}u^{\gamma}\Delta\varphi^{\gamma}\varphi^{\gamma} + u^{2\gamma}(\Delta\varphi^{\gamma})^{2}.$$

We now need to deal with the third and fifth terms on the right side of this equality, up to the integration of both sides.

For the third term, we have

$$\begin{split} \int_{\mathbb{R}^n} \Delta u^{\gamma} \nabla u^{\gamma} \nabla \varphi^{\gamma} \varphi^{\gamma} &= -\int_{\mathbb{R}^n} (u^{\gamma})_i (u^{\gamma})_{ij} (\varphi^{\gamma})_j \varphi^{\gamma} \\ &- \int_{\mathbb{R}^n} (u^{\gamma})_i (u^{\gamma})_j (\varphi^{\gamma})_{ij} \varphi^{\gamma} - \int_{\mathbb{R}^n} (u^{\gamma})_i (u^{\gamma})_j (\varphi^{\gamma})_j (\varphi^{\gamma})_i, \end{split}$$

where $f_i = \partial f / \partial x_i$ and $f_{ij} = \partial^2 f / \partial x_j \partial x_i$. (Here and in the sequel, we use the Einstein summation convention, so for example $\partial_i (u_i u_j \varphi_j) = \sum_{1 \le i, j \le n} \partial_i (u_i u_j \varphi_j)$.) The first term on the right side of the previous equation can be estimated as

$$2\int_{\mathbb{R}^n} (u^{\gamma})_i (u^{\gamma})_{ij} (\varphi^{\gamma})_j \varphi^{\gamma} = \int_{\mathbb{R}^n} \partial_j ((u^{\gamma})_i (u^{\gamma})_i (\varphi^{\gamma})_j \varphi^{\gamma}) - \int_{\mathbb{R}^n} ((u^{\gamma})_i)^2 (\varphi^{\gamma})_{jj} \varphi^{\gamma} - \int_{\mathbb{R}^n} ((u^{\gamma})_i)^2 (\varphi^{\gamma})_j (\varphi^{\gamma$$

Combining these two equalities, we get

$$2\int_{\mathbb{R}^{n}}\Delta u^{\gamma}\nabla u^{\gamma}\nabla \varphi^{\gamma}\varphi^{\gamma} = -\int_{\mathbb{R}^{n}}\partial_{j}\left((u^{\gamma})_{i}(u^{\gamma})_{i}(\varphi^{\gamma})_{j}\varphi^{\gamma}\right) \\ -\int_{\mathbb{R}^{n}}2(u^{\gamma})_{i}(u^{\gamma})_{j}(\varphi^{\gamma})_{ij}\varphi^{\gamma} - \int_{\mathbb{R}^{n}}2(u^{\gamma})_{i}(u^{\gamma})_{j}(\varphi^{\gamma})_{j}(\varphi^{\gamma})_{i} \\ +\int_{\mathbb{R}^{n}}((u^{\gamma})_{i})^{2}(\varphi^{\gamma})_{jj}\varphi^{\gamma} + \int_{\mathbb{R}^{n}}((u^{\gamma})_{i})^{2}(\varphi^{\gamma})_{j}(\varphi^{\gamma})_{j}.$$

Rewriting this equality we have

$$(2-5) \quad 4\int_{\mathbb{R}^n} \Delta u^{\gamma} \nabla u^{\gamma} \nabla \varphi^{\gamma} \varphi^{\gamma} = 2\int_{\mathbb{R}^n} |\nabla u^{\gamma}|^2 \Delta \varphi^{\gamma} \varphi^{\gamma} + 2\int_{\mathbb{R}^n} |\nabla u^{\gamma}|^2 |\nabla \varphi^{\gamma}|^2 -4\int_{\mathbb{R}^n} (u^{\gamma})_i (u^{\gamma})_j (\varphi^{\gamma})_{ij} \varphi^{\gamma} - 4\int_{\mathbb{R}^n} \langle \nabla u^{\gamma}, \nabla \varphi^{\gamma} \rangle^2.$$

For the fifth term on the right side of (2-4) we have

(2-6)
$$\int_{\mathbb{R}^{n}} \Delta u^{\gamma} u^{\gamma} \Delta \varphi^{\gamma} \varphi^{\gamma} = -\int_{\mathbb{R}^{n}} u^{\gamma} \langle \nabla u^{\gamma}, \nabla (\Delta \varphi^{\gamma}) \rangle \varphi^{\gamma} - \int_{\mathbb{R}^{n}} \langle \nabla u^{\gamma}, \nabla \varphi^{\gamma} \rangle u^{\gamma} \Delta \varphi^{\gamma} - \int_{\mathbb{R}^{n}} |\nabla u^{\gamma}|^{2} \Delta \varphi^{\gamma} \varphi^{\gamma}.$$

Combining (2-4), (2-5) and (2-6), one obtains

$$(2-7) \quad \int_{\mathbb{R}^{n}} (\Delta(u^{\gamma}\varphi^{\gamma}))^{2} - \int_{\mathbb{R}^{n}} (\Delta u^{\gamma})^{2} \varphi^{2\gamma}$$

$$= 2 \int_{\mathbb{R}^{n}} |\nabla u^{\gamma}|^{2} |\nabla \varphi^{\gamma}|^{2} - 4 \int_{\mathbb{R}^{n}} \varphi^{\gamma} (\nabla^{2}\varphi^{\gamma} (\nabla u^{\gamma}, \nabla u^{\gamma}))$$

$$+ \int_{\mathbb{R}^{n}} u^{2\gamma} \varphi^{\gamma} \Delta^{2} (\varphi^{\gamma}) - 2 \int_{\mathbb{R}^{n}} u^{2\gamma} (\Delta \varphi^{\gamma})^{2}.$$

Now by the Young equality, for any $\varepsilon > 0$, there exists a constant $C = C(\gamma, \varepsilon)$ such that

$$|\nabla u^{\gamma}|^{2}|\nabla \varphi^{\gamma}|^{2} \leq \frac{\varepsilon}{4}|\nabla u^{\gamma}|^{4}u^{-2\gamma}\varphi^{2\gamma} + C|\nabla \varphi^{\gamma}|^{4}u^{2\gamma}\varphi^{-2\gamma}$$

and

$$|\varphi^{\gamma}(\nabla^{2}\varphi^{\gamma}(\nabla u^{\gamma},\nabla u^{\gamma}))| \leq \frac{\varepsilon}{8} |\nabla u^{\gamma}|^{4} u^{-2\gamma} \varphi^{2\gamma} + C u^{2\gamma} |\nabla^{2}\varphi^{\gamma}|^{2}.$$

Thus by (2-7), together with the two estimates above, one gets

$$\left|\int_{\mathbb{R}^n} (\Delta(u^{\gamma}\varphi^{\gamma}))^2 - \int_{\mathbb{R}^n} (\Delta u^{\gamma})^2 \varphi^{2\gamma}\right| \leq \varepsilon \int_{\mathbb{R}^n} |\nabla u^{\gamma}|^4 u^{-2\gamma} \varphi^{2\gamma} + 6C \int_{\mathbb{R}^n} u^{2\gamma} \|\nabla^4 \varphi^{\gamma}\|^2.$$

The estimates (2-1) and (2-2) follow from this easily.

Next we observe that $|\nabla^2 u^{\gamma}|^2 \varphi^{2\gamma} = \left[\frac{1}{2}\Delta |\nabla u^{\gamma}|^2 - \langle \nabla u^{\gamma}, \nabla \Delta u^{\gamma} \rangle\right] \varphi^{2\gamma}$. Thus up to the integration by parts, with the help of (2-5) and the estimates we just proved, the estimate (2-3) also follows by noticing the identity $\int_{\mathbb{R}^n} (\Delta (u^{\gamma} \varphi^{\gamma}))^2 = \int_{\mathbb{R}^n} |\nabla^2 (u^{\gamma} \varphi^{\gamma})|^2$. The proof of Lemma 2.1 is thus completed.

