



Geometry & Topology

Volume 29 (2025)

**Random unitary representations of surface groups
II: The large n limit**

MICHAEL MAGEE

Random unitary representations of surface groups

II: The large n limit

MICHAEL MAGEE

Let Σ_g be a closed surface of genus $g \geq 2$ and Γ_g denote the fundamental group of Σ_g . We establish a generalization of Voiculescu’s theorem on the asymptotic $*$ -freeness of Haar unitary matrices from free groups to Γ_g . We prove that, for a random representation of Γ_g into $SU(n)$, with law given by the volume form arising from the Atiyah–Bott–Goldman symplectic form on moduli space, the expected value of the trace of a fixed nonidentity element of Γ_g is bounded as $n \rightarrow \infty$. The proof involves an interplay between Dehn’s work on the word problem in Γ_g and classical invariant theory.

14H60, 22D10, 46L54; 20C30, 20C35, 32G15, 70S15

To Oona

1. Introduction	1237
2. Background	1244
3. Combinatorial integration	1251
4. Topology	1264
5. Proof of the main theorem	1279
References	1279

1 Introduction

In a foundational series of papers, Voiculescu [1985; 1986; 1987; 1990; 1991] developed a robust theory of noncommuting random variables that became known as *free probability*. One of the initial landmarks of this theory is the following result. Let F_r denote the noncommutative free group of rank r . Let $U(n)$ denote the group of $n \times n$ complex unitary matrices. For any $w \in F_r$ we obtain a *word map* $w: U(n)^r \rightarrow U(n)$ by substituting matrices for generators of F_r . Let $\mu_{U(n)^r}^{\text{Haar}}$ denote the probability Haar measure on $U(n)^r$ and $\text{Tr}: U(n) \rightarrow \mathbb{C}$ the standard trace. Any integral over a compact group will be done with respect to the probability Haar measure, denoted by $d\mu$.

A simplified version of Voiculescu's result [1991, Theorem 3.8] can be formulated as follows:¹

Theorem 1.1 (Voiculescu) *For any nonidentity $w \in \mathbf{F}_r$, as $n \rightarrow \infty$,*

$$(1-1) \quad \int_{\mathbf{U}(n)^r} \text{Tr}(w(x)) d\mu(x) = o_w(n).$$

We describe the interpretation of Theorem 1.1 as convergence of noncommutative random variables in a moment. Before this, we explain the main result of the current paper.

Another way to think about the integral (1-1), which invites generalization, is to identify $\mathbf{U}(n)^r$ with $\text{Hom}(\mathbf{F}_r, \mathbf{U}(n))$ and Haar measure as a natural probability measure on this *representation variety*. Now it is natural to ask whether there are other infinite discrete groups G besides \mathbf{F}_r such that $\text{Hom}(G, \mathbf{U}(n))$ has a natural measure, and whether similar phenomena as in Theorem 1.1 may hold. *The main point of this paper is to establish the analog of Theorem 1.1 when \mathbf{F}_r is replaced by the fundamental group of a compact surface of genus at least 2.*

We now explain this generalization of Theorem 1.1; for technical reasons it superficially looks slightly different, as follows:

- (1) The integral (1-1) is equal to 0 if $w \notin [\mathbf{F}_r, \mathbf{F}_r]$, the commutator subgroup of \mathbf{F}_r [Magee and Puder 2015, Claim 3.1], and, if $w \in [\mathbf{F}_r, \mathbf{F}_r]$, the value of (1-1) is, for $n \geq n_0(w)$, the same as the corresponding integral over $\text{SU}(n)^r \leq \mathbf{U}(n)^r$, where $\text{SU}(n)$ is the subgroup of determinant one matrices [Magee 2022, Proposition 3.1]. So in all cases of interest we can replace $\mathbf{U}(n)$ by $\text{SU}(n)$ in (1-1).
- (2) Since $\text{Tr} \circ w$ is invariant under the diagonal conjugation action of $\text{SU}(n)$ on $\text{Hom}(\mathbf{F}_r, \text{SU}(n)) \cong \text{SU}(n)^r$, the integral $\int_{\text{SU}(n)^r} \text{Tr}(w(x)) d\mu(x)$ can be written as one over $\text{Hom}(\mathbf{F}_r, \text{SU}(n))/\text{PSU}(n)$. Here $\text{PSU}(n)$ is $\text{SU}(n)$ modulo its center.

For $g \geq 2$ let Σ_g denote a closed topological surface of genus g . We let Γ_g denote the fundamental group of Σ_g with explicit presentation

$$\Gamma_g = \langle a_1, b_1, \dots, a_g, b_g \mid [a_1, b_1] \cdots [a_g, b_g] \rangle.$$

The most natural measure on $\text{Hom}(\Gamma_g, \text{SU}(n))/\text{PSU}(n)$ to replace the measure induced by Haar measure on $\text{Hom}(\mathbf{F}_r, \text{SU}(n))/\text{PSU}(n)$ is called the Atiyah–Bott–Goldman measure. The definition of this measure involves removing singular parts of $\text{Hom}(\Gamma_g, \text{SU}(n))/\text{PSU}(n)$. Indeed, let $\text{Hom}(\Gamma_g, \text{SU}(n))^{\text{irr}}$ denote the collection of homomorphisms that are irreducible as linear representations. Then

$$\mathcal{M}_{g,n} := \text{Hom}(\Gamma_g, \text{SU}(n))^{\text{irr}}/\text{PSU}(n)$$

¹Voiculescu's result [1991, Theorem 3.8] is more general than what we state here, also involving a deterministic sequence of unitary matrices.

is a smooth manifold [Goldman 1984]. Moreover there is a symplectic form $\omega_{g,n}$ on $\mathcal{M}_{g,n}$, called the Atiyah–Bott–Goldman form after [Atiyah and Bott 1983; Goldman 1984]. This symplectic form gives, in the usual way, a volume form on $\mathcal{M}_{g,n}$ denoted by $\text{Vol}_{\mathcal{M}_{g,n}}$. For many more details, see [Goldman 1984] or our prequel paper [Magee 2022, Section 2.7].

For any $\gamma \in \Gamma$, we obtain a function $\text{Tr}_\gamma : \text{Hom}(\Gamma_g, \text{SU}(n)) \rightarrow \mathbb{C}$ defined by

$$\text{Tr}_\gamma(\phi) := \text{Tr}(\phi(\gamma)).$$

This function descends to a function $\text{Tr}_\gamma : \mathcal{M}_{g,n} \rightarrow \mathbb{C}$. We are interested in the expected value

$$\mathbb{E}_{g,n}[\text{Tr}_\gamma] := \frac{\int_{\mathcal{M}_{g,n}} \text{Tr}_\gamma d\text{Vol}_{\mathcal{M}_{g,n}}}{\int_{\mathcal{M}_{g,n}} d\text{Vol}_{\mathcal{M}_{g,n}}}.$$

The main theorem of this paper is the following:

Theorem 1.2 *Let $g \geq 2$. If $\gamma \in \Gamma_g$ is not the identity, then $\mathbb{E}_{g,n}[\text{Tr}_\gamma] = O_\gamma(1)$ as $n \rightarrow \infty$.*

The noncommutative probabilistic consequences of Theorem 1.2 will be discussed in the next section.

1.1 Noncommutative probability

We follow [Voiculescu et al. 1992]. A *noncommutative probability space* is a pair (\mathcal{B}, τ) where \mathcal{B} is a complex unital algebra and τ is a linear functional on \mathcal{B} such that $\tau(1) = 1$. Let $\mathbb{C}\langle x_1, \dots, x_r \rangle$ denote the free noncommutative unital algebra in indeterminates x_1, \dots, x_r . A *random variable* in (\mathcal{B}, τ) is an element of \mathcal{B} . If $(X_1, \dots, X_r) \in \mathcal{B}^r$ are random variables in (\mathcal{B}, τ) , their *joint distribution* is defined to be the linear functional

$$\tilde{\tau} : \mathbb{C}\langle x_1, \dots, x_r \rangle \rightarrow \mathbb{C}$$

given by $\tilde{\tau}(z) := \tau(\Phi(z))$, where $\Phi : \mathbb{C}\langle x_1, \dots, x_r \rangle \rightarrow \mathcal{B}$ is the linear map defined by $\Phi(x_i) = X_i$. For a linear functional $\tilde{\tau}_\infty : \mathbb{C}\langle x_1, \dots, x_r \rangle \rightarrow \mathbb{C}$ with $\tilde{\tau}_\infty(1) = 1$, we say that a sequence of random variables $(X_1^{(n)}, \dots, X_r^{(n)}) \in (\mathcal{B}_n, \tau_n)$ converges in distribution as $n \rightarrow \infty$ to $\tilde{\tau}_\infty$ if $\tilde{\tau}_n$ converges pointwise to $\tilde{\tau}_\infty$ on $\mathbb{C}\langle x_1, \dots, x_r \rangle$.

A very concrete example of this phenomenon is as follows. The function

$$\tau_n : \mathbf{F}_r \rightarrow \mathbb{C}, \quad \tau_n(w) := \frac{1}{n} \int_{\text{U}(n)^r} \text{Tr}(w(x)) d\mu(x)$$

extends to a linear functional τ_n on the algebra $\mathbb{C}[\mathbf{F}_r]$ with $\tau_n(\text{id}) = 1$. From this point of view, Theorem 1.1 implies the following statement:

Theorem 1.3 (Voiculescu) *Let $r \geq 0$ and X_1, \dots, X_r denote fixed generators of F_r , and $\bar{X}_1, \dots, \bar{X}_r$ denote their inverses, ie $\bar{X}_i = X_i^{-1}$. The random variables $X_1, \dots, X_r, \bar{X}_1, \dots, \bar{X}_r$ in the noncommutative probability spaces $(\mathbb{C}[F_r], \tau_n)$ converge as $n \rightarrow \infty$ to a limiting distribution*

$$\tilde{\tau}_\infty: \mathbb{C}\langle x_1, \dots, x_r, \bar{x}_1, \dots, \bar{x}_r \rangle \rightarrow \mathbb{C}$$

that is completely determined by (1-1). Indeed, if w is any monomial in $x_1, \dots, x_r, \bar{x}_1, \dots, \bar{x}_r$, then $\tilde{\tau}_\infty(w) = 1$ if and only if, after identifying \bar{x}_i with x_i^{-1} , w reduces to the identity in $F_r = \langle x_1, \dots, x_r \rangle$, and $\tilde{\tau}_\infty(w) = 0$ otherwise.

In the language of [Voiculescu 1991], in the limiting noncommutative probability space

$$(\mathbb{C}\langle x_1, \dots, x_r, \bar{x}_1, \dots, \bar{x}_r \rangle, \tilde{\tau}_\infty),$$

the subalgebras

$$\mathcal{A}_1 := \mathbb{C}\langle x_1, \bar{x}_1 \rangle, \quad \dots, \quad \mathcal{A}_r := \mathbb{C}\langle x_r, \bar{x}_r \rangle$$

are a *free family of subalgebras*: if $a_j \in \mathcal{A}_{i_j}$ for $j \in [q]$ with $i_1 \neq i_2 \neq \dots \neq i_q$, and $\tilde{\tau}_\infty(a_j) = 0$ for $j \in [q]$, then

$$\tilde{\tau}_\infty(a_1 a_2 \dots a_q) = 0.$$

Accordingly [Voiculescu 1991, Theorem 3.8], if $\{u_j(n) : 1 \leq j \leq r\}$ are independent Haar-random elements of $U(n)$, the family $\{\{u_j(n), u_j^*(n)\} : 1 \leq j \leq r\}$ of sets of random variables are *asymptotically free*.

Because Γ_g is not free, asymptotic freeness does not correctly capture the asymptotic behavior of the expected values $\mathbb{E}_{g,n}[\text{Tr}_\gamma]$; however, an analog of Theorem 1.3 is implied by Theorem 1.2. For $\gamma \in \Gamma_g$ let

$$\tau_{g,n}(\gamma) := \frac{1}{n} \mathbb{E}_{g,n}[\text{Tr}_\gamma].$$

Corollary 1.4 *Let $g \geq 2$, $a_1, b_1, \dots, a_g, b_g$ denote the previously fixed generators of Γ_g , and $\bar{a}_1, \bar{b}_1, \dots, \bar{a}_g, \bar{b}_g$ denote their inverses. The random variables $a_1, b_1, \dots, a_g, b_g, \bar{a}_1, \bar{b}_1, \dots, \bar{a}_g, \bar{b}_g$ in the noncommutative probability spaces $(\mathbb{C}[\Gamma_g], \tau_{g,n})$ converge in distribution as $n \rightarrow \infty$ to a limiting distribution*

$$\tilde{\tau}_{g,\infty}: \mathbb{C}\langle x_1, \dots, x_g, y_1, \dots, y_g, \bar{x}_1, \dots, \bar{x}_g, \bar{y}_1, \dots, \bar{y}_g \rangle \rightarrow \mathbb{C},$$

where x_i (resp. $y_i, \bar{x}_i, \bar{y}_i$) corresponds to a_i (resp. $b_i, \bar{a}_i, \bar{b}_i$). This can be described explicitly as follows. If w is any monomial in $x_1, \dots, x_g, y_1, \dots, y_g, \bar{x}_1, \dots, \bar{x}_g, \bar{y}_1, \dots, \bar{y}_g$, then $\tilde{\tau}_{g,\infty}(w) = 1$ if and only if w maps to the identity under the map

$$\mathbb{C}\langle x_1, \dots, x_g, y_1, \dots, y_g, \bar{x}_1, \dots, \bar{x}_g, \bar{y}_1, \dots, \bar{y}_g \rangle \rightarrow \mathbb{C}[\Gamma_g]$$

obtained by identifying x_i, y_i, \bar{x}_i and \bar{y}_i with the corresponding elements of Γ_g . If w does not map to the identity under this map, then $\tilde{\tau}_{g,\infty}(w) = 0$.

Notice that the estimate given in Theorem 1.2 is stronger than needed to establish Corollary 1.4.

1.2 Related works and further questions

The most closely related existing result to Theorem 1.2 is [Magee and Puder 2023, Theorem 1.2], which establishes Theorem 1.2 when the family of groups $SU(n)$ is replaced by the family of symmetric groups S_n , and Tr is replaced by the character fix given by the number of fixed points of a permutation. In this case, the result is phrased in terms of integrating over $\text{Hom}(\Gamma_g, S_n)$ with respect to the uniform probability measure. The corresponding result for $\text{Hom}(F_r, S_n)$ was proved much longer ago [Nica 1994].

The problem of integrating geometric functions like Tr_γ over $\mathcal{M}_{g,n}$ is also connected to the work of Mirzakhani, since, as Goldman [1984, Section 2] explains, the Atiyah–Bott–Goldman symplectic form generalizes the Weil–Petersson symplectic form on the Teichmüller space of genus g Riemann surfaces. Mirzakhani [2007] developed a method for integrating geometric functions on moduli spaces of Riemann surfaces with respect to the Weil–Petersson volume form. Although there is certainly a similarity between [loc. cit.] and the current work, here the emphasis is on $n \rightarrow \infty$, whereas [loc. cit.] caters to the regime $g \rightarrow \infty$; the target group playing the role of $SU(n)$ is always $\text{PSL}(2, \mathbb{R})$.

We now take the opportunity to mention some questions that Theorem 1.2 leads to. Voiculescu [1991] is able to boost Theorem 1.1 from a convergence in distribution result to a result on convergence in probability; that is, for any $\epsilon > 0$ and fixed $w \in F_r$, the Haar measure of the set

$$\{\phi \in \text{Hom}(F_r, U(n)) : |\text{Tr}(\phi(w))| \leq \epsilon n\}$$

tends to one as $n \rightarrow \infty$ [Voiculescu 1991, Theorem 3.9]. To do this, Voiculescu uses that the family of measure spaces $(\text{Hom}(F_r, U(n)), \mu)$ form a *Levy family* in the sense of [Gromov and Milman 1983]. This latter fact relies on an estimate for the first nonzero eigenvalue of the Laplacian on $\text{Hom}(F_r, U(n))$. It is interesting to ask whether a similar phenomenon holds for the family of measure spaces $(\mathcal{M}_{g,n}, \mu_{g,n}^{\text{ABG}})$, where $\mu_{g,n}^{\text{ABG}}$ is the probability measure corresponding to $\text{Vol}_{\mathcal{M}_{g,n}}$. The fact that $\mathcal{M}_{g,n}$ is noncompact seems to be a significant complication in answering this question using isoperimetric inequalities.

On the other hand, as pointed out to us by a referee, the results of this paper can very likely be extended to give bounds on the variance

$$\mathbb{E}_{g,n}[|\text{Tr}_\gamma|^2]$$

that can be used to improve Theorem 1.2 to the result that, for $\gamma \neq \text{id}$, the normalized traces Tr_γ/n converge in probability to zero as $n \rightarrow \infty$. To avoid adding complications to this paper, this will be pursued elsewhere.

In the prequel to this paper [Magee 2022], we proved that, for any fixed $\gamma \in \Gamma_g$, there is an infinite sequence of rational numbers $a_{-1}(\gamma), a_0(\gamma), a_1(\gamma), \dots \in \mathbb{Q}$ such that, for any $M \in \mathbb{N}$,

$$(1-2) \quad \mathbb{E}_{g,n}[\text{Tr}_\gamma] = a_{-1}(\gamma)n + a_0(\gamma) + \frac{a_1(\gamma)}{n} + \dots + \frac{a_{M-1}(\gamma)}{n^{M-1}} + O_{\gamma,M}\left(\frac{1}{n^M}\right)$$

as $n \rightarrow \infty$. Theorem 1.2 implies that $a_{-1}(\gamma) = 0$ if $\gamma \neq \text{id}$. It is also interesting to understand the other coefficients of this series. This has been accomplished when Γ_g is replaced by F_r in [Magee and Puder 2019], where in fact it is proved that

$$\mathbb{E}_{F_r, n}[\text{Tr}_w] := \int_{\text{U}(n)^r} \text{Tr}(w(x)) d\mu(x)$$

is given by a *rational* function of n and, in particular, can be expanded as in (1-2). The corresponding coefficients of the Laurent series of $\mathbb{E}_{F_r, n}[\text{Tr}_w]$ are explained in terms of Euler characteristics of subgroups of mapping class groups. One corollary is that, as $n \rightarrow \infty$,

$$(1-3) \quad \mathbb{E}_{F_r, n}[\text{Tr}_w] = O\left(\frac{1}{n^{2 \text{cl}(w)-1}}\right),$$

where $\text{cl}(w)$ is the *commutator length* of w : the minimal number of commutators that w can be written as a product of, or ∞ if $w \notin [F_r, F_r]$. We guess that an estimate like (1-3) should hold for $\mathbb{E}_{g, n}[\text{Tr}_\gamma]$, where commutator length in F_r is replaced by commutator length in Γ_g .