Let us return to the equation

(2-8)
$$\Delta^2 u = u^p, \quad u > 0 \text{ in } \mathbb{R}^n.$$

Fix $q = 2\gamma - 1 > 0$ and $\gamma > 1$. Let $\varphi \in C_0^{\infty}(\mathbb{R}^n)$. Multiplying (2-8) by $u^q \varphi^{2\gamma}$ and integration by parts, we obtain

(2-9)
$$\int_{\mathbb{R}^n} \Delta u \Delta(u^q \varphi^{2\gamma}) = \int_{\mathbb{R}^n} u^{p+q} \varphi^{2\gamma}.$$

For the left side of (2-9), we have:

Lemma 2.2. For any $\varphi \in C_0^{\infty}(\mathbb{R}^n)$ with $\varphi \ge 0$, for any $\varepsilon > 0$ and γ with q defined above, there exists a positive constant C depends on γ , ε such that

$$(2-10) \quad \int_{\mathbb{R}^n} \frac{\gamma^2}{q} \Delta u \Delta (u^q \varphi^{2\gamma}) \ge \int_{\mathbb{R}^n} (\Delta u^\gamma \varphi^\gamma)^2 - \int_{\mathbb{R}^n} C u^{2\gamma} \|\nabla^4(\varphi^{2\gamma})\| \\ - \int_{\mathbb{R}^n} (\gamma^2(\gamma-1)^2 + \varepsilon) u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma}.$$

Proof. First, by direct computations, we obtain

$$\begin{split} \Delta u \Delta (u^{2\gamma-1}\varphi^{2\gamma}) &= \Delta u \big((2\gamma-1)u^{2\gamma-2} \Delta u \varphi^{2\gamma} + 2(2\gamma-1)u^{2\gamma-2} \nabla u \nabla (\varphi^{2\gamma}) \\ &+ (2\gamma-1)(2\gamma-2)u^{2\gamma-3} |\nabla u|^2 \varphi^{2\gamma} + u^{2\gamma-1} \Delta \varphi^{2\gamma} \big), \\ (\Delta u^{\gamma}\varphi^{\gamma})^2 &= \gamma^2 u^{2\gamma-2} (\Delta u)^2 \varphi^{2\gamma} + \gamma^2 (\gamma-1)^2 u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma} \\ &+ 2(\gamma-1)\gamma^2 u^{2\gamma-3} |\nabla u|^2 \Delta u \varphi^{2\gamma}. \end{split}$$

Combining these two identities, we get

$$(2-11) \quad \frac{\gamma^2}{q} \Delta u \Delta (u^q \varphi^{2\gamma}) = (\Delta u^\gamma \varphi^\gamma)^2 + 2\gamma^2 u^{2\gamma-2} \Delta u \nabla u \nabla \varphi^{2\gamma} + \frac{\gamma^2}{q} u^{2\gamma-1} \Delta u \Delta \varphi^{2\gamma} - \gamma^2 (\gamma-1)^2 u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma}.$$

For the term $u^{2\gamma-2}\Delta u\nabla u\nabla \varphi^{2\gamma}$, we have

$$u^{2\gamma-2}\Delta u \nabla u \nabla \varphi^{2\gamma} = \partial_i (u^{2\gamma-2} u_i u_j (\varphi^{2\gamma})_j) - (2\gamma-2) u^{2\gamma-3} (u_i)^2 u_j (\varphi^{2\gamma})_j - u^{2\gamma-2} u_i u_{ij} (\varphi^{2\gamma})_j - u^{2\gamma-2} u_i u_j (\varphi^{2\gamma})_{ij}.$$

We can regroup the term $u^{2\gamma-2}u_iu_{ii}(\varphi^{2\gamma})_i$ as

$$2u^{2\gamma-2}u_{i}u_{ij}(\varphi^{2\gamma})_{j} = \partial_{j}(u^{2\gamma-2}(u_{i})^{2}(\varphi^{2\gamma})_{j}) - (2\gamma-2)u^{2\gamma-3}u_{j}(u_{i})^{2}(\varphi^{2\gamma})_{j} - u^{2\gamma-2}(u_{i})^{2}(\varphi^{2\gamma})_{jj}.$$

Therefore we get

$$(2-12) \quad 2u^{2\gamma-2}\Delta u \nabla u \nabla \varphi^{2\gamma} = 2\partial_i (u^{2\gamma-2}u_i u_j (\varphi^{2\gamma})_j) - \partial_j (u^{2\gamma-2}(u_i)^2 (\varphi^{2\gamma})_j) - (2\gamma-2)u^{2\gamma-3}(u_i)^2 u_j (\varphi^{2\gamma})_j + u^{2\gamma-2}(u_i)^2 (\varphi^{2\gamma})_{jj} - 2u^{2\gamma-2}u_i u_j (\varphi^{2\gamma})_{ij}.$$

For the last three terms on the right side of (2-12), applying Young's inequality, we get

$$\begin{split} |u^{2\gamma-3}(u_{i})^{2}u_{j}(\varphi^{2\gamma})_{j}| &\leq \frac{\varepsilon}{6\gamma^{2}(\gamma-1)}u^{2\gamma-4}|\nabla u|^{4}\varphi^{2\gamma} + Cu^{2\gamma}\|\nabla^{4}(\varphi^{2\gamma})\|,\\ |u^{2\gamma-2}(u_{i})^{2}(\varphi^{2\gamma})_{jj}| &\leq \frac{\varepsilon}{6\gamma^{2}}u^{2\gamma-4}|\nabla u|^{4}\varphi^{2\gamma} + Cu^{2\gamma}\|\nabla^{4}(\varphi^{2\gamma})\|,\\ |u^{2\gamma-2}u_{i}u_{j}(\varphi^{2\gamma})_{ij}| &\leq \frac{\varepsilon}{6\gamma^{2}}u^{2\gamma-4}|\nabla u|^{4}\varphi^{2\gamma} + Cu^{2\gamma}\|\nabla^{4}(\varphi^{2\gamma})\|. \end{split}$$

These three inequalities and (2-12) imply

$$(2-13) \quad \int_{\mathbb{R}^n} 2\gamma^2 u^{2\gamma-2} \Delta u \nabla u \nabla \varphi^{2\gamma} \ge -\frac{\varepsilon}{2} \int_{\mathbb{R}^n} u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma} - C \int_{\mathbb{R}^n} u^{2\gamma} \|\nabla^4(\varphi^{2\gamma})\|.$$

Similarly we get

$$(2-14) \int_{\mathbb{R}^n} \frac{\gamma^2}{q} u^{2\gamma-1} \Delta u \Delta \varphi^{2\gamma} \ge -\frac{\varepsilon}{2} \int_{\mathbb{R}^n} u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma} - C \int_{\mathbb{R}^n} u^{2\gamma} ||\nabla^4(\varphi^{2\gamma})||.$$

Inequality (2-10) follows from (2-11), (2-13) and (2-14).

Inequality (2-10) follows from (2-11), (2-13) and (2-14).

As a result of (2-1) and (2-10), we have

$$(2-15) \quad \int_{\mathbb{R}^n} \frac{\gamma^2}{q} \Delta u \Delta (u^q \varphi^{2\gamma}) \ge \int_{\mathbb{R}^n} (\Delta (u^\gamma \varphi^\gamma))^2 - \int_{\mathbb{R}^n} C u^{2\gamma} \|\nabla^4 (\varphi^{2\gamma})\| \\ - \int_{\mathbb{R}^n} (\gamma^2 (\gamma - 1)^2 + \varepsilon) u^{2\gamma - 4} |\nabla u|^4 \varphi^{2\gamma}.$$

Next we estimate the most difficult term, $\int_{\mathbb{R}^n} u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma}$, in (2-15). This is the key step in proving Theorem 1.1.

Lemma 2.3. If u is the classical solution to the biharmonic equation (2-8), and φ is defined as above, then for any sufficiently small $\varepsilon > 0$, we have the following inequality

$$(2-16) \quad \left(\frac{1}{2} - \varepsilon\right) \int_{\mathbb{R}^n} u^{2\gamma - 4} |\nabla u|^4 \varphi^{2\gamma} \le \frac{2}{\gamma^2} \int_{\mathbb{R}^n} (\Delta(u^{\gamma} \varphi^{\gamma}))^2 + \int_{\mathbb{R}^n} C u^{2\gamma} ||\nabla^4(\varphi^{2\gamma})|| \\ - \int_{\mathbb{R}^n} \frac{4}{(4\gamma - 3 + p)(p+1)} u^{2\gamma + p - 1} \varphi^{2\gamma}.$$

Proof. It is easy to see that

(2-17)
$$\int_{\mathbb{R}^n} u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma} = \frac{1}{\gamma^4} \int_{\mathbb{R}^n} u^{-2\gamma} |\nabla u^{\gamma}|^4 \varphi^{2\gamma},$$

and

$$(2-18) \int_{\mathbb{R}^{n}} u^{-2\gamma} |\nabla u^{\gamma}|^{4} \varphi^{2\gamma} = \int_{\mathbb{R}^{n}} u^{-2\gamma} |\nabla u^{\gamma}|^{2} \nabla u^{\gamma} \nabla u^{\gamma} \varphi^{2\gamma} = \int_{\mathbb{R}^{n}} -\nabla u^{-\gamma} |\nabla u^{\gamma}|^{2} \nabla u^{\gamma} \varphi^{2\gamma} = \int_{\mathbb{R}^{n}} u^{-\gamma} |\nabla u^{\gamma}|^{2} \Delta u^{\gamma} \varphi^{2\gamma} + \int_{\mathbb{R}^{n}} u^{-\gamma} \nabla (|\nabla u^{\gamma}|^{2}) \nabla u^{\gamma} \varphi^{2\gamma} + \int_{\mathbb{R}^{n}} u^{-\gamma} |\nabla u^{\gamma}|^{2} \nabla u^{\gamma} \nabla \varphi^{2\gamma},$$

where in the last step we used integration by parts. For the first term in the last part of this equality, we have

$$\int_{\mathbb{R}^n} u^{-\gamma} |\nabla u^{\gamma}|^2 \Delta u^{\gamma} \varphi^{2\gamma} = \gamma^3 \int_{\mathbb{R}^n} ((\gamma - 1)u^{2\gamma - 4} |\nabla u|^4 \varphi^{2\gamma} + u^{2\gamma - 3} |\nabla u|^2 \Delta u \varphi^{2\gamma}).$$

Substituting this into (2-18) and combining with (2-17), we obtain

$$(2-19) \quad \int_{\mathbb{R}^n} u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma} = \int_{\mathbb{R}^n} \frac{1}{\gamma^3} u^{-\gamma} \nabla (|\nabla u^{\gamma}|^2) \nabla u^{\gamma} \varphi^{2\gamma} + \int_{\mathbb{R}^n} u^{2\gamma-3} (|\nabla u|^2) \Delta u \varphi^{2\gamma} + \int_{\mathbb{R}^n} \frac{1}{\gamma^3} u^{-\gamma} (|\nabla u^{\gamma}|^2) \nabla u^{\gamma} \nabla \varphi^{2\gamma}.$$

The first term on the right side of (2-19) can be estimated as

$$(2-20) \quad u^{-\gamma} \nabla (|\nabla u^{\gamma}|^{2}) \nabla u^{\gamma} = 2u^{-\gamma} ((u^{\gamma})_{ij}(u^{\gamma})_{i}(u^{\gamma})_{j})$$

$$\leq 2\gamma (u^{\gamma})_{ij}(u^{\gamma})_{ij} + \frac{u^{-2\gamma}}{2\gamma} (u^{\gamma})_{i}(u^{\gamma})_{j}(u^{\gamma})_{i}(u^{\gamma})_{j}$$

$$= 2\gamma |\nabla^{2}u^{\gamma}|^{2} + \frac{u^{-2\gamma}}{2\gamma} |\nabla u^{\gamma}|^{4}.$$

As a consequence, we have

$$(2-21) \int_{\mathbb{R}^{n}} \frac{1}{\gamma^{3}} u^{-\gamma} \nabla (|\nabla u^{\gamma}|^{2}) \nabla u^{\gamma} \varphi^{2\gamma}$$

$$\leq \int_{\mathbb{R}^{n}} \frac{2}{\gamma^{2}} |\nabla^{2} u^{\gamma}|^{2} \varphi^{2\gamma} + \int_{\mathbb{R}^{n}} \frac{1}{2\gamma^{4}} u^{-2\gamma} |\nabla u^{\gamma}|^{4} \varphi^{2\gamma}$$

$$\leq \int_{\mathbb{R}^{n}} \frac{2}{\gamma^{2}} |\nabla^{2} (u^{\gamma} \varphi^{\gamma})|^{2} + \int_{\mathbb{R}^{n}} C u^{2\gamma} ||\nabla^{4} (\varphi^{2\gamma})|| + \int_{\mathbb{R}^{n}} \frac{1+4\gamma^{2}\varepsilon}{2\gamma^{4}} u^{-2\gamma} |\nabla u^{\gamma}|^{4} \varphi^{2\gamma}$$

$$= \int_{\mathbb{R}^{n}} \frac{2}{\gamma^{2}} (\Delta (u^{\gamma} \varphi^{\gamma}))^{2} + \int_{\mathbb{R}^{n}} C u^{2\gamma} ||\nabla^{4} (\varphi^{2\gamma})|| + \int_{\mathbb{R}^{n}} \frac{1+4\gamma^{2}\varepsilon}{2\gamma^{4}} u^{-2\gamma} |\nabla u^{\gamma}|^{4} \varphi^{2\gamma},$$

where we used (2-3) in the last step.

For the second term on the right side of (2-19), applying estimate (2.3) from [Wei and Ye 2010], that is, $(\Delta u)^2 \ge \frac{2}{p+1}u^{p+1}$, and the fact that $\Delta u < 0$ from Theorem 3.1 in [Wei and Xu 1999] or Theorem 2.1 in [Xu 2000], we have

$$(2-22) \quad \int_{\mathbb{R}^{n}} u^{2\gamma-3} (|\nabla u|^{2}) \Delta u \varphi^{2\gamma} \leq -\int_{\mathbb{R}^{n}} \sqrt{\frac{2}{p+1}} u^{2\gamma-3+\frac{p+1}{2}} (|\nabla u|^{2}) \varphi^{2\gamma}$$
$$= \int_{\mathbb{R}^{n}} \frac{\sqrt{\frac{2}{p+1}}}{2\gamma-2+\frac{p+1}{2}} u^{2\gamma-2+\frac{p+1}{2}} \Delta u \varphi^{2\gamma}$$
$$+ \int_{\mathbb{R}^{n}} \frac{\sqrt{\frac{2}{p+1}}}{2\gamma-2+\frac{p+1}{2}} u^{2\gamma-2+\frac{p+1}{2}} \nabla u \nabla \varphi^{2\gamma}.$$

Using the inequality $-\Delta u \ge \sqrt{\frac{2}{p+1}} u^{\frac{p+1}{2}}$, we get

$$(2-23) \quad \int_{\mathbb{R}^n} \frac{\sqrt{\frac{2}{p+1}}}{2\gamma - 2 + \frac{2}{p+1}} u^{2\gamma - 2 + \frac{p+1}{2}} \Delta u \varphi^{2\gamma} \le -\int_{\mathbb{R}^n} \frac{\frac{2}{p+1}}{2\gamma - 2 + \frac{p+1}{2}} u^{2\gamma + p - 1} \varphi^{2\gamma}.$$