Another strengthening of Theorem 1.1 is the *strong asymptotic freeness* of Haar unitaries. This states that, for any complex linear combination

$$\sum_w a_w w \in \mathbb{C}[F_r],$$

almost surely with respect to Haar random $\phi \in \text{Hom}(F_r, \text{U}(n))$ as $n \rightarrow \infty$, we have

$$\left\| \sum_w a_w \phi(w) \right\| \rightarrow \left\| \sum_w a_w w \right\|_{\text{Op}(\ell^2(F_r))},$$

where the left-hand side is the operator norm on \mathbb{C}^n with standard Hermitian inner product and the norm on the right-hand side is the operator norm in the regular representation of F_r . This result was proved in [Collins and Male 2014]. It is probably very hard to extend this result to Γ_g ; the proof of Collins and Male relies on seminal work of Haagerup and Thorbjørnsen [2005] in a way that does not obviously extend to Γ_g .

We finally mention that the expected values $\mathbb{E}_{g, n}[\text{Tr}_\gamma]$ arise as a limiting form of expected values of Wilson loops in 2D Yang–Mills theory, when the coupling constant is set to zero. This will not be discussed in detail here; we refer the reader instead to the introduction of [Magee 2022]. Here we just mention the recent works [Lemoine 2022; Dahlqvist and Lemoine 2023], which make progress on related problems in the Yang–Mills setting.

1.3 Overview of paper

Here we explain the structure of the paper.

In Sections 2.1–2.5, we give some general background to the paper not depending on [Magee 2022]. In Section 2.6, we import results that we proved in the prequel and that are needed here.

At the beginning of Section 3, we state the key result (Theorem 3.1) of the remainder of the paper. To motivate things, Section 3.1 contains a discussion of why the most straightforward approach does not work, and also a discussion of what will follow instead. In the remainder of Section 3, we explain how to augment the Weingarten calculus to arrive to a formula for the key quantity $\mathcal{F}_n(w, \mu, \nu)$ (defined in Proposition 2.9) in combinatorial terms that are “good” for the next part of the argument.

Indeed, in Section 4.1 we explain how each combinatorial datum we encountered in our formula for $\mathcal{F}_n(w, \mu, \nu)$ can be used to build a decorated surface. In Corollary 4.5 we obtain a bound on $\mathcal{F}_n(w, \mu, \nu)$ in terms of the Euler characteristics of some of the surfaces that previously arose. We may restrict to certain surfaces of simplified form by performing two surgery arguments explained in Section 4.2. Given that now we have reduced estimating $\mathcal{F}_n(w, \mu, \nu)$ to estimating Euler characteristics of certain surfaces, in Section 4.3 we formulate a topological result (Proposition 4.8) which suffices to prove Theorem 3.1. Proposition 4.8 is proved in Section 4.5 using arguments related to Dehn’s algorithm and the work of Birman and Series. The necessary additional background for this proof is given in Section 4.4.

In Section 5, we show how Theorem 3.1, in conjunction with the results of [Magee 2022], proves Theorem 1.2.

1.4 Notation

We write \mathbb{N} for the natural numbers $\{1, 2, 3, \dots\}$ and $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$. We write $[n] := \{1, \dots, n\}$ for $n \in \mathbb{N}$ and $[k, l] := \{k, k + 1, \dots, l\}$ for $k, l \in \mathbb{N}$. If A and B are two sets, we write $A \setminus B$ for the elements of A not in B . If H is a group and $h_1, h_2 \in H$, we write $[h_1, h_2] := h_1 h_2 h_1^{-1} h_2^{-1}$. We let id denote the identity element of a group. We let $[H, H]$ be the subgroup of H generated by elements of the form $[h_1, h_2]$; this is called the commutator subgroup of H . If V is a complex vector space, for $q \in \mathbb{N}_0$ we let

$$V^{\otimes q} := \underbrace{V \otimes V \otimes \dots \otimes V}_q.$$

We use Vinogradov notation as follows. If f and h are functions of $n \in \mathbb{N}$, we write $f \ll h$ to mean that there are constants $n_0 \geq 0$ and $C_0 \geq 0$ such that, for $n \geq n_0$, $|f(n)| \leq C_0 h(n)$. We write $f = O(h)$ to mean $f \ll h$. We write $f \asymp h$ to mean both $f \ll h$ and $h \ll f$. If in any of these statements the implied constants depend on additional parameters, we add these parameters as subscripts to \ll , O or \asymp . Throughout the paper we view the genus g as fixed and so any implied constant may depend on g .

In this paper, Tr denotes the standard (unnormalized) trace on square complex matrices.

Acknowledgments We thank Benoît Collins, Antoine Dahlqvist, Doron Puder, Sanjaye Ramgoolam, Calum Shearer and Henry Wilton for valuable discussions about this work. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement 949143).

2 Background

2.1 Representation theory of symmetric groups

Let S_k denote the symmetric group of permutations of $[k] := \{1, \dots, k\}$, and $\mathbb{C}[S_k]$ denote its group algebra. The group S_0 is by definition the group with one element.

If we refer to $S_l \leq S_k$ with $l \leq k$, we always view S_l as the subgroup of permutations that fix every element of $[l+1, k] := \{l+1, \dots, k\}$. We write $S'_r \leq S_k$ for the subgroup of permutations that fix every element of $[k-r]$. As a consequence, we obtain fixed inclusions $\mathbb{C}[S_l] \subset \mathbb{C}[S_k]$ for l and k as above. When we write $S_l \times S_{k-l} \leq S_k$, the first factor is S_l and the second factor is S'_{k-l} .

A *Young diagram* λ is a left-aligned contiguous collection of identical square boxes in the plane such that the number of boxes in each row is nonincreasing from top to bottom. We write λ_i for the number of boxes in the i^{th} row of λ and say $\lambda \vdash k$ if λ has k boxes. We write $\ell(\lambda)$ for the number of rows of λ . For each $\lambda \vdash k$, there is a *Young subgroup*

$$S_\lambda := S_{\lambda_1} \times S_{\lambda_2} \times \cdots \times S_{\lambda_{\ell(\lambda)}} \leq S_k,$$

where the factors are subgroups in the obvious way, according to the increasing order of $[k]$.

The equivalence classes of irreducible representations of S_k are in one-to-one correspondence with Young diagrams $\lambda \vdash k$. Given λ , the construction of the corresponding irreducible representation V^λ can be done, for example, using Young symmetrizers as in [Fulton and Harris 1991, Lecture 4]. We write χ_λ for the character of S_k associated to V^λ and $d_\lambda := \chi_\lambda(\text{id}) = \dim V^\lambda$. Given $\lambda \vdash k$, the element

$$p_\lambda := \frac{d_\lambda}{k!} \sum_{\sigma \in S_k} \chi_\lambda(\sigma) \sigma \in \mathbb{C}[S_k]$$

is a central idempotent in $\mathbb{C}[S_k]$.

If G is a compact group, (ρ, W) is an irreducible representation of G , and (π, V) is any finite-dimensional representation of G , the (ρ, W) -isotypic subspace of (π, V) is the invariant subspace of V spanned by all irreducible direct summands of (π, V) that are isomorphic to (ρ, W) . When ρ and π can be inferred from W and V , we call this simply the *W-isotypic subspace of V*. If $H \leq G$ is a subgroup and (ρ, W) is an irreducible representation of H , then the W -isotypic subspace of V for H is the W -isotypic subspace of the restriction of (π, V) to H .

If (π, V) is any finite-dimensional unitary representation of S_k , and $\lambda \vdash k$, then V is also a module for $\mathbb{C}[S_k]$ by linear extension of π and $\pi(p_\lambda)$ is the orthogonal projection onto the V^λ -isotypic subspace of V .

For any compact group G , we write $(\text{triv}_G, \mathbb{C})$ for the trivial representation of G . The following lemma can be deduced for example by combining Young's rule [Fulton and Harris 1991, Corollary 4.39] with Frobenius reciprocity.

Lemma 2.1 *Let $k \in \mathbb{N}_0$ and $\lambda \vdash k$. The space of vectors in V^λ fixed by S_λ is one-dimensional.*

2.2 Representation theory of $U(n)$ and $SU(n)$

Every irreducible representation of $U(n)$ restricts to an irreducible representation of $SU(n)$, and all equivalence classes of irreducible representations of $SU(n)$ arise in this way. The equivalence classes of irreducible representations of $U(n)$ are parametrized by dominant weights, which can be thought of as nonincreasing sequences

$$\Lambda = (\Lambda_1, \dots, \Lambda_n) \in \mathbb{Z}^n,$$

also known as *signatures*. We write W^Λ for the irreducible representation of $U(n)$ corresponding to the signature Λ . Two irreducible representations of $U(n)$ restrict to the same one of $SU(n)$ if and only if their signatures differ by a constant vector. Let $\mathbb{T}(n)$ denote the maximal torus of $U(n)$ consisting of diagonal matrices. Any matrix of $\mathbb{T}(n)$ has the form $\text{diag}(\exp(i\theta_1), \dots, \exp(i\theta_n))$, where all $\theta_j \in \mathbb{R}$. Associated to the signature Λ is the character ξ_Λ of $\mathbb{T}(n)$ given by

$$\xi_\Lambda(\text{diag}(\exp(i\theta_1), \dots, \exp(i\theta_n))) := \exp\left(i\left(\sum_{j=1}^n \Lambda_j \theta_j\right)\right).$$

The highest weight theory says among other things that the ξ_Λ -isotypic subspace of W^Λ for $\mathbb{T}(n)$ is one-dimensional. Any vector in this subspace is called a *highest weight vector* of W^Λ .

Given $k, l \in \mathbb{N}_0$ and fixed Young diagrams $\mu \vdash k$ and $\nu \vdash l$, we define a family of representations of $U(n)$ as follows. For $n \geq \ell(\mu) + \ell(\nu)$, define

$$\Lambda_{\mu, \nu}(n) := (\mu_1, \mu_2, \dots, \mu_{\ell(\mu)}, \underbrace{0, \dots, 0}_{n - \ell(\mu) - \ell(\nu)}, -\nu_{\ell(\nu)}, -\nu_{\ell(\nu)-1}, \dots, -\nu_1).$$

We let $(\rho_n^{\mu, \nu}, W_n^{\mu, \nu})$ denote the irreducible representation of $U(n)$ corresponding to $\Lambda_{\mu, \nu}(n)$ when $n \geq \ell(\mu) + \ell(\nu)$. We let $D_{\mu, \nu}(n) := \dim W_n^{\mu, \nu}$ and $s_{\mu, \nu}(g) := \text{Tr}(\rho_n^{\mu, \nu}(g))$ for $g \in U(n)$. If $\mu \vdash k$ and $\nu \vdash l$, then, as $n \rightarrow \infty$,

$$(2-1) \quad D_{\mu, \nu}(n) \asymp n^{k+l}$$

by [Magee 2022, Corollary 2.3] (alternatively [Enomoto and Izumi 2016, Lemma 3.5]).

We now present a version of Schur–Weyl duality for mixed tensors due to Koike [1989]. The very definition of $U(n)$ makes \mathbb{C}^n into a unitary representation of $U(n)$ for the standard Hermitian inner product. We let $\{e_1, \dots, e_n\}$ denote the standard basis of \mathbb{C}^n . If (ρ, W) is any finite-dimensional representation of $U(n)$, we write (ρ^\vee, W^\vee) for the dual representation, where W^\vee is the space of complex linear functionals on W . The vector space $(\mathbb{C}^n)^\vee$ has a dual basis $\{\check{e}_1, \dots, \check{e}_n\}$ given by $\check{e}_j(v) := \langle v, e_j \rangle$. Throughout the paper we frequently use certain canonical isomorphisms, eg

$$((\mathbb{C}^n)^{\otimes p})^\vee \cong ((\mathbb{C}^n)^\vee)^{\otimes p}, \quad \text{End}(W) \cong W \otimes W^\vee,$$

to change points of view on representations; if we use noncanonical isomorphisms, we point them out.

Let $\mathcal{T}_n^{k,l} := (\mathbb{C}^n)^{\otimes k} \otimes ((\mathbb{C}^n)^\vee)^{\otimes l}$, with the convention that $(\mathbb{C}^n)^{\otimes 0} := \mathbb{C}$. With the natural inner product induced by that on \mathbb{C}^n , this is a unitary representation of $U(n)$ under the diagonal action and also a unitary representation of $S_k \times S_l$, where S_k acts by permuting the indices of $(\mathbb{C}^n)^{\otimes k}$ and S_l acts by permuting the indices of $((\mathbb{C}^n)^\vee)^{\otimes l}$. We write $\pi_n^{k,l}: U(n) \rightarrow \text{End}[\mathcal{T}_n^{k,l}]$ and $\rho_n^{k,l}: \mathbb{C}[S_k \times S_l] \rightarrow \text{End}[\mathcal{T}_n^{k,l}]$ for these representations. The actions of $U(n)$ and $S_k \times S_l$ on $\mathcal{T}_n^{k,l}$ commute. We use the notation, for $I = (i_1, \dots, i_k) \in [n]^k$ and $J = (j_1, \dots, j_l) \in [n]^l$,

$$e_I := e_{i_1} \otimes \dots \otimes e_{i_k} \in (\mathbb{C}^n)^{\otimes k}, \quad \check{e}_J := \check{e}_{j_1} \otimes \dots \otimes \check{e}_{j_l} \in ((\mathbb{C}^n)^\vee)^{\otimes l}, \quad e_I^J := e_I \otimes \check{e}_J \in \mathcal{T}_n^{k,l}.$$

We write $I \sqcup J$ for the concatenation $(i_1, \dots, i_k, j_1, \dots, j_l)$.

For $k, l \geq 1$, let $\dot{\mathcal{T}}_n^{k,l}$ denote the intersection of the kernels of the mixed contractions $c_{pq}: \mathcal{T}_n^{k,l} \rightarrow \mathcal{T}_n^{k-1,l-1}$ for $p \in [k]$ and $q \in [l]$ given by

$$(2-2) \quad c_{pq}(e_{i_1} \otimes \dots \otimes e_{i_k} \otimes \check{e}_{j_1} \otimes \dots \otimes \check{e}_{j_l}) \\ := \delta_{ipjq} e_{i_1} \otimes \dots \otimes e_{i_{p-1}} \otimes e_{i_{p+1}} \otimes \dots \otimes e_{i_k} \otimes \check{e}_{j_1} \otimes \dots \otimes \check{e}_{j_{q-1}} \otimes \check{e}_{j_{q+1}} \otimes \dots \otimes \check{e}_{j_l},$$

where δ_{ipjq} is the Kronecker delta. If $k = 1$ or $l = 1$, then the definition is extended in the natural way, interpreting an empty tensor of e_i or \check{e}_i as 1. If either $k = 0$ or $l = 0$, then $\dot{\mathcal{T}}_n^{k,l} = \mathcal{T}_n^{k,l}$ by convention. The space $\dot{\mathcal{T}}_n^{k,l}$ is an invariant subspace under $U(n) \times S_k \times S_l$ and hence a unitary subrepresentation of $\mathcal{T}_n^{k,l}$. On $\dot{\mathcal{T}}_n^{k,l}$ there is an analog of Schur–Weyl duality due to Koike.

Theorem 2.2 [Koike 1989, Theorem 1.1] *There is an isomorphism of unitary representations of $U(n) \times S_k \times S_l$*

$$(2-3) \quad \dot{\mathcal{T}}_n^{k,l} \cong \bigoplus_{\substack{\mu \vdash k, \nu \vdash l \\ \ell(\mu) + \ell(\nu) \leq n}} W_n^{\mu, \nu} \otimes V^\mu \otimes V^\nu.$$

Next we explain how to construct $U(n)$ -subrepresentations of $\dot{\mathcal{T}}_n^{k,l}$ isomorphic to $W_n^{\mu, \nu}$. Suppose that $\xi \in \dot{\mathcal{T}}_n^{k,l}$ is a nonzero vector such that, under the isomorphism (2-3),

$$(2-4) \quad \xi \cong w \otimes v$$

for $w \in W_n^{\mu, \nu}$ and $v \in V^\mu \otimes V^\nu$. Then $U(n) \cdot \xi$ linearly spans a $U(n)$ -subrepresentation of $\dot{\mathcal{T}}_n^{k,l}$ isomorphic to $W_n^{\mu, \nu}$. The following argument to construct such a vector ξ , given $\mu \vdash k$ and $\nu \vdash l$, appears implicitly in [Koike 1989] and is elaborated in [Benkart et al. 1994]. For $n \geq \ell(\mu) + \ell(\nu)$, let

$$(2-5) \quad \tilde{\theta}_{\mu, \nu}^n := e_1^{\otimes \mu_1} \otimes \dots \otimes e_{\ell(\mu)}^{\otimes \mu_{\ell(\mu)}} \otimes (\check{e}_n)^{\otimes \nu_1} \otimes \dots \otimes (\check{e}_{n-\ell(\nu)+1})^{\otimes \nu_{\ell(\nu)}}.$$

This vector is in the $\xi_{\mu, \nu}$ -isotypic subspace of $\dot{\mathcal{T}}_n^{k,l}$ for the maximal torus $\mathbb{T}(n)$ of $U(n)$, where $\xi_{\mu, \nu}$ is the character of $\mathbb{T}(n)$ corresponding to the highest weight in $W_n^{\mu, \nu}$.