On the other hand, for the second term on the right side of (2-22), we have

(2-24)
$$\int_{\mathbb{R}^{n}} u^{2\gamma - 2 + \frac{p+1}{2}} \nabla u \nabla \varphi^{2\gamma} = -\int_{\mathbb{R}^{n}} \frac{1}{L} u^{2\gamma - 1 + \frac{p+1}{2}} \Delta \varphi^{2\gamma}$$
$$= -\int_{\{x \mid \Delta \varphi^{2\gamma} > 0\}} \frac{1}{L} u^{2\gamma - 1 + \frac{p+1}{2}} \Delta \varphi^{2\gamma}$$
$$-\int_{\{x \mid \Delta \varphi^{2\gamma} \le 0\}} \frac{1}{L} u^{2\gamma - 1 + \frac{p+1}{2}} \Delta \varphi^{2\gamma},$$

where the first equality follows from integration by parts and $L = 2\gamma - 1 + \frac{p+1}{2}$. As for the first term on the last part of (2-24), using the inequality

$$\Delta u \leq -\sqrt{\frac{2}{p+1}}u^{\frac{p+1}{2}} < 0,$$

we have

(2-25)
$$\frac{\sqrt{\frac{p+1}{2}}}{L} \int_{\{x \mid \Delta \varphi^{2\gamma} > 0\}} u^{2\gamma - 1} \Delta u \,\Delta \varphi^{2\gamma} \le -\int_{\{x \mid \Delta \varphi^{2\gamma} > 0\}} \frac{1}{L} u^{2\gamma - 1 + \frac{p+1}{2}} \Delta \varphi^{2\gamma}.$$

Similarly to the proof of Lemma 2.1, it is easy to get

$$\left|\int_{\{x\mid\Delta\varphi^{2\gamma}>0\}}\frac{\sqrt{\frac{p+1}{2}}}{L}u^{2\gamma-1}\Delta u\Delta\varphi^{2\gamma}\right|\leq\varepsilon\int_{\mathbb{R}^n}u^{2\gamma-4}|\nabla u|^4\varphi^{2\gamma}+\int_{\mathbb{R}^n}Cu^{2\gamma}\|\nabla^4(\varphi^{2\gamma})\|.$$

From this and (2-25), we have

$$\left| \int_{\{x \mid \Delta \varphi^{2\gamma} > 0\}} \frac{1}{L} u^{2\gamma - 1 + \frac{p+1}{2}} \Delta \varphi^{2\gamma} \right| \le \varepsilon \int_{\mathbb{R}^n} u^{2\gamma - 4} |\nabla u|^4 \varphi^{2\gamma} + \int_{\mathbb{R}^n} C u^{2\gamma} \|\nabla^4(\varphi^{2\gamma})\|.$$

Similarly, we also obtain

$$\left|\int_{\{x\mid\Delta\varphi^{2\gamma}\leq 0\}}\frac{1}{L}u^{2\gamma-1+\frac{p+1}{2}}\Delta\varphi^{2\gamma}\right|\leq \varepsilon\int_{\mathbb{R}^n}u^{2\gamma-4}|\nabla u|^4\varphi^{2\gamma}+\int_{\mathbb{R}^n}Cu^{2\gamma}\|\nabla^4(\varphi^{2\gamma})\|.$$

From the last two inequalities and (2-24), we have

$$(2-26) \left| \int_{\mathbb{R}^n} u^{2\gamma-2+\frac{p+1}{2}} \nabla u \nabla \varphi^{2\gamma} \right| \le \varepsilon \int_{\mathbb{R}^n} u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma} + \int_{\mathbb{R}^n} C u^{2\gamma} \|\nabla^4(\varphi^{2\gamma})\|.$$

Combining (2-22), (2-23) and (2-26), we get the inequality

$$(2-27) \quad \int_{\mathbb{R}^{n}} u^{2\gamma-3} |\nabla u|^{2} \Delta u \varphi^{2\gamma} \leq \varepsilon \int_{\mathbb{R}^{n}} u^{2\gamma-4} |\nabla u|^{4} \varphi^{2\gamma} + \int_{\mathbb{R}^{n}} C u^{2\gamma} \|\nabla^{4}(\varphi^{2\gamma})\| \\ - \int_{\mathbb{R}^{n}} \frac{4}{(4\gamma - 3 + p)(p+1)} u^{2\gamma + p - 1} \varphi^{2\gamma}.$$

Finally, we apply Young's inequality to the third term on the right side of (2-19), and get

$$(2-28) \quad \int_{\mathbb{R}^n} \frac{1}{\gamma^3} u^{-\gamma} (|\nabla u^{\gamma}|^2) \nabla u^{\gamma} \nabla \varphi^{2\gamma} = \int_{\mathbb{R}^n} u^{2\gamma-3} |\nabla u|^2 \nabla u \nabla (\varphi^{2\gamma}) \leq \varepsilon \int_{\mathbb{R}^n} u^{2\gamma-4} |\nabla u|^4 \varphi^{2\gamma} + \int_{\mathbb{R}^n} C u^{2\gamma} ||\nabla^4(\varphi^{2\gamma})||.$$

By (2-19), (2-21), (2-27) and (2-28), we finally obtain

$$\begin{split} \left(\frac{1}{2} - \varepsilon\right) \int_{\mathbb{R}^n} u^{2\gamma - 4} |\nabla u|^4 \varphi^{2\gamma} &\leq \frac{2}{\gamma^2} \int_{\mathbb{R}^n} (\Delta (u^\gamma \varphi^\gamma))^2 + \int_{\mathbb{R}^n} C u^{2\gamma} \|\nabla^4 (\varphi^{2\gamma})\| \\ &- \int_{\mathbb{R}^n} \frac{4}{(4\gamma - 3 + p)(p+1)} u^{2\gamma + p - 1} \varphi^{2\gamma}. \quad \Box \end{split}$$

By (2-9), (2-15) and (2-16), since the number ε is arbitrary small in those three places, we have, for $\delta > 0$ sufficiently small,

$$(2-29) \int_{\mathbb{R}^{n}} \left(1 - 4(\gamma - 1)^{2} - \delta \right) (\Delta(u^{\gamma} \varphi^{\gamma}))^{2} \\ - \int_{\mathbb{R}^{n}} \left(\frac{\gamma^{2}}{2\gamma - 1} - \frac{8\gamma^{2}(\gamma - 1)^{2}}{(4\gamma - 3 + p)(p + 1)} \right) u^{p + 2\gamma - 1} \varphi^{2\gamma} \leq \int_{\mathbb{R}^{n}} C_{\delta} u^{2\gamma} \| \nabla^{4}(\varphi^{2\gamma}) \|,$$

where C_{δ} is a positive constant that depends on δ only. Here, we need to require $1 - 4(\gamma - 1)^2 > 0$, since we have assumed that $\gamma > 1$ in Lemma 2.1. So γ is required be in $(1, \frac{3}{2})$. If we can choose δ small enough to make $1 - 4(\gamma - 1)^2 - \delta$ positive, by the stability property of function u, we obtain

(2-30)
$$\int_{\mathbb{R}^n} (E - p\delta) u^{p+q} \varphi^{2\gamma} \leq \int_{\mathbb{R}^n} C_\delta u^{2\gamma} \| \nabla^4(\varphi^{2\gamma}) \|,$$

where E is defined to be

(2-31)
$$E = p(1 - 4(\gamma - 1)^2) - \frac{\gamma^2}{q} + \frac{8\gamma^2(\gamma - 1)^2}{(4\gamma - 3 + p)(p + 1)}$$

Now we take $\varphi = \eta^m$ with *m* sufficiently large, and choose η a cut-off function satisfying $0 \le \eta \le 1$, $\eta = 1$ for |x| < R and $\eta = 0$ for |x| > 2R. By Young's inequality again, we have

$$(2-32) \quad \int_{\mathbb{R}^n} u^{2\gamma} \|\nabla^4(\varphi^{2\gamma})\| \le C_{\delta} R^{-4} \int_{\mathbb{R}^n} u^{2\gamma} \eta^{2\gamma m-4} \\ \le C_{\delta,\varepsilon} R^{-\frac{4}{1-\theta}} \int_{\mathbb{R}^n} u^2 \eta^{2\gamma m-\frac{4}{1-\theta}} + \varepsilon C_{\delta} \int_{\mathbb{R}^n} u^{2\gamma+p-1} \eta^{2\gamma m},$$

where $C_{\delta,\varepsilon}$ is a positive constant depends on δ and ε , and θ is a number such that $2(1-\theta) + (2\gamma + p - 1)\theta = 2\gamma$, so that $0 < \theta < 1$ for $2 < 2\gamma < 2\gamma + p - 1$. By (2-30) and (2-32), we get

$$(2-33) \qquad (E-p\delta-\varepsilon C_{\delta})\int_{\mathbb{R}^{n}}u^{p+2\gamma-1}\eta^{2\gamma m} \leq C_{\delta,\varepsilon}R^{-\frac{4}{1-\theta}}\int_{\mathbb{R}^{n}}u^{2}\eta^{2\gamma m-\frac{4}{1-\theta}}$$

Since θ is strictly less than 1 and will be fixed for given γ , p, we can choose m sufficiently large to make $2\gamma m - \frac{4}{1-\theta} > 0$. On the other hand, if E > 0, we can find small δ and then small ε , such that $E - p\delta - \varepsilon C_{\delta} > 0$. Therefore, by the definition of function η and (2-33), we obtain

(2-34)
$$(E - p\delta - \varepsilon C_{\delta}) \int_{B_R} u^{p+2\gamma-1} \le C_{\delta,\varepsilon} R^{-\frac{4}{1-\theta}} \int_{B_{2R}} u^2.$$

By (2.10) of [Wei and Ye 2010], we have $\int_{B_{2R}} u^2 \leq C R^{n-\frac{8}{p-1}}$, as a result, the

left side of (2-34) is less equal than $C_{\delta,\varepsilon}R^{n-\frac{8}{p-1}-\frac{4}{1-\theta}}$, which tends to 0 as *R* tends to ∞ , provided the power $n - \frac{8}{p-1} - \frac{4}{1-\theta}$ is negative, which is equivalent to $(p+2\gamma-1) > (p-1)\frac{n}{4}$ according to the definition of θ . So, if $(p+2\gamma-1) > (p-1)\frac{n}{4}$ and $E - p\delta - C_{\delta}\varepsilon > 0$, we have $u \equiv 0$.

Thus, we have proved the nonexistence of stable solution to (2-8) if p satisfies the condition $(p + 2\gamma - 1) > (p - 1)\frac{n}{4}$ and E > 0 (for δ , ε are arbitrary small). By Lemma A.2 in the Appendix, the power p can be in the interval $(\frac{n}{n-8}, 1 + \frac{8p^*}{n-4})$. Combining with Theorem 1.1 of [Wei and Ye 2010], we have proved Theorem 1.1, that is, for any $1 , <math>n \ge 20$, (2-8) has no stable solution.

3. Proof of Theorem 1.2

In proving Theorem 1.2, it is enough to consider stable solutions u_{λ} to (1-4), since $u^* = \lim_{\lambda \to \lambda^*} u_{\lambda}$. Now we give a uniform bound for the stable solutions to (1-4) when $0 < d < \lambda < \lambda^*$, where *d* is a fixed positive constant from $(0, \lambda^*)$.

First, we need to analyze the solution near the boundary. Specifically, we need the regularity of the stable solutions of the equation

(3-1)
$$\begin{cases} \Delta^2 u = \lambda (u+1)^p, \ \lambda > 0 \text{ in } \Omega, \\ u > 0 \text{ in } \Omega, \\ u = \Delta u = 0 \text{ on } \partial \Omega \end{cases}$$

near the boundary (as well as their derivatives; see remark after the next theorem).

Theorem 3.1. Let Ω be a bounded, smooth, and convex domain. There exists a constant *C* (independent of λ , *u*) and small positive number ε , such that for stable solutions *u* to (3-1) we have

$$(3-2) u(x) < C for all x \in \Omega_{\varepsilon} := \{z \in \Omega : d(z, \partial \Omega) < \varepsilon\}.$$

Proof. This result is well known. See [Guo and Wei 2009]. For the sake of completeness, we include a proof here. By Lemma 3.5 of [Cowan et al. 2010], we see that there exists a constant *C* independent of λ , *u*, such that

(3-3)
$$\int_{\Omega} (1+u)^p \, dx \le C$$

We write (3-1) as

$$\begin{cases} \Delta u + v = 0, & \text{in } \Omega, \\ \Delta v + \lambda (1+u)^p = 0, & \text{in } \Omega, \\ u = v = 0, & \text{in } \partial \Omega \end{cases}$$

If we set $f_1(u, v) = v$, $f_2(u, v) = \lambda(u+1)^p$, we see that $\partial f_1/\partial v = 1 > 0$ and $\partial f_2/\partial u = \lambda p(u+1)^{p-1} > 0$. Therefore, the convexity of Ω , Lemma 5.1 of [Troy 1981], and the moving plane method near $\partial \Omega$ (as in the appendix of [Guo and
Webb 2002]) imply that there exist $t_0 > 0$ and α which depends only on the domain Ω , such that $u(x - t\nu)$ and $v(x - t\nu)$ are nondecreasing for $t \in [0, t_0]$, $\nu \in \mathbb{R}^n$ satisfying $|\nu| = 1$ and $(\nu, n(x)) \ge \alpha$ and $x \in \partial \Omega$. Therefore, we can find $\rho, \varepsilon > 0$ such that for any $x \in \Omega_{\varepsilon} := \{z \in \Omega : d(z, \partial \Omega) < \varepsilon\}$ there exists a fixed-sized cone Γ_x (with x as its vertex) with

- meas $(\Gamma_x) \ge \rho$,
- $\Gamma_x \subset \{z \in \Omega : d(z, \partial \Omega) < 2\varepsilon\}$, and
- $u(y) \ge u(x)$ for any $y \in \Gamma_x$.

Then, for any $x \in \Omega_{\varepsilon}$, we have

$$(1+u(x))^p \le \frac{1}{\operatorname{meas}(\Gamma_x)} \int_{\Gamma_x} (1+u)^p \le \frac{1}{\rho} \int_{\Omega} (1+u)^p \le C.$$

This implies that $(1 + u(x))^p \le C$, therefore $u(x) \le C$.

Remark. By classical elliptic regularity theory, u(x) and its derivatives up to fourth order are bounded on the boundary by a constant independent of u. See [Wei 1996] for more details.