Let $\mathfrak{p}_\mu \in \mathbb{C}[S_k]$ and $\mathfrak{p}_\nu \in \mathbb{C}[S_l]$ be the projections defined in Section 2.1. Let $\rho_n^k: S_k \rightarrow \text{End}(\mathcal{T}_n^{k,l})$ denote the representation of S_k described above and $\hat{\rho}_n^l: S_l \rightarrow \text{End}(\mathcal{T}_n^{k,l})$ that of S_l . Clearly these two

representations commute. Now let

$$(2-6) \quad \theta_{\mu, \nu}^n := \rho_n^k(\mathfrak{p}_\mu) \hat{\rho}_n^l(\mathfrak{p}_\nu) \tilde{\theta}_{\mu, \nu}^n \in \dot{\mathcal{T}}_n^{k, l}.$$

Now this is in the same isotypic subspace for $\mathbb{T}(n)$ as before since $S_k \times S_l$ commutes with $U(n)$. Moreover, it is in the subspace of $\dot{\mathcal{T}}_n^{k, l}$ corresponding to $W_n^{\mu, \nu} \otimes V^\mu \otimes V^\nu$ under the isomorphism (2-3). The intersection of the two subspaces of $\dot{\mathcal{T}}_n^{k, l}$ just discussed corresponds via (2-3) to $\mathbb{C}w \otimes V^\mu \otimes V^\nu$, where w is a highest weight vector in $W_n^{\mu, \nu}$, and hence $\theta_{\mu, \nu}^n$ takes the form of (2-4), as we desired.

Of course, we also want to know $\theta_{\mu, \nu}^n \neq 0$.

Lemma 2.3 *Suppose that $k, l \in \mathbb{N}_0$, $\mu \vdash k$, $\nu \vdash l$, and $\theta_{\mu, \nu}^n$ is as in (2-6) for $n \geq \ell(\mu) + \ell(\nu)$. We have*

$$\|\theta_{\mu, \nu}^n\|^2 = \frac{d_\mu d_\nu}{[S_k : S_\mu][S_l : S_\nu]}.$$

Proof Recall the definition of Young subgroups S_μ and S_ν from Section 2.1. Letting $\tilde{\theta} = \tilde{\theta}_{\mu, \nu}^n$ (as in (2-5)) and $\theta = \theta_{\mu, \nu}^n$, we have

$$\begin{aligned} \theta &= \rho_n^k(\mathfrak{p}_\mu) \hat{\rho}_n^l(\mathfrak{p}_\nu) \tilde{\theta} = \frac{d_\mu d_\nu}{k!l!} \sum_{\sigma = (\sigma_1, \sigma_2) \in S_k \times S_l} \chi_\mu(\sigma_1) \chi_\nu(\sigma_2) \rho_n^k(\sigma_1) \hat{\rho}_n^l(\sigma_2) \tilde{\theta} \\ &= \frac{d_\mu d_\nu}{k!l!} \sum_{\substack{[\sigma_1] \in S_k/S_\mu \\ [\sigma_2] \in S_l/S_\nu}} \left(\sum_{\tau_1 \in S_\mu} \chi_\mu(\sigma_1 \tau_1) \right) \left(\sum_{\tau_2 \in S_\nu} \chi_\nu(\sigma_2 \tau_2) \right) \rho_n^k(\sigma_1) \hat{\rho}_n^l(\sigma_2) \tilde{\theta}. \end{aligned}$$

The second equality used that $\tilde{\theta}$ is invariant under $S_\mu \times S_\nu$.

By Lemma 2.1, there is a one-dimensional subspace of invariant vectors for S_μ in V^μ . If $v_\mu \in V^\mu$ is a unit vector in this space, then

$$(2-7) \quad \sum_{\tau_1 \in S_\mu} \chi_\mu(\sigma_1 \tau_1) = |S_\mu| \langle \sigma_1 v_\mu, v_\mu \rangle.$$

Since the vectors $\rho_n^k(\sigma_1) \hat{\rho}_n^l(\sigma_2) \tilde{\theta}$ for $[\sigma_1] \in S_k/S_\mu$ and $[\sigma_2] \in S_l/S_\nu$ are orthogonal unit vectors, this gives

$$\begin{aligned} \|\theta\|^2 &= \left(\frac{d_\mu d_\nu}{k!l!} \right)^2 \sum_{\substack{[\sigma_1] \in S_k/S_\mu \\ [\sigma_2] \in S_l/S_\nu}} \left(\sum_{\tau_1 \in S_\mu} \chi_\mu(\sigma_1 \tau_1) \right)^2 \left(\sum_{\tau_2 \in S_\nu} \chi_\nu(\sigma_2 \tau_2) \right)^2 \\ &= \left(\frac{d_\mu d_\nu}{k!l!} \right)^2 |S_\mu|^2 |S_\nu|^2 \sum_{\substack{[\sigma_1] \in S_k/S_\mu \\ [\sigma_2] \in S_l/S_\nu}} |\langle \sigma_1 v_\mu, v_\mu \rangle|^2 |\langle \sigma_2 v_\nu, v_\nu \rangle|^2 \quad (\text{by (2-7)}) \\ &= \left(\frac{d_\mu d_\nu}{k!l!} \right)^2 |S_\mu| |S_\nu| \sum_{\substack{\sigma_1 \in S_k \\ \sigma_2 \in S_l}} |\langle \sigma_1 v_\mu, v_\mu \rangle|^2 |\langle \sigma_2 v_\nu, v_\nu \rangle|^2 = \frac{d_\mu d_\nu}{[S_k : S_\mu][S_l : S_\nu]}. \end{aligned}$$

The last inequality used the orthogonality relations for matrix coefficients. □

Recall that we write $\pi_n^{k,l} : U(n) \rightarrow \text{End}(\mathcal{T}_n^{k,l})$ for the diagonal representation of $U(n)$ on $\mathcal{T}_n^{k,l}$. Lemma 2.3 implies that $\theta_{\mu,v}^n$ is a nonzero vector. By the remarks following (2-6), it is of the pure tensor form $w \otimes v$ under the Schur–Weyl isomorphism (2-3), with $w \in W_n^{\mu,v}$, and hence we obtain the following corollary:

Corollary 2.4 *Suppose $n \geq \ell(\mu) + \ell(v)$. The subspace*

$$W_n(\theta_{\mu,v}^n) := \text{span}\{\pi_n^{k,l}(u)\theta_{\mu,v}^n : u \in U(n)\} \subset \dot{\mathcal{T}}_n^{k,l}$$

is, under $\pi_n^{k,l}$, a $U(n)$ -subrepresentation of $\dot{\mathcal{T}}_n^{k,l}$ isomorphic to $W_n^{\mu,v}$.

2.3 The Weingarten calculus

The Weingarten calculus is a method based on Schur–Weyl duality that allows one to calculate integrals of products of matrix coefficients in the defining representation of $U(n)$ in terms of sums over permutations. It was discovered initially by Weingarten [1978], and developed further in [Xu 1997; Collins 2003; Collins and Śniady 2006].

We present two formulations of the Weingarten calculus. Given $k \in \mathbb{N}$ and $n \in \mathbb{N}$, the *Weingarten function* with parameters n and k is the element² of $\mathbb{C}[S_k]$ [Collins and Śniady 2006, equation (9)]

$$(2-8) \quad \text{Wg}_{n,k} := \frac{1}{(k!)^2} \sum_{\substack{\lambda \vdash k \\ \ell(\lambda) \leq n}} \frac{d_\lambda^2}{D_\lambda(n)} \sum_{\sigma \in S_k} \chi_\lambda(\sigma)\sigma.$$

We write $\text{Wg}_{n,k}(\sigma)$ for the coefficient of σ in (2-8). The following theorem was proved by Collins and Śniady [2006, Corollary 2.4]:

Theorem 2.5 *For $k \in \mathbb{N}$ and for $i_1, i'_1, j_k, j'_k, \dots, i_k, i'_k, j_k, j'_k \in [n]$,*

$$(2-9) \quad \int_{u \in U(n)} u_{i_1 j_1} \cdots u_{i_k j_k} \bar{u}_{i'_1 j'_1} \cdots \bar{u}_{i'_k j'_k} d\mu(u) = \sum_{\sigma, \tau \in S_k} \delta_{i_1 i'_{\sigma(1)}} \cdots \delta_{i_k i'_{\sigma(k)}} \delta_{j_1 j'_{\tau(1)}} \cdots \delta_{j_k j'_{\tau(k)}} \text{Wg}_{n,k}(\tau\sigma^{-1}),$$

where δ_{pq} is the Kronecker delta function.

It is sometimes more flexible to reformulate Theorem 2.5 in terms of projections. Here $u \in U(n)$ acts on $A \in \text{End}((\mathbb{C}^n)^{\otimes k})$ by $A \mapsto \pi_n^k(u)A\pi_n^k(u^{-1})$, where $\pi_n^k : U(n) \rightarrow \text{End}((\mathbb{C}^n)^{\otimes k})$ is the diagonal action. Write $P_{n,k}$ for the orthogonal projection in $\text{End}((\mathbb{C}^n)^{\otimes k})$ onto the $U(n)$ -invariant vectors. The following proposition is due to [Collins and Śniady 2006, Proposition 2.3]:

²Although not relevant here, classically the Weingarten function arises as the multiplicative inverse of $\sum_{\sigma \in S_k} n^{\#\text{cycles}(\sigma)}\sigma$ in $\mathbb{C}[S_k]$ whenever $n \geq k$.

Proposition 2.6 (Collins and Śniady) *Let $n, k \in \mathbb{N}$. Suppose $A \in \text{End}((\mathbb{C}^n)^{\otimes k})$. Then*

$$P_{n,k}[A] = \rho_n^k(\Phi[A] \cdot W_{g_{n,k}}),$$

where

$$\Phi[A] := \sum_{\sigma \in S_k} \text{Tr}(A \rho_n^k(\sigma^{-1})) \sigma.$$

Later we will need the following bound for the Weingarten function due to [Collins and Śniady 2006, Proposition 2.6]. For a permutation σ , let $|\sigma|$ denote the minimum number of transpositions that σ can be written as a product of.

Proposition 2.7 *For any fixed $\sigma \in S_k$, $W_{g_{n,k}}(\sigma) \ll_k n^{-k-|\sigma|}$ as $n \rightarrow \infty$.*

2.4 Free groups and surface groups

Let $F_{2g} := \langle a_1, b_1, \dots, a_g, b_g \rangle$ be the free group on $2g$ generators $a_1, b_1, \dots, a_g, b_g$ and $R_g := [a_1, b_1] \cdots [a_g, b_g] \in F_{2g}$. There is a quotient map $F_{2g} \rightarrow \Gamma_g$ given by reduction modulo R_g . We say that $w \in F_{2g}$ represents the conjugacy class of $\gamma \in \Gamma_g$ if the projection of w to Γ_g is in the conjugacy class of γ in Γ_g .

Given $w \in F_{2g}$, we view w as a combinatorial word in $a_1, a_1^{-1}, b_1, b_1^{-1}, \dots, a_g, a_g^{-1}, b_g, b_g^{-1}$ by writing it in reduced (shortest) form; ie a_1 does not follow a_1^{-1} etc. We say that w is *cyclically reduced* if the first letter of its reduced word is not the inverse of the last letter. The length $|w|$ of $w \in F_{2g}$ is the length of its reduced form word. We say $w \in F_{2g}$ is a *shortest element* representing the conjugacy class of $\gamma \in \Gamma_g$ if it has minimal length among all elements representing the conjugacy class of γ . If w is a shortest element representing some conjugacy class in Γ_g , then w is cyclically reduced.

For any group H , the commutator subgroup $[H, H] \leq H$ is the subgroup generated by all elements of the form $[h_1, h_2] := h_1 h_2 h_1^{-1} h_2^{-1}$ with $h_1, h_2 \in H$. If $\gamma \in [\Gamma_g, \Gamma_g]$ and w represents the conjugacy class of γ , then $w \in [F_{2g}, F_{2g}]$ (see [Magee 2022, Section 2.6]).

2.5 Witten zeta functions

Witten zeta functions appeared first in [Witten 1991] and were named by Zagier [1994]. The *Witten zeta function of $SU(n)$* is defined, for s in a half-plane of convergence, by

$$(2-10) \quad \zeta(s; n) := \sum_{(\rho, W) \in \widehat{SU(n)}} \frac{1}{(\dim W)^s},$$

where $\widehat{SU(n)}$ denotes the equivalence classes of irreducible representations of $SU(n)$. Indeed, the series (2-10) converges for $\text{Re}(s) > 2/n$ by [Larsen and Lubotzky 2008, Theorem 5.1] (see also [Häsä and Stasinski 2019, Section 2]). Also relevant to this work is a result of Guralnick, Larsen and Manack [Guralnick et al. 2012, Theorem 2 and equation (7)], which states, for fixed $s > 0$,

$$(2-11) \quad \lim_{n \rightarrow \infty} \zeta(s; n) = 1.$$

2.6 Results of the prequel paper

By [Magee 2022, Proposition 1.5], if $\gamma \notin [\Gamma_g, \Gamma_g]$, then $\mathbb{E}_{g,n}[\text{Tr}_\gamma] = 0$ for $n \geq n_0(\gamma)$. This proves Theorem 1.2 in this case. Hence, in the rest of the paper we need only consider $\gamma \in [\Gamma_g, \Gamma_g]$ and hence $w \in [F_{2g}, F_{2g}]$ if $w \in F_{2g}$ represents the conjugacy class of γ .

For each $w \in F_{2g}$, we have a word map $w : U(n)^{2g} \rightarrow U(n)$ obtained by substituting matrices for the generators of F_{2g} . For example, if $u_1, v_1, \dots, u_g, v_g \in U(n)$ then $R_g(u_1, v_1, \dots, u_g, v_g) = [u_1, v_1] \cdots [u_g, v_g]$. We begin with the following result from [Magee 2022, Corollary 1.8]:

Proposition 2.8 *Suppose that $g \geq 2$, $\gamma \in \Gamma_g$, and $w \in F_{2g}$ represents the conjugacy class of γ . For any $B \in \mathbb{N}$, we have, as $n \rightarrow \infty$,*

$$(2-12) \quad \mathbb{E}_{g,n}[\text{Tr}_\gamma] = \zeta(2g - 2; n)^{-1} \sum_{\substack{\mu, \nu \text{ Young diagrams} \\ \ell(\mu), \ell(\nu) \leq B \\ \mu_1, \nu_1 \leq B^2}} D_{\mu, \nu}(n) \mathcal{F}_n(w, \mu, \nu) + O_{B,w,g}(n^{|w|} n^{-2 \log B}),$$

where

$$(2-13) \quad \mathcal{F}_n(w, \mu, \nu) := \int_{\text{SU}(n)^{2g}} \text{Tr}(w(x)) \overline{s_{\mu, \nu}(R_g(x))} d\mu(x).$$

Notice that, for $n \geq 2B$, the right-hand side of (2-12) makes sense, ie $D_{\mu, \nu}, s_{\mu, \nu}$ are well defined. We also have the following proposition, which follows from [Magee 2022, Proposition 3.1] together with $\bar{s}_{\mu, \nu} = s_{\nu, \mu}$:

Proposition 2.9 *Let $w \in [F_{2g}, F_{2g}]$. Then, for any fixed μ, ν and $n \geq \ell(\mu) + \ell(\nu)$,*

$$\mathcal{F}_n(w, \mu, \nu) = \mathcal{F}_n(w, \nu, \mu) := \int_{U(n)^{2g}} \text{Tr}(w(x)) s_{\nu, \mu}(R_g(x)) d\mu(x).$$

This is convenient as it will allow us to use the Weingarten calculus directly as it is presented in Section 2.3 for $U(n)$ rather than $SU(n)$. By using Proposition 2.9, taking a representative $w \in F_{2g}$ of the conjugacy class of γ and taking B such that $|w| - 2 \log B \leq -1$ in Proposition 2.8, we obtain the following result, from which we begin the new arguments of this paper:

Corollary 2.10 *Let $\gamma \in [\Gamma_g, \Gamma_g]$ and $w \in [F_{2g}, F_{2g}]$ be a representative of the conjugacy class of $\gamma \in \Gamma$. Then there exists a finite set $\tilde{\Omega}$ of pairs (μ, ν) of Young diagrams such that*

$$\mathbb{E}_{g,n}[\text{Tr}_\gamma] = \zeta(2g - 2; n)^{-1} \sum_{(\mu, \nu) \in \tilde{\Omega}} D_{\mu, \nu}(n) \mathcal{F}_n(w, \nu, \mu) + O_{w,g}\left(\frac{1}{n}\right).$$

As we know $\lim_{n \rightarrow \infty} \zeta(2g - 2, n) = 1$ by (2-11), we have now reduced the proof of Theorem 1.2 to establishing suitable bounds for the integrals $\mathcal{F}_n(w, \mu, \nu)$, where we can view μ and ν as fixed Young diagrams since $\tilde{\Omega}$ is finite.

3 Combinatorial integration

3.1 Setup and motivation

The main result of the rest of the paper is the following:

Theorem 3.1 *Let $\gamma \in \Gamma_g$ with $\gamma \neq \text{id}$. Let $w \in F_{2g}$ be a shortest element representing the conjugacy class of γ . For each $k, l \in \mathbb{N}_0$, there is a constant $C(w, k, l) > 0$ such that, for any $\mu \vdash k, \nu \vdash l$*

$$|D_{\mu,\nu}(n)\mathcal{F}_n(w, \mu, \nu)| \leq C(w, k, l)$$

for all $n \in \mathbb{N}$.

Accordingly, since we know the large n behavior of $D_{\mu,\nu}(n)$ from (2-1), in this section we wish to estimate

$$\mathcal{F}_n(w, \mu, \nu) = \int_{U(n)^{2g}} \text{Tr}(w(x))s_{\mu,\nu}(R_g(x)) d\mu(x)$$

for fixed $\mu \vdash k, \nu \vdash l$.

What doesn't work We begin by discussing why the most straightforward approach to this problem leads to serious complications. It is possible to approach the problem by writing $s_{\mu,\nu}(h)$ as a fixed finite linear combination of functions

$$p_{\mu'}(h)p_{\nu'}(h^{-1}),$$

where $p_{\mu'}(h)$ (resp. $p_{\nu'}(h^{-1})$) is a power sum symmetric polynomial of the eigenvalues of h (resp. h^{-1} or \bar{h}). See for example [Magee 2022, Section 3.3] for one way to do this. The coefficients of this expansion are fixed, but not transparent, since they involve Littlewood–Richardson coefficients. In any case, this approach leads to writing $\mathcal{F}_n(w, \mu, \nu)$ as a finite linear combination of integrals of the form

$$(3-1) \quad \int_{U(n)^{2g}} \text{Tr}(w(x))\text{Tr}(R_g(x)^{k_1}) \cdots \text{Tr}(R_g(x)^{k_p})\text{Tr}(R_g(x)^{-l_1}) \cdots \text{Tr}(R_g(x)^{-l_q}) d\mu(x),$$

where $\sum k_j = |\mu|$ and $\sum l_j = |\nu|$.

Magee and Puder [2019] give a full asymptotic expansion for (3-1) as $n \rightarrow \infty$. However, these estimates are not sufficient for the current paper and, to motivate the rest of this section, we explain briefly the issues involved. However, this discussion is not needed to understand the arguments that we will make to prove Theorem 3.1.

The main result of [Magee and Puder 2019] gives a full “genus” expansion of (3-1) in terms of surfaces and maps on surfaces dictated by $w \in F_{2g}$. Roughly speaking, every term in this expansion comes from a homotopy class of map f from an orientable surface Σ_f to $\bigvee_{i=1}^{2g} S^1$; to contribute to (3-1) the surface Σ_f has one boundary component that maps to w at the level of the fundamental groups, p boundary components that map respectively to $R_g^{k_1}, \dots, R_g^{k_p}$ at the level of fundamental groups, and q

boundary components that map respectively to $R_g^{-l_1}, \dots, R_g^{-l_q}$ at the level of fundamental groups. The contribution of the pair (f, Σ_f) to (3-1) is of the form $c(f, \Sigma_f)n^{\chi(\Sigma_f)}$; the coefficient $c(f, \Sigma_f)$ is an Euler characteristic of a symmetry group of (f, Σ_f) and is not easy to calculate in general. However, one could still hope to get decay of (3-1) by controlling the possible $\chi(\Sigma_f)$ that could appear.

There are two issues with this. The first one is that, if w is not the shortest element representing the conjugacy class of γ , then we get bounds that are not helpful. For a very simple example, let $w = R_g^l$ and $\gamma = \text{id}_{\Gamma_g}$, and consider the potential contribution from $p = 0, q = 1$ and $l_1 = l$. Then, for any ν with $|\nu| = l$, there is contribution to $\mathcal{F}_n(w, \emptyset, \nu)$ that is a multiple of

$$\int_{\text{U}(n)^{2g}} \text{Tr}(R_g(x)^l)\text{Tr}(R_g(x)^{-l}) d\mu(x).$$

Here, in the theory of [Magee and Puder 2019], there is a (Σ_f, f) that is an annulus, one boundary component corresponding to $w = R_g^l$ and one corresponding to R_g^{-l} , so we can only bound the corresponding contribution to $D_{\emptyset, \nu}(n)\mathcal{F}_n(w, \emptyset, \nu)$ by using [Magee and Puder 2019] on the order of $D_{\emptyset, \nu}(n) \asymp n^l$. On the other hand, any approach that works to establish Theorem 3.1 (for $\gamma \neq \text{id}$) should extend to show that, when $\gamma = \text{id}$, $D_{\emptyset, \nu}(n)\mathcal{F}_n(w, \emptyset, \nu) \ll n$ as $\mathbb{E}_{g, n}[\text{Tr}_{\text{id}}] = n$.

Indeed, this phenomenon extends to words of the form $w_0 R_g^l$ and more generally to words that are not shortest representatives of some conjugacy class in Γ_g . It means that, even if we use something similar in spirit to [Magee and Puder 2019], to prove Theorem 3.1 we must incorporate the theory of shortest representative words. This indeed takes place in Sections 4.3–4.5; the topological result proved there hinges on this theory.

The second issue is a little more subtle and only appears for “mixed” representations, ie both $\mu, \nu \neq \emptyset$. In this case, suppose w is a shortest element representing some conjugacy class in Γ_g and $w \in [F_{2g}, F_{2g}]$. This means that there is a pair (f_0, Σ_{f_0}) where Σ_{f_0} has one boundary component that maps to w at the level of the fundamental groups. Let us take $\mu, \nu = (k), (k)$, ie each Young diagram has one row of k boxes. This means we get a potential contribution to $D_{\mu, \nu}(n)\mathcal{F}_n(w, \mu, \nu)$ that is a constant multiple of

$$(3-2) \quad D_{(k), (k)}(n) \int_{\text{U}(n)^{2g}} \text{Tr}(w(x))\text{Tr}(R_g(x)^k)\text{Tr}(R_g(x)^{-k}) d\mu(x).$$

Now, for every $k \in \mathbb{N}$, there is (f, Σ_f) contributing to (3-2) with one component that is (f_0, Σ_{f_0}) and the other an annulus with boundary components corresponding to R_g^k and R_g^{-k} . Since the annulus has Euler characteristic 0, and $D_{(k), (k)} \asymp n^{2k}$, the order of this contribution to $D_{(k), (k)}(n)\mathcal{F}_n(w, (k), (k))$ is potentially $\gg n^{2k}n^{\chi(\Sigma_{f_0})}$. For large enough k , the exponent here is arbitrarily large, which is clearly catastrophic. In reality, this contribution must cancel with some other contribution, but we do not know how to see these cancellations.

This ends the discussion of the difficulties of the most straightforward approach to the problems of this paper.

What does work To bypass the previous issues we produce a refined version of the Weingarten calculus that leads to a restricted set of surfaces, for instance not including the ones causing the problem above as well as all generalizations of this issue.

The basic approach is the following. Instead of trying to deal with a complicated formula for $s_{\mu,\nu}(R_2(x))$ (as above), we instead use the copy $W_n(\theta_{\mu,\nu}^n)$ of $W_n^{\mu,\nu}$ in $\dot{\mathcal{J}}_n^{k,l}$ that we found in Corollary 2.4. In Section 3.3, we compute the orthogonal projection q_θ from $\mathcal{F}_n^{k,l}$ (note: not $\dot{\mathcal{J}}_n^{k,l}$) onto $W_n(\theta_{\mu,\nu}^n)$ (Proposition 3.2). In the formula we obtain, we give bounds on the coefficients appearing therein (Lemma 3.3). In addition, we remember that $q_\theta \in \text{End}(\dot{\mathcal{J}}_n^{k,l})$; this fact is not obvious from our formula but turns out to be vital going forward.

The calculation of q_θ is extra to, but in the same spirit as, the vanilla Weingarten calculus, which is why we claim to have refined the Weingarten calculus here.

In the expression for $\mathcal{F}_n(w, \mu, \nu)$, we now write

$$s_{\mu,\nu}(R_2(x)) = \text{Tr}_{\mathcal{F}_n^{k,l}}(Aq_\theta Bq_\theta A^{-1}q_\theta B^{-1}q_\theta Cq_\theta Dq_\theta C^{-1}q_\theta D^{-1}q_\theta),$$

where A, B, C and D are the images of the generators of Γ_2 under x . Then the entire integral of $\text{Tr}(w(x))s_{\mu,\nu}(R_2(x))$ is done using the usual Weingarten calculus. The fact that $q_\theta \in \text{End}(\dot{\mathcal{J}}_n^{k,l})$ intervenes at a critical point to show that certain contributions from the classical Weingarten calculus cancel and lead to restrictions on the nonzero contributions. Precisely, the restriction we obtain is summarized in the *forbidden matching* property below (Section 3.4) and property (P4) (Section 4.3).

3.2 Proof of Theorem 3.1 when $k = l = 0$

Here we give a proof of Theorem 3.1 when $k = l = 0$. This will allow us to bypass the slightly confusing issue of using the Weingarten function $W_{g_n,k+l}$ when $k + l = 0$ in Section 3.3.

If $k = l = 0$, then the only possible $\mu \vdash k$ and $\nu \vdash l$ are empty Young diagrams $\mu = \nu = \emptyset$, and $W_n^{\emptyset,\emptyset}$ is the trivial representation of $U(n)$, so $D_{\emptyset,\emptyset}(n) = 1$ for all $n \geq 1$ and $s_{\emptyset,\emptyset}(h) = 1$ for all $h \in U(n)$. We then have

$$(3-3) \quad D_{\emptyset,\emptyset}(n)\mathcal{F}_n(w, \emptyset, \emptyset) = \mathcal{F}_n(w, \emptyset, \emptyset) = \int_{U(n)^{2g}} \text{Tr}(w(x)) d\mu(x).$$

If $w \in F_{2g}$ is a cyclically shortest word representing the conjugacy class of $\gamma \in \Gamma_g$ with $\gamma \neq \text{id}$, then $w \neq \text{id}$. It then follows from (1-1) that $D_{\emptyset,\emptyset}(n)\mathcal{F}_n(w, \emptyset, \emptyset) = o_w(n)$ as $n \rightarrow \infty$, but, in fact, (3-3) is given by a rational function of n for $n \geq n_0(w)$ by a straightforward application of the Weingarten calculus [Magee and Puder 2019]. This implies $D_{\emptyset,\emptyset}(n)\mathcal{F}_n(w, \emptyset, \emptyset) = O_w(1)$ as $n \rightarrow \infty$, as required.

This proves Theorem 3.1 when $k = l = 0$. Hence, in the rest of Section 3, we can assume $k + l > 0$.

3.3 A projection formula

Here we develop an integral calculus that is more powerful than the usual Weingarten calculus and allows us to directly tackle $\mathcal{J}_n(w, \mu, \nu)$ without writing it in terms of integrals as in (3-1). The key point is that our method leads to the *forbidden matchings* property of Section 3.4 and property (P4) of Section 4.3.

We now view $k, l, \mu \vdash k$ and $\nu \vdash l$ as fixed, assume $k + l > 0$ and $n \geq \ell(\mu) + \ell(\nu)$, and write $\theta = \theta_{\mu, \nu}^n$ as in (2-6), suppressing the dependence on n . Let $W_n(\theta)$ be defined as in Corollary 2.4. Thus $W_n(\theta)$ is an irreducible summand of $\dot{\mathcal{J}}_n^{k, l}$ isomorphic to $W_n^{\mu, \nu}$ for the group $U(n)$.

In the remainder of the paper we drop the dependence of our notation on n whenever it adds clarity.

Our first task is to compute the orthogonal projection q_θ onto $W(\theta)$. Let P_θ denote the orthogonal projection in $\mathcal{T}_n^{k, l}$ onto θ . We also view P_θ as an element of $\text{End}(\dot{\mathcal{J}}_n^{k, l})$ by restriction.

Under the canonical isomorphism $\text{End}(\dot{\mathcal{J}}_n^{k, l}) \cong \dot{\mathcal{J}}_n^{k, l} \otimes (\dot{\mathcal{J}}_n^{k, l})^\vee$, we have $P_\theta \cong (\theta \otimes \theta^\vee) / \|\theta\|^2$, and also, from (2-6),

$$(3-4) \quad P_\theta = \frac{1}{\|\theta\|^2} \rho^k(\mathfrak{p}_\mu) \hat{\rho}^l(\mathfrak{p}_\nu) [\tilde{\theta}_{\mu, \nu} \otimes \tilde{\theta}_{\mu, \nu}^\vee] \rho^k(\mathfrak{p}_\mu) \hat{\rho}^l(\mathfrak{p}_\nu);$$

here the inner square bracket is interpreted as an element of $\text{End}(\dot{\mathcal{J}}_n^{k, l})$. By Schur’s lemma, we have

$$(3-5) \quad q_\theta = D_{\mu, \nu}(n) \int_{h \in U(n)} \pi(h) P_\theta \pi(h^{-1}) d\mu(h)$$

since the right-hand side is an element of $\text{End}(W(\theta)) \subset \text{End}(\mathcal{T}_n^{k, l})$ that commutes with $\pi^{k, l}(U(n))$, so it is a multiple of q_θ , and it has the correct trace.

On the other hand, we can view $\mathcal{T}_n^{k, l} \otimes (\dot{\mathcal{J}}_n^{k, l})^\vee \cong \mathcal{T}_n^{k+l, k+l}$ by the canonical isomorphism

$$\mathcal{T}_n^{k, l} \otimes (\dot{\mathcal{J}}_n^{k, l})^\vee \cong (\mathbb{C}^n)^{\otimes k} \otimes ((\mathbb{C}^n)^{\otimes l})^\vee \otimes ((\mathbb{C}^n)^{\otimes k})^\vee \otimes (\mathbb{C}^n)^{\otimes l}$$

followed by the fixed isomorphism

$$(3-6) \quad \varphi: e_I^J \otimes \check{e}_{I'}^{J'} \mapsto e_{I \sqcup J'} \otimes \check{e}_{I' \sqcup J}.$$

Finally, there is a canonical isomorphism $\mathcal{T}_n^{k+l, k+l} \cong \text{End}((\mathbb{C}^n)^{\otimes k+l})$. So, combining these, we fix isomorphisms

$$(3-7) \quad \text{End}(\mathcal{T}_n^{k, l}) \cong \dot{\mathcal{J}}_n^{k, l} \otimes (\dot{\mathcal{J}}_n^{k, l})^\vee \xrightarrow{\varphi} \mathcal{T}_n^{k+l, k+l} \cong \text{End}((\mathbb{C}^n)^{\otimes k+l}).$$

We view the outer two isomorphisms as fixed identifications. These isomorphisms are of unitary representations of $U(n)$ when everything is given its natural inner product. Moreover, for $\sigma = (\sigma_1, \sigma_2) \in S_k \times S_l$ and $\tau = (\tau_1, \tau_2) \in S_k \times S_l$, we have, for $A \in \text{End}(\mathcal{T}_n^{k, l})$,

$$(3-8) \quad \varphi[\rho^k(\sigma_1) \hat{\rho}^l(\sigma_2) A \rho^k(\tau_1) \hat{\rho}^l(\tau_2)] = \rho^{k+l}(\sigma_1, \tau_2^{-1}) \varphi[A] \rho^{k+l}(\tau_1, \sigma_2^{-1}),$$

recalling that $\rho^{k+l}: \mathbb{C}[S_{k+l}] \rightarrow \text{End}((\mathbb{C}^n)^{\otimes k+l})$ is the representation by permuting coordinates.

We now return to the calculation of q_θ in (3-5). We have

$$(3-9) \quad q_\theta = D_{\mu,v}(n)\varphi^{-1}[P_{n,k+l}[\varphi(P_\theta)]],$$

where $P_{n,k+l}$ is the projection onto the $U(n)$ -invariant vectors (by conjugation) in $\text{End}((\mathbb{C}^n)^{\otimes k+l})$. This can now be done using the classical Weingarten calculus. By Proposition 2.6, we have

$$(3-10) \quad P_{n,k+l}[\varphi(P_\theta)] = \rho^{k+l}(\Phi[\varphi(P_\theta)] \cdot \text{Wg}_{n,k+l}),$$

where

$$\Phi[\varphi(P_\theta)] = \sum_{\sigma \in S_{k+l}} \text{Tr}(\varphi(P_\theta)\rho^{k+l}(\sigma^{-1}))\sigma.$$

By (3-8) and (3-4), and since $\text{eg } \chi_\mu(g) = \chi_\mu(g^{-1})$, we obtain

$$\begin{aligned} \varphi(P_\theta) &= \frac{1}{\|\theta\|^2} \varphi(\rho^k(\mathfrak{p}_\mu)\hat{\rho}^l(\mathfrak{p}_v)[\tilde{\theta}_{\mu,v} \otimes \tilde{\theta}_{\mu,v}^\vee]\rho^k(\mathfrak{p}_\mu)\hat{\rho}^l(\mathfrak{p}_v)) \\ &= \frac{1}{\|\theta\|^2} \rho^{k+l}(\mathfrak{p}_{\mu \otimes v})\varphi(\tilde{\theta}_{\mu,v} \otimes \tilde{\theta}_{\mu,v}^\vee)\rho^{k+l}(\mathfrak{p}_{\mu \otimes v}), \end{aligned}$$

where

$$\mathfrak{p}_{\mu \otimes v} := \frac{d_\mu d_v}{k!l!} \sum_{\sigma=(\sigma_1, \sigma_2) \in S_k \times S_l} \chi_\mu(\sigma_1)\chi_v(\sigma_2)\sigma \in \mathbb{C}[S_{k+l}].$$

Now, using that Φ is a $\mathbb{C}[S_{k+l}]$ -bimodule morphism [Collins and Śniady 2006, Proposition 2.3 (1)], we obtain

$$\begin{aligned} \Phi[\varphi(P_\theta)] &= \frac{1}{\|\theta\|^2} \mathfrak{p}_{\mu \otimes v} \Phi[\varphi(\tilde{\theta}_{\mu,v} \otimes \tilde{\theta}_{\mu,v}^\vee)]\mathfrak{p}_{\mu \otimes v} \\ &= \frac{1}{\|\theta\|^2} \mathfrak{p}_{\mu \otimes v} \left(\sum_{\sigma \in S_{k+l}} \text{Tr}(\varphi(\tilde{\theta}_{\mu,v} \otimes \tilde{\theta}_{\mu,v}^\vee)\rho^{k+l}(\sigma^{-1}))\sigma \right) \mathfrak{p}_{\mu \otimes v}. \end{aligned}$$

Now, $\text{Tr}(\varphi(\tilde{\theta}_{\mu,v} \otimes \tilde{\theta}_{\mu,v}^\vee)\rho^{k+l}(\sigma^{-1}))$ is equal to 1 if and only if σ is in $S_\mu \times S_v \leq S_k \times S_l$, and is 0 otherwise. So we obtain

$$\Phi[\varphi(P_\theta)] = \frac{1}{\|\theta\|^2} \mathfrak{p}_{\mu \otimes v} \left(\sum_{\sigma \in S_\mu \times S_v} \sigma \right) \mathfrak{p}_{\mu \otimes v};$$

hence, from (3-10),

$$P_{n,k+l}[\varphi(P_\theta)] = \rho^{k+l}(z_\theta),$$

where

$$(3-11) \quad z_\theta := \sum_{\tau \in S_{k+l}} z_\theta(\tau)\tau := \frac{1}{\|\theta\|^2} \mathfrak{p}_{\mu \otimes v} \left(\sum_{\sigma \in S_\mu \times S_v} \sigma \right) \mathfrak{p}_{\mu \otimes v} \text{Wg}_{n,k+l} \in \mathbb{C}[S_{k+l}].$$

Therefore we obtain the following proposition:

Proposition 3.2
$$q_\theta = D_{\mu,v}(n)\varphi^{-1}[\rho^{k+l}(z_\theta)].$$

We can use the bound for the coefficients of $W_{g_{n,k+l}}$ from Proposition 2.7 to infer a bound on the coefficients $z_\theta(\tau)$. For $\sigma \in S_{k+l}$, let $\|\sigma\|_{k,l}$ denote the minimum m for which

$$\sigma = \sigma_0 t_1 t_2 \cdots t_m,$$

where $\sigma_0 \in S_k \times S_l$ and t_1, \dots, t_m are transpositions in S_{k+l} .