We now turn to the proof of Theorem 1.2 proper, using the ideas of Section 2. Multiplying (1-4) by $(u + 1)^q$ and integrating by parts, we have

(3-4)
$$\int_{\Omega} \lambda(u+1)^{p+q} = \int_{\Omega} \Delta^2 u(u+1)^q = \int_{\partial\Omega} \frac{\partial(\Delta u)}{\partial n} + \int_{\Omega} \Delta(u+1)\Delta(u+1)^q.$$

Setting v = u + 1, by direct calculation, we get

$$\begin{split} \int_{\Omega} (\Delta v^{\gamma})^2 &= \int_{\Omega} \gamma^2 v^{2\gamma-2} (\Delta v)^2 + \int_{\Omega} \gamma^2 (\gamma-1)^2 v^{2\gamma-4} |\nabla v|^4 \\ &+ 2 \int_{\Omega} \gamma^2 (\gamma-1) v^{2\gamma-3} \Delta v |\nabla v|^2, \\ \int_{\Omega} \Delta v \Delta v^q &= \int_{\Omega} q (\Delta v)^2 v^{q-1} + \int_{\Omega} q (q-1) |\nabla v|^2 \Delta v v^{q-2}. \end{split}$$

From these two equalities and (3-4) we obtain

(3-5)
$$\int_{\Omega} \left(\frac{q}{\gamma^2} (\Delta v^{\gamma})^2 - q(\gamma - 1)^2 |\nabla v|^4 v^{2\gamma - 4} \right) + \int_{\partial \Omega} \frac{\partial (\Delta v)}{\partial n} = \int_{\Omega} \lambda v^{p+q}.$$

For the second term in (3-5), we have

$$(3-6) \quad \int_{\Omega} |\nabla v|^{4} v^{2\gamma-4}$$

$$= \frac{1}{\gamma^{4}} \int_{\Omega} v^{-2\gamma} |\nabla v^{\gamma}|^{4} = \frac{1}{\gamma^{4}} \int_{\Omega} |\nabla v^{\gamma}|^{2} \nabla v^{\gamma} (-\nabla v^{-\gamma})$$

$$= \frac{1}{\gamma^{4}} \int_{\Omega} \left(-\nabla \frac{|\nabla v^{\gamma}|^{2} \nabla v^{\gamma}}{v^{\gamma}} + \frac{\nabla (|\nabla v^{\gamma}|^{2}) \nabla v^{\gamma}}{v^{\gamma}} + \frac{|\nabla v^{\gamma}|^{2} \Delta v^{\gamma}}{v^{\gamma}} \right)$$

$$= \frac{1}{\gamma^{4}} \int_{\Omega} v^{-\gamma} \nabla (|\nabla v^{\gamma}|^{2}) \nabla v^{\gamma} + |\nabla v^{\gamma}|^{2} \Delta v^{\gamma} - \frac{1}{\gamma} \int_{\partial \Omega} v^{2\gamma-3} |\nabla v|^{2} \frac{\partial v}{\partial n}.$$

A simple calculation yields

(3-7)
$$\frac{1}{\gamma^4} \int_{\Omega} v^{-\gamma} |\nabla v^{\gamma}|^2 \Delta v^{\gamma} = \frac{\gamma - 1}{\gamma} \int_{\Omega} v^{2\gamma - 4} |\nabla v|^4 + \frac{1}{\gamma} \int_{\Omega} v^{2\gamma - 3} |\nabla v|^2 \Delta v.$$

Substituting (3-7) into (3-6), we get

(3-8)
$$\int_{\Omega} |\nabla v|^4 v^{2\gamma-4} = \int_{\Omega} v^{2\gamma-3} |\nabla v|^2 \Delta v + \frac{1}{\gamma^3} \int_{\Omega} v^{-\gamma} \nabla (|\nabla v^{\gamma}|^2) \nabla v^{\gamma} - \int_{\partial \Omega} |\nabla v|^2 \frac{\partial v}{\partial n}.$$

We now estimate the second term on the right side of (3-8). From the proof of Lemma 2.3, together with the identity $\frac{1}{2}\Delta |\nabla v^{\gamma}|^2 = |\nabla^2 v^{\gamma}|^2 + \langle \nabla \Delta v^{\gamma}, \nabla v^{\gamma} \rangle$, we have

$$(3-9) \quad \frac{1}{\gamma^3} \int_{\Omega} v^{-\gamma} \nabla (|\nabla v^{\gamma}|^2) \nabla v^{\gamma} \leq \frac{1}{2} \int_{\Omega} |\nabla v|^4 v^{2\gamma-4} + \frac{2}{\gamma^2} \int_{\Omega} (\Delta v^{\gamma})^2 \\ + \frac{1}{\gamma^2} \int_{\partial\Omega} \frac{\partial |\nabla v^{\gamma}|^2}{\partial n} - \frac{2}{\gamma^2} \int_{\partial\Omega} (\Delta v^{\gamma}) \frac{\partial v^{\gamma}}{\partial n}.$$

By (3-8) and (3-9), thanks to the convexity of the domain Ω , we get

$$(3-10) \quad \frac{1}{2} \int_{\Omega} |\nabla v|^4 v^{2\gamma - 4} \\ \leq \int_{\Omega} v^{2\gamma - 3} |\nabla v|^2 \Delta v + \frac{2}{\gamma^2} \int_{\Omega} (\Delta v^{\gamma})^2 - (2\gamma - 1) \int_{\partial \Omega} |\nabla v|^2 \frac{\partial v}{\partial n}.$$

For the first term on the right side of (3-10), since v = u + 1, we have $\Delta v = \Delta u < 0$ by maximal principle, and the inequality

$$(3-11) \qquad \qquad \Delta v < -\sqrt{\frac{2\lambda}{p+1}}v^{\frac{p+1}{2}} < 0,$$

by Lemma 3.2 of [Cowan et al. 2010]. Thus

$$\int_{\Omega} v^{2\gamma-3} |\nabla v|^2 \Delta v \le \int_{\Omega} -\sqrt{\frac{2\lambda}{p+1}} v^{2\gamma-3+\frac{p+1}{2}} |\nabla v|^2.$$

508

Moreover, we have

$$\begin{split} \int_{\Omega} -\sqrt{\frac{2\lambda}{p+1}} v^{2\gamma-3+\frac{p+1}{2}} |\nabla v|^2 &= -\int_{\Omega} \frac{\sqrt{\frac{2\lambda}{p+1}}}{2\gamma-2+\frac{p+1}{2}} \nabla (v^{2\gamma-2+\frac{p+1}{2}} \nabla v) \\ &+ \int_{\Omega} \frac{\sqrt{\frac{2\lambda}{p+1}}}{2\gamma-2+\frac{p+1}{2}} v^{2\gamma-2+\frac{p+1}{2}} \Delta v. \end{split}$$

For the second term on the right, using (3-11) again, we have

$$\int_{\Omega} \frac{\sqrt{\frac{2\lambda}{p+1}}}{2\gamma - 2 + \frac{p+1}{2}} v^{2\gamma - 2 + \frac{p+1}{2}} \Delta v \le -\int_{\Omega} \frac{\frac{2\lambda}{p+1}}{2\gamma - 2 + \frac{p+1}{2}} v^{2\gamma + p - 1}$$

Hence, we obtain

$$(3-12) \quad \int_{\Omega} v^{2\gamma-3} |\nabla v|^2 \Delta v \le -\int_{\partial\Omega} \frac{\sqrt{\frac{2\lambda}{p+1}}}{2\gamma-2+\frac{p+1}{2}} \frac{\partial v}{\partial n} - \int_{\Omega} \frac{\frac{2\lambda}{p+1}}{2\gamma-2+\frac{p+1}{2}} v^{2\gamma+p-1},$$

where we used $v|_{\partial\Omega} = u + 1|_{\partial\Omega} = 1$, for the boundary term in (3-4), (3-10) and (3-12). By the remark after Theorem 3.1, we find that there exists a constant *C* (the constant *C* appeared now and later in this section is independent of *u*), such that

(3-13)
$$\int_{\partial\Omega} \left(|\nabla u|^2 \left| \frac{\partial u}{\partial n} \right| + \left| \frac{\partial (\Delta u)}{\partial n} \right| + \left| \frac{\partial u}{\partial n} \right| \right) \le C.$$

Combining (3-5), (3-10), (3-12) and (3-13), we get

$$\left(1 - 4(\gamma - 1)^2\right) \int_{\Omega} (\Delta(u+1)^{\gamma})^2 + \left(\frac{8\lambda\gamma^2(\gamma - 1)^2}{(4\gamma + p - 3)(p+1)} - \frac{\lambda\gamma^2}{q}\right) \int_{\Omega} (u+1)^{p+q} \le C.$$

If $1 - 4(\gamma - 1)^2 > 0$ and

(3-14)
$$p(1-4(\gamma-1)^2) + \frac{8\gamma^2(\gamma-1)^2}{(4\gamma+p-3)(p+1)} - \frac{\gamma^2}{q} > 0$$

and u is a stable solution to (1-4), we have

$$\left(p(1-4(\gamma-1)^2) + \frac{8\gamma^2(\gamma-1)^2}{(4\gamma+p-3)(p+1)} - \frac{\gamma^2}{2\gamma-1}\right) \int_{\Omega} (u+1)^{p+q} \le \frac{C}{\lambda}.$$

This leads to $u + 1 \in L^{p+q}$.

If p + q > (p - 1)n/4, then classical regularity theory implies that $u \in L^{\infty}(\Omega)$.

Therefore we have established the bound of extremal solutions of (1-4) if (3-14) is satisfied and

$$p < \frac{8\gamma + n - 4}{n - 4}.$$

By Lemma A.2 and Theorem 3.8 of [Wei and Ye 2010], we have proved that the extremal solution u^* , the unique solution of (1-4) (where $\lambda = \lambda^*$), is bounded provided that one of these conditions hold:

- (1) If $n \le 8$, then p > 1.
- (2) If $9 \le n \le 19$, there exists $\varepsilon_n > 0$ such that for any 1 .
- (3) If $n \ge 20$, then $1 , where <math>p^*$ was defined immediately after Theorem 1.1.

Appendix

In this appendix, we study the inequalities

(A-1)
$$p(1-4(\gamma-1)^2) - \frac{\gamma^2}{2\gamma-1} + \frac{8\gamma^2(\gamma-1)^2}{(4\gamma-3+p)(p+1)} > 0$$

and

$$(A-2) p < \frac{8\gamma + n - 4}{n - 4}.$$

In order to get a better range for the power p from (A-1) and (A-2), we must study the following equation obtained by letting $p = \frac{8\gamma + n - 4}{n - 4}$ in (A-1):

(A-3)
$$\frac{8\gamma + n - 4}{n - 4} \left(1 - 4(\gamma - 1)^2 \right) - \frac{\gamma^2}{2\gamma - 1} + \frac{8\gamma^2 (\gamma - 1)^2}{\left(4\gamma - 3 + \frac{8\gamma + n - 4}{n - 4}\right)\left(\frac{8\gamma + n - 4}{n - 4} + 1\right)} = 0.$$

We need only consider the behavior of (A-3) for $\gamma \in (1, \frac{3}{2})$. Through tedious computations, we see that the equation at the bottom of page 495 is the simplified form of (A-3). As a consequence, they have same roots in $(1, \frac{3}{2})$.

We denote the left side of (A-3) by $h(\gamma)$. Notice that if $\gamma = \frac{n-4}{n-8}$, then $p = \frac{n}{n-8}$ and $\gamma - 1 = \frac{4}{n-8}$. Hence

$$h\left(\frac{n-4}{n-8}\right) = \frac{8}{n-8}(n^4 - 18n^3 - 56n^2 + 384n - 512).$$

In fact, if n = 20, then $h\left(\frac{4}{3}\right) = 512 > 0$. On the other hand, it is also easy to see that $h\left(\frac{3}{2}\right) < 0$, while it is obvious that $\left(4\gamma - 3 + \frac{8\gamma + n - 4}{n - 4}\right)\left(\frac{8\gamma + n - 4}{n - 4} + 1\right) > 0$ and $(2\gamma - 1) > 0$ when $\gamma \in \left(\frac{n - 4}{n - 8}, \frac{3}{2}\right)$. Therefore, by continuity, (A-3) possesses a root in $\left(\frac{n - 4}{n - 8}, \frac{3}{2}\right)$. We denote the smallest root of (A-3) greater than $\frac{n - 4}{n - 8}$ by p^* . Once we pick out a γ from the interval $\left(\frac{n - 4}{n - 8}, p^*\right)$, $h(\gamma)$ is of course positive. By continuity, we can find a small positive number δ such that the inequality

$$p(1-4(\gamma-1)^2) - \frac{\gamma^2}{2\gamma-1} + \frac{8\gamma^2(\gamma-1)^2}{(4\gamma-3+p)(p+1)} > 0$$

holds when $p \in \left(\frac{8\gamma+n-4}{n-4} - \delta, \frac{8\gamma+n-4}{n-4}\right)$. So, we conclude that when γ runs in the whole interval $\left(\frac{n-4}{n-8}, p^*\right)$, the power p can be in the whole interval $\left(\frac{n}{n-8}, 1 + \frac{8p^*}{n-4}\right)$. We summarize the result as follows:

Lemma A.2. When $n \ge 20$, the range of p satisfying (A-1) and (A-2) equals $\left(\frac{n}{n-8}, 1+\frac{8p^*}{n-4}\right)$, and this interval is not empty.

References

- [Berchio and Gazzola 2005] E. Berchio and F. Gazzola, "Some remarks on biharmonic elliptic problems with positive, increasing and convex nonlinearities", *Electron. J. Differential Equations* **2005** (2005), Paper No. 34, 1–20. MR 2006e:35067 Zbl 1129.35349
- [Cowan et al. 2010] C. Cowan, P. Esposito, and N. Ghoussoub, "Regularity of extremal solutions in fourth order nonlinear eigenvalue problems on general domains", *Discrete Contin. Dyn. Syst.* 28:3 (2010), 1033–1050. MR 2011e:35085 Zbl 1196.35152
- [Farina 2007] A. Farina, "On the classification of solutions of the Lane–Emden equation on unbounded domains of \mathbb{R}^N ", J. Math. Pures Appl. (9) 87:5 (2007), 537–561. MR 2008c:35070 Zbl 1143.35041
- [Ferrero et al. 2009] A. Ferrero, H.-C. Grunau, and P. Karageorgis, "Supercritical biharmonic equations with power-type nonlinearity", *Ann. Mat. Pura Appl.* (4) **188**:1 (2009), 171–185. MR 2010g: 35087 Zbl 1179.35125
- [Gazzola and Grunau 2006] F. Gazzola and H.-C. Grunau, "Radial entire solutions for supercritical biharmonic equations", *Math. Ann.* **334**:4 (2006), 905–936. MR 2007b:35114 Zbl 1152.35034
- [Gui et al. 1992] C. Gui, W.-M. Ni, and X. Wang, "On the stability and instability of positive steady states of a semilinear heat equation in \mathbb{R}^n ", *Comm. Pure Appl. Math.* **45**:9 (1992), 1153–1181. MR 93h:35095 Zbl 0811.35048
- [Guo and Webb 2002] Z. Guo and J. R. L. Webb, "Large and small solutions of a class of quasilinear elliptic eigenvalue problems", *J. Differential Equations* **180**:1 (2002), 1–50. MR 2002k:35246 Zbl 1014.35030
- [Guo and Wei 2009] Z. Guo and J. Wei, "On a fourth order nonlinear elliptic equation with negative exponent", *SIAM J. Math. Anal.* **40**:5 (2009), 2034–2054. MR 2010b:35125 Zbl 1175.35144
- [Guo and Wei 2010] Z. Guo and J. Wei, "Qualitative properties of entire radial solutions for a biharmonic equation with supercritical nonlinearity", *Proc. Amer. Math. Soc.* **138**:11 (2010), 3957–3964. MR 2012a:35054 Zbl 1203.35105
- [Lin 1998] C.-S. Lin, "A classification of solutions of a conformally invariant fourth order equation in \mathbb{R}^n ", *Comment. Math. Helv.* **73**:2 (1998), 206–231. MR 99c:35062 Zbl 0933.35057
- [Troy 1981] W. C. Troy, "Symmetry properties in systems of semilinear elliptic equations", *J. Differential Equations* **42**:3 (1981), 400–413. MR 83b:35051 Zbl 0486.35032
- [Wei 1996] J. Wei, "Asymptotic behavior of a nonlinear fourth order eigenvalue problem", *Comm. Partial Differential Equations* **21**:9-10 (1996), 1451–1467. MR 97h:35066 Zbl 0872.35013
- [Wei and Xu 1999] J. Wei and X. Xu, "Classification of solutions of higher order conformally invariant equations", *Math. Ann.* **313**:2 (1999), 207–228. MR 2000a:58093 Zbl 0940.35082
- [Wei and Ye 2010] J. Wei and D. Ye, "Liouville theorems for finite Morse index solutions of biharmonic problem", preprint, 2010, http://www.math.cuhk.edu.hk/~wei/4thstable29-10-10.pdf. To appear in *Math. Ann.*