Lemma 3.3 For all $\tau \in S_{k+l}$ and $\theta = \theta_{\mu,\nu}$ as above, $z_\theta(\tau) = O_{k,l}(n^{-k-l-\|\tau\|_{k,l}})$ as $n \rightarrow \infty$.

Proof Referring to (3-11), as $n \rightarrow \infty$, $\|\theta\|^{-2} = O_{k,l}(1)$ by Lemma 2.3 and the coefficients of $\mathfrak{p}_{\mu \otimes \nu}(\sum_{\sigma \in S_\mu \times S_\nu} \sigma) \mathfrak{p}_{\mu \otimes \nu}$ are clearly $O_{k,l}(1)$, so z_θ has the form

$$\left(\sum_{\sigma \in S_k \times S_l} A(\sigma) \sigma \right) W_{g_{n,k+l}},$$

where each $A(\sigma)$ is $O_{k,l}(1)$. This means

$$z_\theta(\tau) = \sum_{\substack{\sigma \in S_k \times S_l \\ \sigma' \in S_{k+l} \\ \sigma \sigma' = \tau}} A(\sigma) W_{g_{n,k+l}}(\sigma').$$

The order of any of the finitely many summands above is $n^{-k-l-|\sigma'|}$ by Proposition 2.7, and the minimum possible value of $|\sigma'|$ is $\|\tau\|_{k,l}$. □

Before moving on, it is useful to explain the operator $\varphi^{-1}[\rho^{k+l}(\pi)]$ for $\pi \in S_{k+l}$. For $I = (i_1, \dots, i_{k+l})$, let $I'(I; \pi) := i_{\pi(1)}, \dots, i_{\pi(k)}$ and $J'(I; \pi) := i_{\pi(k+1)}, \dots, i_{\pi(k+l)}$. As an element of

$$(\mathbb{C}^n)^{\otimes k+l} \otimes ((\mathbb{C}^n)^\vee)^{\otimes k+l},$$

$\rho^{k+l}(\pi)$ is given by

$$\sum_{\substack{I=(i_1, \dots, i_k) \\ J=(j_{k+1}, \dots, j_{k+l})}} e^{I'(I \sqcup J; \pi) \sqcup J'(I \sqcup J; \pi)} \otimes \check{e}_{I \sqcup J},$$

so, from (3-6),

$$(3-12) \quad \varphi^{-1}[\rho^{k+l}(\pi)] = \sum_{\substack{I=(i_1, \dots, i_k) \\ J=(j_{k+1}, \dots, j_{k+l})}} e^{J_{I'(I \sqcup J; \pi)}} \otimes \check{e}_I^{J'(I \sqcup J; \pi)}.$$

3.4 A combinatorial integration formula

In this rest of Section 3, we assume $g = 2$. All proofs extend to $g \geq 3$. We write $\{a, b, c, d\}$ for the generators of F_4 and $R := [a, b][c, d]$. Assume both γ and w are not the identity and $w \in [F_4, F_4]$ according to the remarks at the beginning of Section 2.6. We write w in the reduced form

$$(3-13) \quad w = f_1^{\epsilon_1} f_2^{\epsilon_2} \cdots f_{|w|}^{\epsilon_{|w|}}, \quad \epsilon_u \in \{\pm 1\}, f_u \in \{a, b, c, d\},$$

where, if $f_u = f_{u+1}$, then $\epsilon_u = \epsilon_{u+1}$. For $f \in \{a, b, c, d\}$, let p_f denote the number of occurrences of f^{+1} in (3-13). The expression (3-13) implies that, for $h := (h_a, h_b, h_c, h_d) \in U(n)^4$,

$$(3-14) \quad \text{Tr}(w(h)) = \sum_{i_j \in [n]} (h_{f_1}^{\epsilon_1})_{i_1 i_2} (h_{f_2}^{\epsilon_2})_{i_2 i_3} \cdots (h_{f_{|w|}}^{\epsilon_{|w|}})_{i_{|w|} i_1}.$$

Working with this expression will be cumbersome so we explain a diagrammatic way to think about (3-14). This will be the starting point for how we eventually understand $\mathcal{F}_n(w, \mu, \nu)$ in terms of decorated surfaces. We begin with a collection of intervals as follows:

w-intervals and the w-loop Firstly, for every $j \in [w]$ with $f_j = f$ as in (3-13) and $\epsilon_j = 1$, we take a copy of $[0, 1]$ and direct it from 0 to 1.

In our constructions, every interval will have two directions: the *intrinsic direction* (which is the direction from 0 to 1) and the *assigned direction*. In the case just discussed, these agree, but in general they will not.

We write $[0, 1]_{f,j,w}$ for such an interval and $\mathcal{I}_{f,w}^+$ for the collection of these intervals.

For every $j \in [w]$ with $f_j = f$ as in (3-13) and $\epsilon_j = -1$, we take a copy of $[0, 1]$ and direct this interval from 1 to 0. We write $[0, 1]_{f^{-1},j,w}$ for such an interval and $\mathcal{I}_{f,w}^-$ for the collection of these intervals.

All the intervals described above are called *w-intervals*. There are $|w|$ of these intervals in total.

w-intermediate-intervals Between each $[0, 1]_{f_j^{\epsilon_j},j,w}$ and $[0, 1]_{f_{j+1}^{\epsilon_{j+1}},j+1,w}$ we add a new interval connecting $1_{f_j^{\epsilon_j},j,w}$ to $0_{f_{j+1}^{\epsilon_{j+1}},j+1,w}$, where the indices j run mod $|w|$. These intervals added are called *w-intermediate-intervals*. Note that these intervals together with the *w-intervals* now form a closed cycle that is paved by $2|w|$ intervals alternating between *w-intervals* and *w-intermediate-intervals*. Starting at $[0, 1]_{f_1^{\epsilon_1},1,w}$, reading the directions and *f*-labels of the *w-intervals* so that every *w-interval* is traversed from 0 to 1 spells out the word *w*. The resulting circle is called the *w-loop* and the previously defined orientation of this loop is now fixed. See Figure 1 for an illustration of the *w-loop* in a particular example.

We now view the indices i_j as an assignment

$$\mathbf{a}: \{\text{endpoints of } w\text{-intervals}\} \rightarrow [n],$$

$$\mathbf{a}(0_{f,j,w}) := i_j, \quad \mathbf{a}(1_{f,j,w}) = i_{j+1}, \quad \mathbf{a}(0_{f^{-1},j,w}) = i_j, \quad \mathbf{a}(1_{f^{-1},j,w}) = i_{j+1}.$$

The condition that \mathbf{a} comes from a single collection of i_j is precisely that *if two endpoints of w-intervals are connected by a w-intermediate-interval, they are assigned the same value by a*. Let $\mathcal{A}(w)$ denote the collection of such \mathbf{a} . If I is any copy of $[0, 1]$, we write 0_I for the copy of 0 and 1_I for the copy of 1 in I . We can now write

$$\text{Tr}(w(h)) = \sum_{\mathbf{a} \in \mathcal{A}(w)} \prod_{f \in \{a,b,c,d\}} \left(\prod_{i \in \mathcal{I}_{f,w}^+} h_{\mathbf{a}(0_i)\mathbf{a}(1_i)} \right) \left(\prod_{j \in \mathcal{I}_{f,w}^-} \bar{h}_{\mathbf{a}(1_j)\mathbf{a}(0_j)} \right).$$

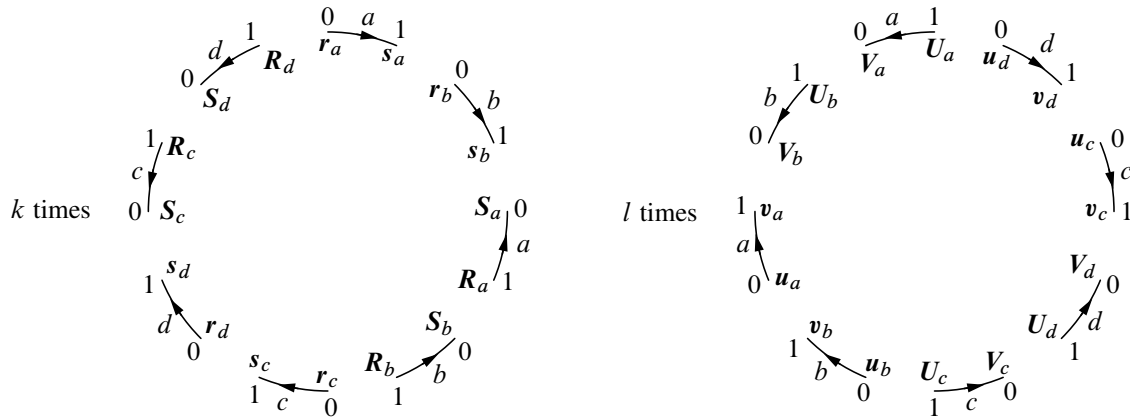


Figure 2: The R -intervals (left) and the R^{-1} -intervals (right). We have indicated their assigned direction and label (which f they correspond to). We have also, for each endpoint of an interval, indicated which index function, eg r_a , has this endpoint in its domain.

For each $j \in [k]$ and $f \in \{a, b, c, d\}$, we make a copy of $[0, 1]$, direct it from 1 to 0, label it by f , and also number it by j . We write $\mathcal{J}_{f,R}^-$ for the collection of these intervals. These correspond to occurrences of f^{-1} in R .

(These two constructions of k intervals correspond to the presence of f and f^{-1} each exactly once in R .)

These intervals are called R -intervals. There are $8k$ R -intervals in total (for general g , there are $4gk$ of these intervals).

R^{-1} -intervals For each $j \in [k + 1, k + l]$ and $f \in \{a, b, c, d\}$, we make a copy of $[0, 1]$, direct it from 0 to 1, label it by f , and also number it by j . We write $\mathcal{J}_{f,R^{-1}}^+$ for the collection of these intervals. These correspond to occurrences of f in R^{-1} .

For each $j \in [k + 1, k + l]$ and $f \in \{a, b, c, d\}$, we make a copy of $[0, 1]$, direct it from 1 to 0, label it by f , and also number it by j . We write $\mathcal{J}_{f,R^{-1}}^-$ for the collection of these intervals. These correspond to occurrences of f^{-1} in R^{-1} .

These intervals are called R^{-1} -intervals. There are $8l$ R^{-1} -intervals in total (for general g , there are $4gl$ of these intervals). See Figure 2 for an illustration of the R - and R^{-1} -intervals.

We now view (by identifying endpoints of intervals with the given numbers of intervals in $[k + l]$)

$$\begin{aligned}
 r_f &: \{0_i : i \in \mathcal{J}_{f,R}^+\} \rightarrow [n], & R_f &: \{1_i : i \in \mathcal{J}_{f,R}^-\} \rightarrow [n], \\
 s_f &: \{1_i : i \in \mathcal{J}_{f,R}^+\} \rightarrow [n], & S_f &: \{0_i : i \in \mathcal{J}_{f,R}^-\} \rightarrow [n], \\
 U_f &: \{1_i : i \in \mathcal{J}_{f,R^{-1}}^-\} \rightarrow [n], & u_f &: \{0_i : i \in \mathcal{J}_{f,R^{-1}}^+\} \rightarrow [n], \\
 V_f &: \{0_i : i \in \mathcal{J}_{f,R^{-1}}^-\} \rightarrow [n], & v_f &: \{1_i : i \in \mathcal{J}_{f,R^{-1}}^+\} \rightarrow [n].
 \end{aligned}$$

We obtain, from (3-16),

$$\begin{aligned} &\langle h_a e_{s_a}^{V_a}, e_{r_a}^{U_a} \rangle \langle h_b e_{s_b}^{V_b}, e_{r_b}^{U_b} \rangle \langle h_a^{-1} e_{R_a}^{u_a}, e_{S_a}^{v_a} \rangle \langle h_b^{-1} e_{R_b}^{u_b}, e_{S_b}^{v_b} \rangle \\ &\quad \cdot \langle h_c e_{s_c}^{V_c}, e_{r_c}^{U_c} \rangle \langle h_d e_{s_d}^{V_d}, e_{r_d}^{U_d} \rangle \langle h_c^{-1} e_{R_c}^{u_c}, e_{S_c}^{v_c} \rangle \langle h_d^{-1} e_{R_d}^{u_d}, e_{S_d}^{v_d} \rangle \\ &= \prod_f \prod_{\substack{i^+ \in \mathfrak{J}_{f,R}^+ \\ i^- \in \mathfrak{J}_{f,R}^-}} \prod_{\substack{j^+ \in \mathfrak{J}_{f,R-1}^+ \\ j^- \in \mathfrak{J}_{f,R-1}^-}} h_{r_f(0_{i^+})} s_f(1_{i^+}) h_{u_f(0_{j^+})} v_f(1_{j^+}) \bar{h}_{R_f(1_{i^-})} S_f(0_{i^-}) \bar{h}_{U_f(1_{j^-})} V_f(0_{j^-}). \end{aligned}$$

With this formalism, we obtain

$$\begin{aligned} (3-17) \quad \mathcal{F}_n(w, \mu, \nu) &= \sum_{\substack{P_i \\ \mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f \\ \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f, \mathbf{S}_f}} \sum_{\mathbf{a} \in \mathfrak{sl}(w)} \beta_{p_2 s_a}^{V_a} \bar{\beta}_{p_1 r_a}^{U_a} \beta_{p_3 s_b}^{V_b} \bar{\beta}_{p_2 r_b}^{U_b} \beta_{p_4 R_a}^{u_a} \bar{\beta}_{p_3 S_a}^{v_a} \beta_{p_5 R_b}^{u_b} \bar{\beta}_{p_4 S_b}^{v_b} \beta_{p_6 s_c}^{V_c} \\ &\quad \cdot \bar{\beta}_{p_5 r_c}^{U_c} \beta_{p_7 s_d}^{V_d} \bar{\beta}_{p_6 r_d}^{U_d} \beta_{p_8 R_c}^{u_c} \bar{\beta}_{p_7 S_c}^{v_c} \beta_{p_1 R_d}^{u_d} \bar{\beta}_{p_8 S_d}^{v_d} \\ &\quad \cdot \prod_{f \in \{a, b, c, d\}} \int_{h \in \mathbb{U}(n)} \prod_{\substack{i \in \mathfrak{J}_{f,w}^+, j \in \mathfrak{J}_{f,w}^- \\ i^+ \in \mathfrak{J}_{f,R}^+, i^- \in \mathfrak{J}_{f,R}^- \\ j^+ \in \mathfrak{J}_{f,R-1}^+, j^- \in \mathfrak{J}_{f,R-1}^-}} h_{r_f(0_{i^+})} s_f(1_{i^+}) h_{u_f(0_{j^+})} v_f(1_{j^+}) \bar{h}_{R_f(1_{i^-})} S_f(0_{i^-}) \\ &\quad \cdot \bar{h}_{U_f(1_{j^-})} V_f(0_{j^-}) dh. \end{aligned}$$

For each f , the integral in (3-17) can be done using the Weingarten calculus (Theorem 2.5). To do this, fix bijections for each $f \in \{a, b, c, d\}$

$$\begin{aligned} \mathfrak{J}_f^+ &:= \mathfrak{J}_{f,R}^+ \cup \mathfrak{J}_{f,R-1}^+ \cup \mathfrak{J}_{f,w}^+ \cong [k + l + p_f], \\ \mathfrak{J}_f^- &:= \mathfrak{J}_{f,R}^- \cup \mathfrak{J}_{f,R-1}^- \cup \mathfrak{J}_{f,w}^- \cong [k + l + p_f] \end{aligned}$$

such that

$$\mathfrak{J}_{f,w}^+ \cong [k + l + 1, k + l + p_f], \quad \mathfrak{J}_{f,w}^- \cong [k + l + 1, k + l + p_f]$$

and

$$(3-18) \quad \mathfrak{J}_{f,R}^+ \cong [k], \quad \mathfrak{J}_{f,R}^- \cong [k], \quad \mathfrak{J}_{f,R-1}^+ \cong [k + 1, k + l], \quad \mathfrak{J}_{f,R-1}^- \cong [k + 1, k + l]$$

correspond to the original numberings of $\mathfrak{J}_{f,R}^+$, $\mathfrak{J}_{f,R}^-$, $\mathfrak{J}_{f,R-1}^+$ and $\mathfrak{J}_{f,R-1}^-$.