[Xu 2000] X. Xu, "Uniqueness theorem for the entire positive solutions of biharmonic equations in \mathbb{R}^{n} ", *Proc. Roy. Soc. Edinburgh Sect. A* **130**:3 (2000), 651–670. MR 2001f:35143 Zbl 0961.35037

Received November 1, 2011. Revised April 24, 2012.

JUNCHENG WEI DEPARTMENT OF MATHEMATICS THE CHINESE UNIVERSITY OF HONG KONG SHATIN, NT HONG KONG wei@math.cuhk.edu.hk

XINGWANG XU DEPARTMENT OF MATHEMATICS NATIONAL UNIVERSITY OF SINGAPORE BLOCK S17 (SOC1) 10 LOWER KENT RIDGE ROAD SINGAPORE 119076 SINGAPORE

matxuxw@nus.edu.sg

WEN YANG DEPARTMENT OF MATHEMTICS CHINESE UNIVERSITY OF HONG KONG SHATIN, NT HONG KONG wyang@math.cuhk.edu.hk

512

CONTENTS

Volume 263, no. 1 and no. 2

Luis J. Alías, S. Carolina García-Martínez and Marco Rigoli: <i>Biharmonic</i> hypersurfaces in complete Riemannian manifolds	1
G. A. Zimmermann with G. A. T. F. da Costa	475
Sandeep Bhargava and Yun Gao: <i>Realizations of BC_r-graded intersection matrix</i> algebras with grading subalgebras of type B_r , $r \ge 3$	257
Julien Bichon and Michel Dubois-Violette: <i>Half-commutative orthogonal Hopf</i> algebras	13
Claudio Carmeli and Rita Fioresi: <i>Superdistributions, analytic and algebraic super</i> <i>Harish-Chandra pairs</i>	29
Angel Carocca , Rubén A. Hidalgo and Rubí E. Rodríguez: <i>Orbifolds with signature</i> $(0; k, k^{n-1}, k^n, k^n)$	53
Ercai Chen with Xianfeng Ma	117
Imin Chen and Yoonjin Lee: Explicit isogeny theorems for Drinfeld modules	87
Elie Compoint and Eduardo Corel: <i>Stable flags, trivializations and regular connections</i>	283
Eduardo Corel with Elie Compoint	283
Michel Dubois-Violette with Julien Bichon	13
Rita Fioresi with Claudio Carmeli	29
Yun Gao with Sandeep Bhargava	257
S. Carolina García-Martínez with Luis J. Alías and Marco Rigoli	1
Rubén A. Hidalgo with Angel Carocca and Rubí E. Rodríguez	53
Nathan Jones: Elliptic aliquot cycles of fixed length	353
Daniel J. Katz : Asymptotic L^4 norm of polynomials derived from characters	373
Junho Lee: Degree-three spin Hurwitz numbers	399
Yoonjin Lee with Imin Chen	87
Hong-Quan Li: Remark on "Maximal functions on the unit n-sphere" by Peter M. Knopf (1987)	253
Youlin Li and Jiming Ma: $(\mathbb{Z}_2)^3$ -colorings and right-angled hyperbolic 3-manifolds	419

514

419
117
143
171
435
453
1
53
207
475
435
453
495
241
495
453
495

Guidelines for Authors

Authors may submit manuscripts at msp.berkeley.edu/pjm/about/journal/submissions.html and choose an editor at that time. Exceptionally, a paper may be submitted in hard copy to one of the editors; authors should keep a copy.

By submitting a manuscript you assert that it is original and is not under consideration for publication elsewhere. Instructions on manuscript preparation are provided below. For further information, visit the web address above or write to pacific@math.berkeley.edu or to Pacific Journal of Mathematics, University of California, Los Angeles, CA 90095–1555. Correspondence by email is requested for convenience and speed.

Manuscripts must be in English, French or German. A brief abstract of about 150 words or less in English must be included. The abstract should be self-contained and not make any reference to the bibliography. Also required are keywords and subject classification for the article, and, for each author, postal address, affiliation (if appropriate) and email address if available. A home-page URL is optional.

Authors are encouraged to use LAT_EX , but papers in other varieties of T_EX , and exceptionally in other formats, are acceptable. At submission time only a PDF file is required; follow the instructions at the web address above. Carefully preserve all relevant files, such as LAT_EX sources and individual files for each figure; you will be asked to submit them upon acceptance of the paper.

Bibliographical references should be listed alphabetically at the end of the paper. All references in the bibliography should be cited in the text. Use of BibT_EX is preferred but not required. Any bibliographical citation style may be used but tags will be converted to the house format (see a current issue for examples).

Figures, whether prepared electronically or hand-drawn, must be of publication quality. Figures prepared electronically should be submitted in Encapsulated PostScript (EPS) or in a form that can be converted to EPS, such as GnuPlot, Maple or Mathematica. Many drawing tools such as Adobe Illustrator and Aldus FreeHand can produce EPS output. Figures containing bitmaps should be generated at the highest possible resolution. If there is doubt whether a particular figure is in an acceptable format, the authors should check with production by sending an email to pacific@math.berkeley.edu.

Each figure should be captioned and numbered, so that it can float. Small figures occupying no more than three lines of vertical space can be kept in the text ("the curve looks like this:"). It is acceptable to submit a manuscript will all figures at the end, if their placement is specified in the text by means of comments such as "Place Figure 1 here". The same considerations apply to tables, which should be used sparingly.

Forced line breaks or page breaks should not be inserted in the document. There is no point in your trying to optimize line and page breaks in the original manuscript. The manuscript will be reformatted to use the journal's preferred fonts and layout.

Page proofs will be made available to authors (or to the designated corresponding author) at a website in PDF format. Failure to acknowledge the receipt of proofs or to return corrections within the requested deadline may cause publication to be postponed.

PACIFIC JOURNAL OF MATHEMATICS

Volume 263 No. 2 June 2013

Realizations of BC_r -graded intersection matrix algebras with grading subalgebras of type B_r , $r \ge 3$	257
SANDEEP BHARGAVA and YUN GAO	
Stable flags, trivializations and regular connections	283
ELIE COMPOINT and EDUARDO COREL	
Elliptic aliquot cycles of fixed length	353
NATHAN JONES	
Asymptotic L^4 norm of polynomials derived from characters	373
DANIEL J. KATZ	
Degree-three spin Hurwitz numbers	399
JUNHO LEE	
$(\mathbb{Z}_2)^3$ -colorings and right-angled hyperbolic 3-manifolds	419
YOULIN LI and JIMING MA	
Real closed separation theorems and applications to group algebras	435
TIM NETZER and ANDREAS THOM	
Uniqueness theorem for ordinary differential equations with Hölder	453
continuity	
YIFEI PAN, MEI WANG and YU YAN	
An analogue to the Witt identity	475
G. A. T. F. DA COSTA and G. A. ZIMMERMANN	
On the classification of stable solutions to biharmonic problems in	495
large dimensions	
JUNCHENG WEL XINGWANG XU and WEN YANG	