Hence, if $\sigma_f, \tau_f \in S_{k+l+p_f}$ we view $\sigma_f, \tau_f: \mathfrak{J}_f^+ \rightarrow \mathfrak{J}_f^-$ by the above fixed bijections. For each $f \in \{a, b, c, d\}$, we say $(\mathbf{a}, \mathbf{r}_f, \mathbf{u}_f, \mathbf{R}_f, \mathbf{U}_f) \rightarrow \sigma_f$ if, for all $i \in \mathfrak{J}_f^+$ and $i' \in \mathfrak{J}_f^-$ with $\sigma_f(i) = i'$, we have

$$[\mathbf{r}_f \sqcup \mathbf{u}_f \sqcup \mathbf{a}](0_i) = [\mathbf{R}_f \sqcup \mathbf{U}_f \sqcup \mathbf{a}](1_{i'});$$

here we wrote eg $[\mathbf{r}_f \sqcup \mathbf{u}_f \sqcup \mathbf{a}]$ for the function that \mathbf{a}, \mathbf{r}_f and \mathbf{u}_f induce on $\{0_i : i \in \mathfrak{J}_f^+\}$. Similarly, we say $(\mathbf{a}, \mathbf{s}_f, \mathbf{v}_f, \mathbf{S}_f, \mathbf{V}_f) \rightarrow \tau_f$ if, for all $i \in \mathfrak{J}_f^+, i' \in \mathfrak{J}_f^-$ with $\tau_f(i) = i'$, we have

$$[\mathbf{s}_f \sqcup \mathbf{v}_f \sqcup \mathbf{a}](1_i) = [\mathbf{S}_f \sqcup \mathbf{V}_f \sqcup \mathbf{a}](0_{i'}).$$

Theorem 2.5 translates to

$$\int_{h \in U(n)} \prod_{\substack{i \in \mathcal{J}_{f,w}^+, j \in \mathcal{J}_{f,w}^- \\ i^+ \in \mathcal{J}_{f,R}^+, i^- \in \mathcal{J}_{f,R}^- \\ j^+ \in \mathcal{J}_{f,R-1}^+, j^- \in \mathcal{J}_{f,R-1}^-}} h_{\mathbf{r}_f(0_{i^+})s_f(1_{i^+})} h_{\mathbf{u}_f(0_{j^+})\mathbf{v}_f(1_{j^+})} \bar{h}_{\mathbf{R}_f(1_{i^-})\mathbf{S}_f(0_{i^-})} \bar{h}_{\mathbf{U}_f(1_{j^-})\mathbf{V}_f(0_{j^-})} dh$$

$$= \sum_{\sigma_f, \tau_f \in S_{k+l+p_f}} W_{\mathbf{g}_{n,k+l+p_f}}(\sigma_f \tau_f^{-1}) \mathbb{1}\{(\mathbf{a}, \mathbf{r}_f, \mathbf{u}_f, \mathbf{R}_f, \mathbf{U}_f) \rightarrow \sigma_f, (\mathbf{a}, \mathbf{s}_f, \mathbf{v}_f, \mathbf{S}_f, \mathbf{V}_f) \rightarrow \tau_f\},$$

so putting this into (3-17) gives

$$\mathcal{J}_n(w, \mu, \nu) = \sum_{\sigma_f, \tau_f \in S_{k+l+p_f}} \left(\prod_{f \in \{a,b,c,d\}} W_{\mathbf{g}_{n,k+l+p_f}}(\sigma_f \tau_f^{-1}) \right)$$

$$\cdot \sum_{\substack{p_i \mathbf{a} \in \mathfrak{sl}(w), \mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f, \mathbf{S}_f \\ (\mathbf{a}, \mathbf{r}_f, \mathbf{u}_f, \mathbf{R}_f, \mathbf{U}_f) \rightarrow \sigma_f \\ (\mathbf{a}, \mathbf{s}_f, \mathbf{v}_f, \mathbf{S}_f, \mathbf{V}_f) \rightarrow \tau_f}} \beta_{p_2 s_a}^{V_a} \bar{\beta}_{p_1 r_a}^{U_a} \beta_{p_3 s_b}^{V_b} \bar{\beta}_{p_2 r_b}^{U_b} \beta_{p_4 R_a}^{u_a} \bar{\beta}_{p_3 S_a}^{v_a} \beta_{p_5 R_b}^{u_b} \bar{\beta}_{p_4 S_b}^{v_b}$$

$$\cdot \beta_{p_6 s_c}^{V_c} \bar{\beta}_{p_5 r_c}^{U_c} \beta_{p_7 s_d}^{V_d} \bar{\beta}_{p_6 r_d}^{U_d} \beta_{p_8 R_c}^{u_c} \bar{\beta}_{p_7 S_c}^{v_c} \beta_{p_1 R_d}^{u_d} \bar{\beta}_{p_8 S_d}^{v_d}.$$

Here we make our main improvement over the classical Weingarten calculus. We introduce the following beneficial property that the σ_f and τ_f possibly have:

Forbidden matchings property For every $f \in \{a, b, c, d\}$, the following hold: neither σ_f nor τ_f map any element of $\mathcal{J}_{f,R}^+$ to an element of $\mathcal{J}_{f,R-1}^-$, or map an element of $\mathcal{J}_{f,R-1}^+$ to an element of $\mathcal{J}_{f,R}^-$.

We have the following key lemma:

Lemma 3.4 If for some $f \in \{a, b, c, d\}$, σ_f and τ_f do **not** have the **forbidden matchings** property, then, for any choice of p_1, \dots, p_8 ,

$$(3-19) \quad \sum_{\substack{\mathbf{a} \in \mathfrak{sl}(w), \mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f, \mathbf{S}_f \\ (\mathbf{a}, \mathbf{r}_f, \mathbf{u}_f, \mathbf{R}_f, \mathbf{U}_f) \rightarrow \sigma_f \\ (\mathbf{a}, \mathbf{s}_f, \mathbf{v}_f, \mathbf{S}_f, \mathbf{V}_f) \rightarrow \tau_f}} \beta_{p_2 s_a}^{V_a} \bar{\beta}_{p_1 r_a}^{U_a} \beta_{p_3 s_b}^{V_b} \bar{\beta}_{p_2 r_b}^{U_b} \beta_{p_4 R_a}^{u_a} \bar{\beta}_{p_3 S_a}^{v_a} \beta_{p_5 R_b}^{u_b} \bar{\beta}_{p_4 S_b}^{v_b}$$

$$\cdot \beta_{p_6 s_c}^{V_c} \bar{\beta}_{p_5 r_c}^{U_c} \beta_{p_7 s_d}^{V_d} \bar{\beta}_{p_6 r_d}^{U_d} \beta_{p_8 R_c}^{u_c} \bar{\beta}_{p_7 S_c}^{v_c} \beta_{p_1 R_d}^{u_d} \bar{\beta}_{p_8 S_d}^{v_d} = 0.$$

Proof Indeed, suppose σ_a matches an element $i \in \mathcal{J}_{a,R}^+$ with $j \in \mathcal{J}_{a,R-1}^-$; $\sigma_a(i) = j$. With our given fixed bijections (3-18), i corresponds to an element of $[k]$ and j corresponds to an element of $[k + 1, k + l]$. Without loss of generality in the argument suppose that 0_i corresponds to 1 and 0_j corresponds to $k + 1$. The condition $\sigma_a(i) = j$ and $(\mathbf{a}, \mathbf{r}_a, \mathbf{u}_a, \mathbf{R}_a, \mathbf{U}_a) \rightarrow \sigma_f$ means that, as functions on $[k]$ and $[k + 1, k + l]$, $\mathbf{r}_a(1) = \mathbf{U}_a(k + 1)$. There are no other constraints on these values.

Then, for all variables in (3-19) fixed apart from \mathbf{r}_a and \mathbf{U}_a , and all values of \mathbf{r}_a and \mathbf{U}_a fixed other than $\mathbf{r}_a(1)$ and $\mathbf{U}_a(k + 1)$, the ensuing sum over \mathbf{r}_a and \mathbf{U}_a is

$$\sum_{\mathbf{r}_a(1)=\mathbf{U}_a(k+1)} \beta_{p_2 r_a}^{U_a}.$$

But, recalling the contraction operators from (2-2), this sum is the coordinate of

$$e_{\mathbf{r}_a(2)} \otimes \cdots \otimes e_{\mathbf{r}_a(k)} \otimes \check{e}_{U_a(k+2)} \otimes \cdots \otimes \check{e}_{U_a(k+l)}$$

in $c_{1,1}(v_{p_2})$. But $c_{1,1}(v_{p_2}) = 0$ because $v_{p_2} \in \dot{J}_n^{k,l}$. □

We henceforth write $\sum_{\sigma_f, \tau_f}^*$ to mean the sum is restricted to σ_f and τ_f satisfying the **forbidden matchings property**. Lemma 3.4 now implies

$$\begin{aligned} (3-20) \quad \mathcal{J}_n(w, \mu, v) &= \sum_{\sigma_f, \tau_f \in S_{k+l+p_f}}^* \left(\prod_{f \in \{a,b,c,d\}} W_{\mathfrak{g}_n, k+l+p_f}(\sigma_f \tau_f^{-1}) \right) \\ &\cdot \sum_{\substack{p_i \quad \mathbf{a} \in \mathfrak{sl}(w), \mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f, \mathbf{S}_f \\ (\mathbf{a}, \mathbf{r}_f, \mathbf{u}_f, \mathbf{R}_f, \mathbf{U}_f) \rightarrow \sigma_f \\ (\mathbf{a}, \mathbf{s}_f, \mathbf{v}_f, \mathbf{S}_f, \mathbf{V}_f) \rightarrow \tau_f}} \beta_{p_2 s_a}^{V_a} \bar{\beta}_{p_1 r_a}^{U_a} \beta_{p_3 s_b}^{V_b} \bar{\beta}_{p_2 r_b}^{U_b} \beta_{p_4 R_a}^{u_a} \bar{\beta}_{p_3 S_a}^{v_a} \beta_{p_5 R_b}^{u_b} \bar{\beta}_{p_4 S_b}^{v_b} \\ &\cdot \beta_{p_6 s_c}^{V_c} \bar{\beta}_{p_5 r_c}^{U_c} \beta_{p_7 s_d}^{V_d} \bar{\beta}_{p_6 r_d}^{U_d} \beta_{p_8 R_c}^{u_c} \bar{\beta}_{p_7 S_c}^{v_c} \beta_{p_1 R_d}^{u_d} \bar{\beta}_{p_8 S_d}^{v_d}. \end{aligned}$$

Moreover, we can significantly tidy up (3-20). For everything in (3-20) fixed except for eg p_2 , the ensuing sum over p_2 is

$$\sum_{p_2} \beta_{p_2 s_a}^{V_a} \bar{\beta}_{p_2 r_b}^{U_b} = \sum_{p_2} \langle e_{\mathbf{r}_b}^{U_b}, v_{p_2} \rangle \langle v_{p_2}, e_{s_a}^{V_a} \rangle = \langle q_\theta e_{\mathbf{r}_b}^{U_b}, e_{s_a}^{V_a} \rangle.$$

Therefore, executing the sums over p_i in (3-20), we replace the sum over p_i and the product over β -terms by

$$(3-21) \quad \langle q_\theta e_{\mathbf{r}_b}^{U_b}, e_{s_a}^{V_a} \rangle \langle q_\theta e_{s_a}^{v_a}, e_{s_b}^{V_b} \rangle \langle q_\theta e_{s_b}^{v_b}, e_{\mathbf{R}_a}^{u_a} \rangle \langle q_\theta e_{\mathbf{r}_c}^{U_c}, e_{\mathbf{R}_b}^{u_b} \rangle \langle q_\theta e_{\mathbf{r}_d}^{U_d}, e_{s_c}^{V_c} \rangle \langle q_\theta e_{s_c}^{v_c}, e_{s_d}^{V_d} \rangle \langle q_\theta e_{s_d}^{v_d}, e_{\mathbf{R}_c}^{u_c} \rangle \cdot \langle q_\theta e_{\mathbf{r}_a}^{U_a}, e_{\mathbf{R}_d}^{u_d} \rangle.$$

By Proposition 3.2, we have eg

$$\langle q_\theta e_{\mathbf{r}_b}^{U_b}, e_{s_a}^{V_a} \rangle = D_{\mu, v}(n) \sum_{\pi \in S_{k+l}} z_\theta(\pi) \langle \varphi^{-1}[\rho_n^{k,l}(\pi)] e_{\mathbf{r}_b}^{U_b}, e_{s_a}^{V_a} \rangle.$$

Now recall from (3-12) that

$$\varphi^{-1}[\rho_n^{k+l}(\pi)] = \sum_{\substack{I=(i_1, \dots, i_k) \\ J=(j_{k+1}, \dots, j_{k+l})}} e_{I'(I \sqcup J; \pi)}^J \otimes \check{e}_I^{J'(I \sqcup J; \pi)}.$$

This means that $\langle \varphi^{-1}[\rho_n^{k+l}(\pi)] e_{\mathbf{r}_b}^{U_b}, e_{s_a}^{V_a} \rangle$ is equal to either 0 or 1 and $\langle \varphi^{-1}[\rho_n^{k+l}(\pi)] e_{\mathbf{r}_b}^{U_b}, e_{s_a}^{V_a} \rangle = 1$ if and only if, letting (3-18) induce identifications

$$\begin{aligned} \{1_i : i \in \mathcal{I}_{a, R}^+\} &\cong [k], & \{1_i : i \in \mathcal{I}_{b, R-1}^-\} &\cong [k+1, k+l], \\ \{0_i : i \in \mathcal{I}_{b, R}^+\} &\cong [k], & \{0_i : i \in \mathcal{I}_{a, R-1}^-\} &\cong [k+1, k+l] \end{aligned}$$

via their given indexing of intervals, we have $[s_a \sqcup U_b] \circ \pi = [r_b \sqcup V_a]$, where eg $s_a \sqcup U_b$ is the function either on endpoints of intervals or on $[k + l]$ induced by the union of s_a and U_b . Hence, repeating this argument,

$$(3-21) = D_{\mu, \nu}(n)^8 \sum_{\pi_1, \dots, \pi_8 \in S_{k+l}} \left(\prod_{i=1}^8 z_\theta(\pi_i) \right) \mathbb{1} \{ [s_a \sqcup U_b] \circ \pi_1 = [r_b \sqcup V_a], [s_b \sqcup v_a] \circ \pi_2 = [S_a \sqcup V_b], \\ [R_a \sqcup v_b] \circ \pi_3 = [S_b \sqcup u_a], [R_b \sqcup U_c] \circ \pi_4 = [r_c \sqcup u_b], \\ [s_c \sqcup U_d] \circ \pi_5 = [r_d \sqcup V_c], [s_d \sqcup v_c] \circ \pi_6 = [S_c \sqcup V_d], \\ [R_c \sqcup v_d] \circ \pi_7 = [S_d \sqcup u_c], [R_d \sqcup U_a] \circ \pi_8 = [r_a \sqcup u_d] \}.$$

Putting all these arguments together gives

$$\mathcal{J}_n(w, \mu, \nu) = D_{\mu, \nu}(n)^8 \sum_{\sigma_f, \tau_f \in S_{p_f+k+l}}^* \sum_{\pi_1, \dots, \pi_8 \in S_{k+l}} \left(\prod_{f \in \{a, b, c, d\}} W_{g_{n, k+l+p_f}}(\sigma_f \tau_f^{-1}) \right) \left(\prod_{i=1}^8 z_\theta(\pi_i) \right) \\ \cdot \sum_{P_i} \sum_{\substack{\mathbf{a} \in \mathcal{A}(w), \mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f, \mathbf{S}_f \\ (\mathbf{a}, \mathbf{r}_f, \mathbf{u}_f, \mathbf{R}_f, \mathbf{U}_f) \rightarrow \sigma_f \\ (\mathbf{a}, \mathbf{s}_f, \mathbf{v}_f, \mathbf{S}_f, \mathbf{V}_f) \rightarrow \tau_f}} \mathbb{1} \{ [s_a \sqcup U_b] \circ \pi_1 = [r_b \sqcup V_a], [s_b \sqcup v_a] \circ \pi_2 = [S_a \sqcup V_b], \\ [R_a \sqcup v_b] \circ \pi_3 = [S_b \sqcup u_a], [R_b \sqcup U_c] \circ \pi_4 = [r_c \sqcup u_b], \\ [s_c \sqcup U_d] \circ \pi_5 = [r_d \sqcup V_c], [s_d \sqcup v_c] \circ \pi_6 = [S_c \sqcup V_d], \\ [R_c \sqcup v_d] \circ \pi_7 = [S_d \sqcup u_c], [R_d \sqcup U_a] \circ \pi_8 = [r_a \sqcup u_d] \}.$$

This formula says that we can calculate $\mathcal{J}_n(w, \mu, \nu)$ by summing over some combinatorial data of matchings (the σ_f, τ_f and π_i) a quantity that we can understand well times a count of the number of indices that satisfy the prescribed matchings. To formalize this point of view we make the following definition:

Definition 3.5 A matching datum of the triple (w, k, l) is a pair $(\sigma_f, \tau_f) \in S_{k+l+p_f} \times S_{k+l+p_f}$ as above, satisfying the forbidden matchings property for each $f \in \{a, b, c, d\}$, together with $(\pi_1, \dots, \pi_8) \in (S_{k+l})^8$. We write

$$\text{MATCH}(w, k, l)$$

for the finite collection of all matching data for (w, k, l) .

Given a matching datum $\{\sigma_f, \tau_f, \pi_i\}$, we write $\mathcal{N}(\{\sigma_f, \tau_f, \pi_i\})$ for the number of choices of $\mathbf{a} \in \mathcal{A}(w)$, $\mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f$ and \mathbf{S}_f such that

$$(3-22) \quad \begin{aligned} (\mathbf{a}, \mathbf{r}_f, \mathbf{u}_f, \mathbf{R}_f, \mathbf{U}_f) &\rightarrow \sigma_f, & (\mathbf{a}, \mathbf{s}_f, \mathbf{v}_f, \mathbf{S}_f, \mathbf{V}_f) &\rightarrow \tau_f, \\ [s_a \sqcup U_b] \circ \pi_1 &= [r_b \sqcup V_a], & [s_b \sqcup v_a] \circ \pi_2 &= [S_a \sqcup V_b], \\ [R_a \sqcup v_b] \circ \pi_3 &= [S_b \sqcup u_a], & [R_b \sqcup U_c] \circ \pi_4 &= [r_c \sqcup u_b], \\ [s_c \sqcup U_d] \circ \pi_5 &= [r_d \sqcup V_c], & [s_d \sqcup v_c] \circ \pi_6 &= [S_c \sqcup V_d], \\ [R_c \sqcup v_d] \circ \pi_7 &= [S_d \sqcup u_c], & [R_d \sqcup U_a] \circ \pi_8 &= [r_a \sqcup u_d]. \end{aligned}$$

With this notation, we have proved the following theorem:

Theorem 3.6 For $k + l > 0$ with $\mu \vdash k$ and $\nu \vdash l$ and $w \in [F_4, F_4]$, we have

$$(3-23) \quad \mathcal{F}_n(w, \mu, \nu) = D_{\mu, \nu}(n)^8 \sum_{\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)} \left(\prod_{i=1}^8 z_{\theta}(\pi_i) \right) \left(\prod_{f \in \{a, b, c, d\}} \text{Wg}_{n, k+l+p_f}(\sigma_f \tau_f^{-1}) \right) \cdot \mathcal{N}(\{\sigma_f, \tau_f, \pi_i\}).$$

We conclude this section by bounding the terms $z_{\theta}(\pi_i)$ and $\text{Wg}_{n, k+l+p_f}(\sigma_f \tau_f^{-1})$ using Proposition 2.7 and Lemma 3.3, recalling also (2-1). Note that $\sum_{f \in \{a, b, c, d\}} p_f = \frac{1}{2}|w|$. This yields:

Corollary 3.7 For $k + l > 0$ with $\mu \vdash k$ and $\nu \vdash l$ and $w \in [F_4, F_4]$, we have

$$(3-24) \quad \mathcal{F}_n(w, \mu, \nu) \ll_{k, l, w} n^{-4k-4l-|w|/2} \sum_{\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)} n^{-\sum_f |\sigma_f \tau_f^{-1}| - \sum_{i=1}^8 \|\pi_i\|_{k, l}} \mathcal{N}(\{\sigma_f, \tau_f, \pi_i\}).$$

We will proceed in the next section to understand all the quantities in (3-24) in topological terms by constructing a surface from each $\{\sigma_f, \tau_f, \pi_i\}$.

4 Topology

4.1 Construction of surfaces from matching data

We now show how a datum in $\text{MATCH}(w, k, l)$ can be used to construct a surface such that the terms appearing in (3-23) can be bounded by topological features of the surface. This construction is similar to the constructions of [Magee and Puder 2019; 2015], but with the presence of additional π_i adding a new aspect. We continue to assume $g = 2$ for simplicity. We can still assume that $\gamma \in [\Gamma_2, \Gamma_2]$ and hence $w \in [F_4, F_4]$.

Construction of the 1-skeleton

π -intervals The identifications of the previous section mean that we view

$$(4-1) \quad \begin{aligned} \pi_1 : \{0_i : i \in \mathcal{J}_{b,R}^+ \cup \mathcal{J}_{a,R-1}^-\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{a,R}^+ \cup \mathcal{J}_{b,R-1}^-\}, \\ \pi_2 : \{0_i : i \in \mathcal{J}_{a,R}^- \cup \mathcal{J}_{b,R-1}^+\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{b,R}^+ \cup \mathcal{J}_{a,R-1}^+\}, \\ \pi_3 : \{0_i : i \in \mathcal{J}_{b,R}^- \cup \mathcal{J}_{a,R-1}^+\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{a,R}^- \cup \mathcal{J}_{b,R-1}^+\}, \\ \pi_4 : \{0_i : i \in \mathcal{J}_{c,R}^+ \cup \mathcal{J}_{b,R-1}^+\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{b,R}^- \cup \mathcal{J}_{c,R-1}^-\}, \\ \pi_5 : \{0_i : i \in \mathcal{J}_{d,R}^+ \cup \mathcal{J}_{c,R-1}^-\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{c,R}^+ \cup \mathcal{J}_{d,R-1}^-\}, \\ \pi_6 : \{0_i : i \in \mathcal{J}_{c,R}^- \cup \mathcal{J}_{d,R-1}^-\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{d,R}^+ \cup \mathcal{J}_{c,R-1}^+\}, \\ \pi_7 : \{0_i : i \in \mathcal{J}_{d,R}^- \cup \mathcal{J}_{c,R-1}^+\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{c,R}^- \cup \mathcal{J}_{d,R-1}^+\}, \\ \pi_8 : \{0_i : i \in \mathcal{J}_{a,R}^+ \cup \mathcal{J}_{d,R-1}^+\} &\rightarrow \{1_{i'} : i' \in \mathcal{J}_{d,R}^- \cup \mathcal{J}_{a,R-1}^-\}. \end{aligned}$$

We add an arc between any two interval endpoints that are mapped to one another by some π_i . All the intervals added here are called π -intervals. The purpose of this construction is that the conditions concerning π_i in (3-22) correspond to the fact that *two endpoints of intervals connected by a π -interval are assigned the same value in $[n]$ by the relevant functions out of $\mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f$ and \mathbf{S}_f (at most one of these functions has any given interval endpoint in its domain).*

The π -intervals together with the R -intervals and R^{-1} -intervals form a collection of loops, which we call R^\pm - π -loops.

σ -arcs and τ -arcs Recall from the previous sections that we view

$$\sigma_f, \tau_f: \mathcal{I}_f^+ \rightarrow \mathcal{I}_f^-.$$

We add an arc between each 0_i and $1_{i'}$ with $\sigma_f(i) = i'$ and between each 1_i and $0_{i'}$ with $\tau_f(i) = i'$. These arcs are called σ_f -arcs and τ_f -arcs, respectively. Any σ_f -arc (resp. τ_f -arc) is also called a σ -arc (resp. τ -arc). Notice even though an arc is formally the same as an interval, we distinguish these types of objects. The only arcs that exist are σ -arcs and τ -arcs. The purpose of this construction is that the conditions pertaining to σ_f and τ_f in (3-22) are equivalent to the fact that *two endpoints of intervals connected by a σ -arc or τ -arc are assigned the same value in $[n]$ by the relevant functions out of $\mathbf{a}, \mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f$ and \mathbf{S}_f .*

After adding these arcs, every endpoint of an interval has exactly one arc emanating from it. We have therefore now constructed a trivalent graph

$$G(\{\sigma_f, \tau_f, \pi_i\}).$$

Each vertex of the graph is an endpoint of two intervals and one arc. The number of vertices of this graph is twice the total number of w -intervals, R -intervals and R^{-1} -intervals, which is $2(|w| + 8(k + l))$. Therefore we have

$$(4-2) \quad \chi(G(\{\sigma_f, \tau_f, \pi_i\})) = -(|w| + 8(k + l)).$$

(For general g , we have $\chi(G(\{\sigma_f, \tau_f, \pi_i\})) = -(|w| + 4g(k + l))$.) Moreover, *the conditions in (3-22) are now interpreted purely in terms of the combinatorics of this graph.*

Gluing in discs There are two types of cycles in $G(\{\sigma_f, \tau_f, \pi_i\})$ that we wish to consider:

- Cycles that alternate between following either a w -intermediate-interval or a π -interval and then either a σ -arc or a τ -arc. These cycles are disjoint from one another, and every σ - or τ -arc is contained in exactly one such cycle. We call these cycles *type I cycles*. For every type I cycle, we glue a disc to $G(\{\sigma_f, \tau_f, \pi_i\})$ along its boundary, following the cycle. These discs will be called *type I discs*. (These are analogous to the o -discs of [Magee and Puder 2019].)

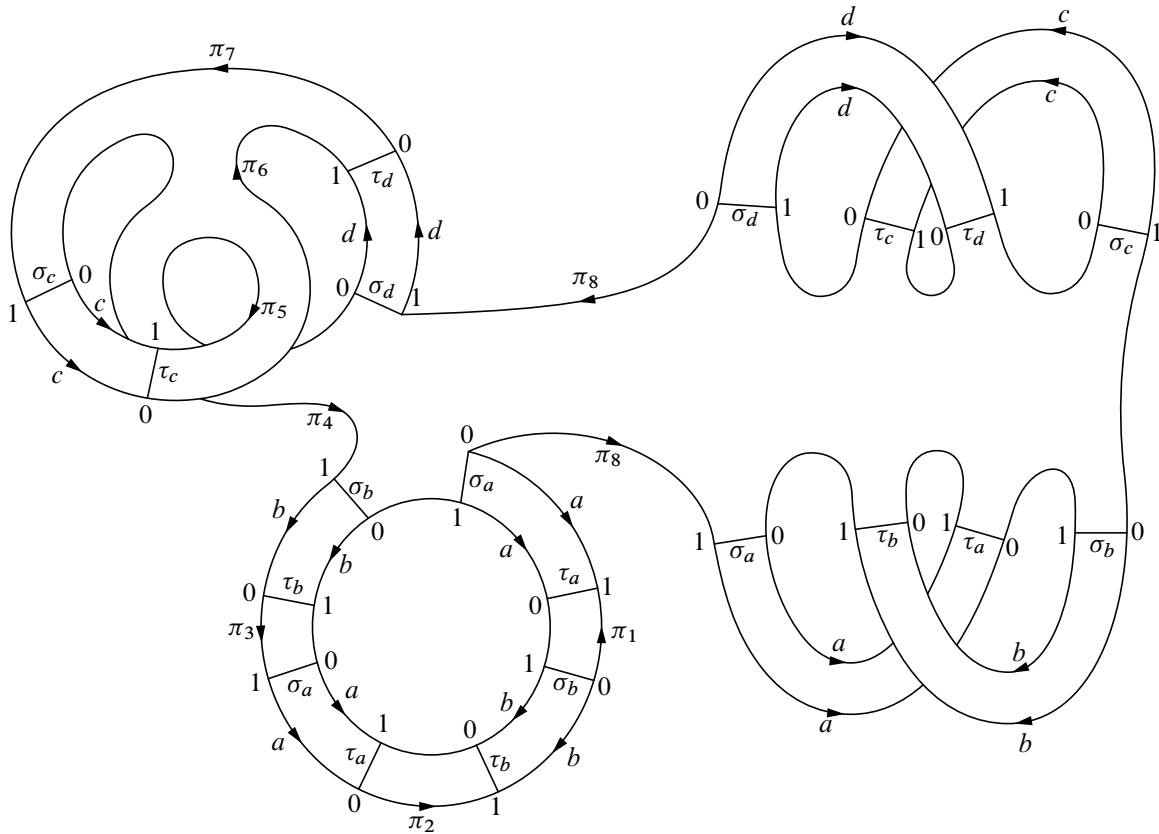


Figure 3: An example $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$ for $w = ab^{-1}a^{-1}b$. The σ , τ and some of the π_i -arcs are labeled along with the numbers (0 or 1) of the points being matched in the w -intervals. Each w -interval is also labeled with its corresponding letter. Here $k = l = 1$; π_8 is a transposition and all other π_i are the identity. There is one resulting R^\pm - π -loop. In this example, for each $f \in \{a, b, c, d\}$, $\sigma_f = \tau_f$. This means that all type II discs are rectangles.

- Cycles that alternate between following either a w -interval, an R -interval or an R^{-1} -interval and then either a σ -arc or a τ -arc. Again, these cycles are disjoint, and every σ - or τ -arc is contained in exactly one such cycle. We call these cycles *type II cycles*. For every type II cycle, we glue a disc to $G(\{\sigma_f, \tau_f, \pi_i\})$ identifying the boundary of the disc with the cycle. These discs will be called *type II discs*. (These are similar to the z -discs of [Magee and Puder 2019].)

Because every interior of an interval meets exactly one of the glued-in discs, and every arc has two boundary segments of discs glued to it, the object resulting from gluing in these discs is a decorated topological surface, which we denote by

$$\Sigma(\{\sigma_f, \tau_f, \pi_i\}).$$

An example of this construction is depicted in Figure 3.

The boundary components of $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$ consist of the w -loop and the R^\pm - π -loops. It is not hard to check that $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$ is orientable with an orientation compatible with the fixed orientations of the boundary loops corresponding to traversing every w -interval or R^\pm -interval from 0 to 1.

We view the given CW-complex structure and the assigned labelings and directions of the intervals that now pave $\partial\Sigma$ as part of the data of $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$. The number of discs of $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$ is connected to the quantities appearing in Theorem 3.6 as follows:

Lemma 4.1
$$\mathcal{N}(\{\sigma_f, \tau_f, \pi_i\}) = n^{\#\{\text{type I discs of } \Sigma(\{\sigma_f, \tau_f, \pi_i\})\}}.$$

Proof The constraints on the functions $\mathbf{a}, \mathbf{r}_f, \mathbf{R}_f, \mathbf{V}_f, \mathbf{v}_f, \mathbf{U}_f, \mathbf{u}_f, \mathbf{s}_f$ and \mathbf{S}_f in (3-22) now correspond to the fact that, altogether, they assign the same value in $[n]$ to every interval endpoint in the same type I cycle, and there are no other constraints between them. \square

The quantities $|\sigma_f \tau_f^{-1}|$ in (3-24) can also be related to $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$ as follows:

Lemma 4.2
$$\prod_{f \in \{a, b, c, d\}} n^{-|\sigma_f \tau_f^{-1}|} = n^{-4(k+l) - |w|/2} n^{\#\{\text{type II discs of } \Sigma(\{\sigma_f, \tau_f, \pi_i\})\}}.$$

Proof Recalling the definition of $|\sigma_f \tau_f^{-1}|$ from Proposition 2.7, we can also write

$$|\sigma_f \tau_f^{-1}| = k + l + p_f - \#\{\text{cycles of } \sigma_f \tau_f^{-1}\}.$$

The cycles of $\{\sigma_f \tau_f^{-1} : f \in \{a, b, c, d\}\}$ are in one-to-one correspondence with the type II cycles of $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$ and hence also the type II discs. Therefore,

$$\begin{aligned} \prod_{f \in \{a, b, c, d\}} n^{-|\sigma_f \tau_f^{-1}|} &= n^{-4(k+l)} n^{\sum_{f \in \{a, b, c, d\}} (-p_f + \#\{\text{cycles of } \sigma_f \tau_f^{-1}\})} \\ &= n^{-4(k+l) - |w|/2} n^{\#\{\text{type II discs of } \Sigma(\{\sigma_f, \tau_f, \pi_i\})\}}. \end{aligned} \quad \square$$

We are now able to prove the following:

Theorem 4.3 For $k + l > 0$ with $\mu \vdash k$ and $\nu \vdash l$ and $w \in [\mathbf{F}_4, \mathbf{F}_4]$, we have

$$\mathcal{I}_n(w, \mu, \nu) \ll_{w, k, l} \sum_{\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)} n^{-\sum_{i=1}^8 \|\pi_i\|_{k, l}} n^{\chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\}))}.$$

Proof Combining Lemmas 4.1 and 4.2 with Corollary 3.7 gives

$$\mathcal{I}_n(w, \mu, \nu) \ll_{w, k, l} n^{-8k - 8l - |w|} \sum_{\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)} n^{-\sum_{i=1}^8 \|\pi_i\|_{k, l}} n^{\#\{\text{discs of } \Sigma(\{\sigma_f, \tau_f, \pi_i\})\}}.$$

Then, from (4-2), we obtain

$$\begin{aligned} \mathcal{I}_n(w, \mu, \nu) &\ll_{w, k, l} \sum_{\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)} n^{-\sum_{i=1}^8 \|\pi_i\|_{k, l}} n^{\chi(G(\{\sigma_f, \tau_f, \pi_i\})) + \#\{\text{discs of } \Sigma(\{\sigma_f, \tau_f, \pi_i\})\}} \\ &= \sum_{\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)} n^{-\sum_{i=1}^8 \|\pi_i\|_{k, l}} n^{\chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\}))}. \end{aligned} \quad \square$$

4.2 Two simplifying surgeries

Theorem 4.3 suggests that we now bound

$$\chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\})) - \sum_{i=1}^8 \|\pi_i\|_{k,l}$$

for all $\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)$. To do this, we make some observations that simplify the task. If C is a simple closed curve in a surface S , then *compressing S along C* means that we cut S along C and then glue discs to cap off any new boundary components created by the cut.

Suppose that we are given $\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)$. Then $\{\sigma_f, \sigma_f, \pi_i\}$ is also in $\text{MATCH}(w, k, l)$ (the *forbidden matching* property continues to hold). It is not hard to see that

$$\chi(\Sigma(\{\sigma_f, \sigma_f, \pi_i\})) \geq \chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\})).$$

Indeed, the τ_f -arcs can be replaced by σ_f -parallel arcs inside the type II discs of $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$. The resulting surface's arcs may not cut the surface into discs, but this can be fixed by (possibly repeatedly) compressing the surface along simple closed curves disjoint from the arcs, leaving the combinatorial data of the arcs unchanged but only potentially increasing the Euler characteristic.

It remains to deal with the sum $\sum_{i=1}^8 \|\pi_i\|_{k,l}$.

Suppose again that an arbitrary $\{\sigma_f, \tau_f, \pi_i\} \in \text{MATCH}(w, k, l)$ is given. For each $i \in [8]$, write

$$\pi_i = \pi_i^* \sigma_i,$$

where $\pi_i^* \in S_k \times S_l$, $\sigma_i = (\pi_i^*)^{-1} \pi_i \in S_{k+l}$ and $|\sigma_i| = \|\pi_i\|_{k,l}$. Let $X_0 := \Sigma(\{\sigma_f, \tau_f, \pi_i\})$.

Take $\Sigma(\{\sigma_f, \tau_f, \pi_i\})$ and add to it all the π_i^* -intervals that would have been added if π_i was replaced by π_i^* for each $i \in [8]$ in its construction. The resulting object X_1 is the decorated surface X_0 together with a collection of π_i^* -intervals with endpoints in the boundary of X_0 , and interiors disjoint from X_0 . This adds $8(k + l)$ edges to X_0 and hence

$$\chi(X_1) = \chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\})) - 8(k + l).$$

Now we consider all cycles that for any fixed $i \in [8]$, alternate between π_i -intervals and π_i^* -intervals. The number of these cycles is the total number of cycles of the permutations $\{(\pi_i^*)^{-1} \pi_i : i \in [8]\}$. On the other hand, the number of cycles of $(\pi_i^*)^{-1} \pi_i$ is

$$k + l - |(\pi_i^*)^{-1} \pi_i| = k + l - |\sigma_i| = k + l - \|\pi_i\|_{k,l}.$$

So in total there are $8(k + l) - \sum_i \|\pi_i\|_{k,l}$ of these cycles. For every such cycle, we glue a disc along its boundary to the cycle. The resulting object is denoted X_2 . Now, X_2 is a topological surface, and we added $8(k + l) - \sum_i \|\pi_i\|_{k,l}$ discs to X_1 to form X_2 , so

$$\chi(X_2) = \chi(X_1) + 8(k + l) - \sum_i \|\pi_i\|_{k,l} = \chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\})) - \sum_i \|\pi_i\|_{k,l}.$$

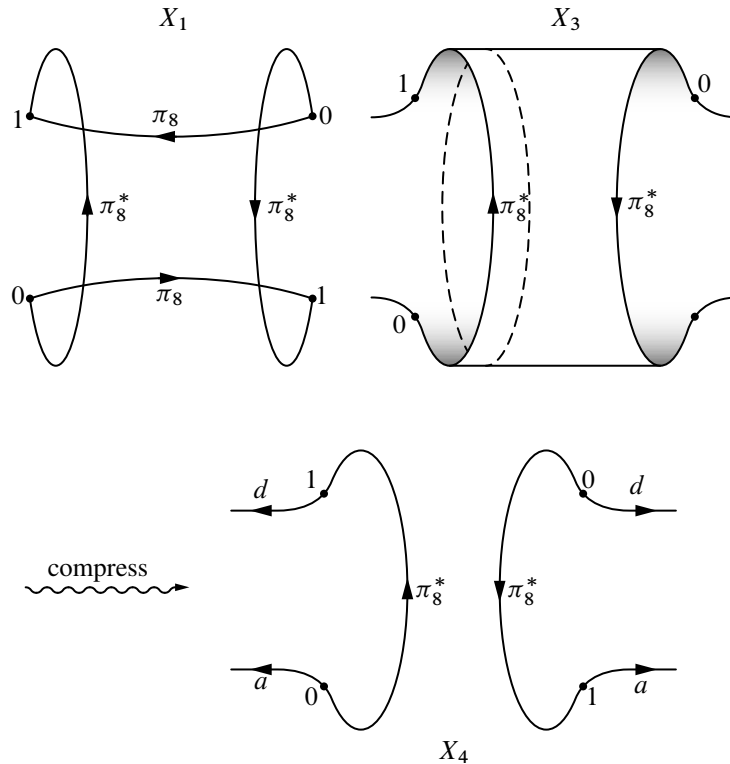


Figure 4: A local illustration of the second type of simplifying surgery, precisely in the context of Figure 3. The dashed simple closed curve in X_3 is disjoint from any arcs, and cutting along this curve and gluing in two discs yields X_4 . Going back to Figure 3 again, the net effect of this surgery is to cut the left half from the right half.

Now “forget” all the original π_i -intervals from X_2 to form X_3 . The surface X_3 is a decorated surface in the same sense as X_0 , except the connected components of $X_3 - \{\text{arcs}\}$ may not be discs. Similarly to before, by sequentially compressing X_3 along nonnullhomotopic simple closed curves disjoint from arcs, if they exist, we obtain a new decorated surface X_4 . See Figure 4 for an illustration of this surgery taking place. Moreover, and this is the main point, X_4 is the same as $\Sigma(\{\sigma_f, \tau_f, \pi_i^*\})$ in the sense that they are related by a decoration-respecting cellular homeomorphism. Compression can only increase the Euler characteristic, so we obtain

$$\chi(\Sigma(\{\sigma_f, \tau_f, \pi_i^*\})) \geq \chi(X_3) = \chi(X_2) = \chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\})) - \sum_i \|\pi_i\|_{k,l}.$$

Combining these two arguments proves the following proposition:

Proposition 4.4 For any given $\{\sigma_f, \tau_f, \pi_i\}$, there exist $\pi_i^* \in S_k \times S_l$ for $i \in [8]$ such that

$$\chi(\Sigma(\{\sigma_f, \sigma_f, \pi_i^*\})) - \sum_{i=1}^8 \|\pi_i^*\|_{k,l} = \chi(\Sigma(\{\sigma_f, \sigma_f, \pi_i^*\})) \geq \chi(\Sigma(\{\sigma_f, \tau_f, \pi_i\})) - \sum_{i=1}^8 \|\pi_i\|_{k,l}.$$

This has the following immediate corollary when combined with Theorem 4.3. Let

$$\text{MATCH}^*(w, k, l)$$

denote the subset of $\text{MATCH}(w, k, l)$ consisting of $\{\sigma_f, \sigma_f, \pi_i\}$ (ie $\sigma_f = \tau_f$ for each $f \in \{a, b, c, d\}$) with $\pi_i \in S_k \times S_l$ for each $i \in [8]$.

Corollary 4.5 For $k + l > 0$ with $\mu \vdash k$ and $\nu \vdash l$ and $w \in [F_4, F_4]$, we have

$$\mathcal{F}_n(w, \mu, \nu) \ll_{w,k,l} n^{\max_{\{\sigma_f, \sigma_f, \pi_i\} \in \text{MATCH}^*(w,k,l)} \chi(\Sigma(\{\sigma_f, \sigma_f, \pi_i\}))}.$$

The benefit to having $\pi_i \in S_k \times S_l$ for $i \in [8]$ is the following. Suppose now that $\{\sigma_f, \sigma_f, \pi_i\} \in \text{MATCH}^*(w, k, l)$. Recall that the boundary loops of $\Sigma(\{\sigma_f, \sigma_f, \pi_i\})$ consist of one w -loop and some number of R^\pm - π -loops. The condition that each $\pi_i \in S_k \times S_l$ means that no π -interval ever connects an endpoint of a R -interval with an endpoint of an R^{-1} -interval. So every boundary component of $\Sigma(\{\sigma_f, \sigma_f, \pi_i\})$ that is not the w -loop contains either only R -intervals or only R^{-1} -intervals, and, in fact, when following the boundary component and reading the directions and labels of the intervals according to traversing each from 0 to 1, reads out a positive power of R (in the former case of only R -intervals) or a negative power of R^{-1} (in the latter case of only R^{-1} -intervals). The sum of the positive powers of R in boundary loops is k , and the sum of the negative powers of R is $-l$. Knowing this boundary structure is extremely important for the arguments in the next sections.

4.3 A topological result that proves Theorem 3.1

Here, in the spirit of [Culler 1981], we explain another way to think about the surfaces $\Sigma(\{\sigma_f, \sigma_f, \pi_i\})$ for $\{\sigma_f, \sigma_f, \pi_i\} \in \text{MATCH}^*(w, k, l)$ that is easier to work with than the construction we gave. At this point we also show how things work for general $g \geq 2$. An *arc* in a surface Σ is a properly embedded interval in Σ with endpoints in the boundary $\partial\Sigma$.

Definition 4.6 For $w \in F_{2g}$, we define $\text{surfaces}(w, k, l)$ to be the set of all decorated surfaces Σ^* as follows. A decorated surface $\Sigma^* \in \text{surfaces}(w, k, l)$ is an oriented surface with boundary, with compatibly oriented boundary components, together with a collection of disjoint embedded arcs that cut Σ^* into topological discs. One boundary component is assigned to be a w -loop, and every other boundary component is assigned to be either a R -loop or an R^{-1} -loop. Each arc is assigned a transverse direction and a label in $\{a_1, b_1, \dots, a_g, b_g\}$. Every arc-endpoint in $\partial\Sigma^*$ inherits a transverse direction and label from the assigned direction and label of its arc. We require that Σ^* satisfy the following properties:

- (P1) When one follows the w -loop according to its assigned orientation, and reads f when an f -labeled arc-endpoint is traversed in its given direction, and f^{-1} when an f -labeled arc-endpoint is traversed counter to its given direction, one reads a cyclic rotation of w in reduced form, depending on where one begins to read.

- (P2) When one follows any R -loop according to its assigned orientation in the same way as before, one reads (a cyclic rotation) of some positive power of R_g in reduced form. The sum of these positive powers over all R -loops is k .
- (P3) When one follows any R^{-1} -loop according to its assigned orientation in the same way as before, one reads (a cyclic rotation) of some negative power of R_g in reduced form. The sum of these negative powers over all R^{-1} -loops is $-l$.
- (P4) No arc connects an R -loop to an R^{-1} -loop.

Given a surface $\Sigma(\{\sigma_f, \sigma_g, \pi_i\})$ with $\{\sigma_f, \sigma_g, \pi_i\} \in \text{MATCH}^*(w, k, l)$, all the type II discs of the surface are rectangles. Hence, by collapsing each w -interval, R -interval and R^{-1} -interval to a point, and collapsing every type II rectangle to an arc, we obtain a CW-complex that is a surface with boundary, cut into discs by arcs. Every arc inherits a transverse direction and label from the compatible assigned directions and labels of the intervals in the boundary of its originating type II rectangle. We call this modified surface $\Sigma^* = \Sigma^*(\{\sigma_f, \pi_i\})$. It clearly satisfies (P1)–(P3) and (P4) follows from the *forbidden matchings* property. (Of course, when $g = 2$, we identify $\{a, b, c, d\}$ with $\{a_1, b_1, a_2, b_2\}$.) We also have $\chi(\Sigma(\{\sigma_f, \sigma_g, \pi_i\})) = \chi(\Sigma^*(\{\sigma_f, \pi_i\}))$. With Definition 4.6 and the remarks proceeding it, we can now state a further consequence of Corollary 4.5 as it extends to general $g \geq 2$:

Corollary 4.7 For $k + l > 0$ with $\mu \vdash k$ and $\nu \vdash l$ and $w \in [F_{2g}, F_{2g}]$, as $n \rightarrow \infty$,

$$\mathcal{J}_n(w, \mu, \nu) \ll_{w,k,l} n^{\max\{\chi(\Sigma^*): \Sigma^* \in \text{surfaces}(w,k,l)\}}.$$

In order for Corollary 4.7 to give us strong enough results, it needs to be combined with the following nontrivial topological bound:

Proposition 4.8 If $w \in [F_{2g}, F_{2g}]$ is a shortest element representing the conjugacy class of $\gamma \in \Gamma_g$, $w \neq \text{id}$ and $\Sigma^* \in \text{surfaces}(w, k, l)$, then $\chi(\Sigma^*) \leq -(k + l)$.

Remark 4.9 Proposition 4.8 is by no means a trivial statement and one has to use that w is a shortest element representing the conjugacy class of some element of Γ_g . For example, if $w = R_g$, then w represents the conjugacy class of id_{Γ_g} , but for $k = 0$ and $l = 1$ there is an “obvious” annulus in $\text{surfaces}(w, 0, 1)$. This has $\chi = 0 > -(k + l) = -1$. Proposition 4.8 also requires $w \neq \text{id}$; if $w = \text{id}$ then for $k = 0$ and $l = 1$ one can take a disc with no arcs as a valid element of $\text{surfaces}(\text{id}, 0, 0)$. This has $\chi = 1 > -(k + l) = 0$. In fact this disc is ultimately responsible for $\mathbb{E}_{g,n}[\text{Tr}_{\text{id}}] = n$.

The proof of Proposition 4.8 is self-contained and given in Section 4.5. Before doing this, we prove Theorem 3.1.

Proof of Theorem 3.1 given Proposition 4.8 Since Theorem 3.1 was proved when $k = l = 0$ in Section 3.2, we can assume $k + l > 0$. Then combining Corollary 4.7 and Proposition 4.8 gives

$$\mathcal{J}_n(w, \mu, \nu) \ll_{w,k,l} n^{-(k+l)}.$$

On the other hand, $D_{\mu,\nu}(n) = O(n^{k+l})$ from (2-1). Therefore, $D_{\mu,\nu}(n)\mathcal{J}_n(w, \mu, \nu) \ll_{w,k,l} 1$. \square

4.4 Work of Dehn and Birman–Series

As we mentioned in Section 3.1, to prove Proposition 4.8 we have to use the fact that $w \in [F_{2g}, F_{2g}]$ is a shortest element representing the conjugacy class of $\gamma \in \Gamma_g$. We use a combinatorial characterization of such words that stems from Dehn’s algorithm [1912] for solving the problem of whether a given word represents the identity in Γ_g . The ideas of Dehn’s algorithm were refined in [Birman and Series 1987]. Magee and Puder [2023] used Birman and Series’ results (alongside other methods) to obtain the analog of Theorem 1.2 when the family of groups $SU(n)$ is replaced by the family of symmetric groups S_n . Similar consequences of the work of Dehn, Birman and Series that we used in [loc. cit.] will be used here.

We now follow the language of [Magee and Puder 2023] to state the results we need in this paper. These results are simple and direct consequences of the work of Birman and Series.

We view the universal cover of Σ_g as a disc tiled by $4g$ -gons that we call U . We assume every edge of this tiling is directed and labeled by some element of $\{a_1, b_1, \dots, a_g, b_g\}$ such that when we read counterclockwise along the boundary of any octagon we read the reduced cyclic word $[a_1, b_1] \cdots [a_g, b_g]$. By fixing a basepoint $u \in U$, we obtain a free cellular action of Γ_g on U that respects the labels and directions of edges and identifies the quotient $\Gamma_g \backslash U$ with Σ_g ; this gives a description of Σ_g as a $4g$ -gon with glued sides, as is typical.

Now suppose that $\gamma \in \Gamma$ is not the identity. The quotient $A_\gamma := \langle \gamma \rangle \backslash U$ of U by the cyclic group generated by γ is an open annulus tiled by infinitely many $4g$ -gons. The edges of A_γ inherit directions and labels from those of the edges of U . The point $u \in U$ maps to some point, denoted by $x_0 \in A_\gamma$.

Now let $w \in F_{2g}$ be an element that represents γ , and identify w with a combinatorial word by writing w in reduced form. Beginning at x_0 , and following the path spelled out by w beginning at x_0 , we obtain an oriented closed loop L_w in the one-skeleton of A_γ . If w is a shortest element representing the conjugacy class of γ , then this loop L_w must not have self-intersections. In this case, which we from now assume, L_w is therefore a topologically embedded circle in the annulus A_γ that is nonnullhomotopic and cuts A_γ into two annuli A_γ^\pm .

Every vertex of A_γ has $4g$ incident half-edges each of which has an orientation and direction given by the edge they are in. Going clockwise, the cyclic order of the half-edges incident at any vertex is:

a_1 -outgoing, b_1 -incoming, a_1 -incoming, b_1 -outgoing, \dots , a_g -outgoing, b_g -incoming, a_g -incoming, b_g -outgoing.

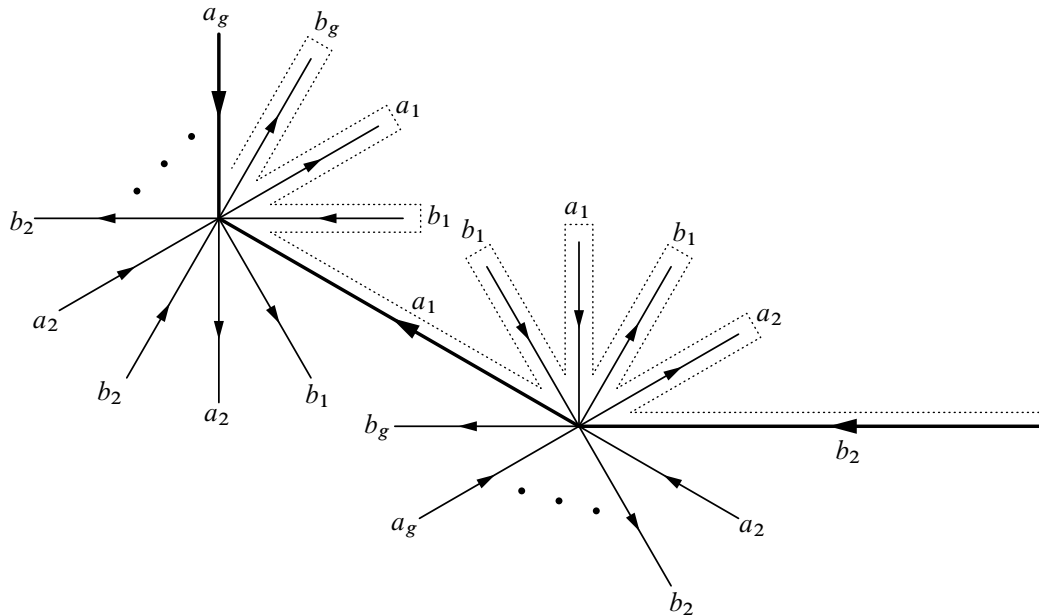


Figure 5: A piece P of \hat{L}_w in the case when the reduced form of w contains $a_g a_1^{-1} b_2^{-1}$ as a subword. The edges of L_w are in bold. The piece is indicated by the dotted lines. This piece P has $\epsilon(P) = 2$, $\mathfrak{h}\epsilon(P) = 7$ and $\chi(P) = 1$. Note that a piece may also run along the other side of L_w .

We define \hat{L}_w to be the loop L_w with all incident half edges in A_γ attached. We call the new half-edges added *hanging half-edges*.

Moreover, we thicken up \hat{L}_w by viewing each edge of L_w as a rectangle, each hanging half-edge as a half-rectangle, and each vertex replaced by a disc. In other words, we take a small neighborhood of \hat{L}_w in A_γ . We now think of \hat{L}_w as the thickened version. This is a topological annulus, where the hanging half-edges have become stubs hanging off. A *piece of \hat{L}_w* is a contiguous collection of hanging half-rectangles and rectangle sides following edges of L_w in the boundary of \hat{L}_w . Such a piece is in either A_γ^+ or A_γ^- . Given a piece P of \hat{L}_w we write $\epsilon(P)$ for the number of rectangle sides following edges of L_w , and $\mathfrak{h}\epsilon(P)$ for the number of hanging-half edges in P . We say that a piece P has Euler characteristic $\chi(P) = 0$ if it follows an entire boundary component of \hat{L}_w , and $\chi(P) = 1$ otherwise, as we view it as an interval running along the rectangle sides and around the sides of the hanging half-rectangles. See Figure 5 for an illustration of a piece of \hat{L}_w .

Birman and Series [1987, Theorem 2.12(a)] prove that, if w is a shortest element representing the conjugacy class of $\gamma \in \Gamma_g$, then there are strong restrictions on the pieces of \hat{L}_w that can appear. This has the following consequence, which is given by³ [Magee and Puder 2023, Proof of Lemma 5.18]:

³We stress that Lemma 4.10 is a straightforward consequence of Birman and Series' work, so, even though we cite [Magee and Puder 2023], this paper does not depend on that work in any significant way.

Lemma 4.10 *If w is a shortest element representing the conjugacy class of $\gamma \in \Gamma_g$, and both γ and hence w are nonidentity, then for any piece P of \hat{L}_w , we have*

$$\epsilon(P) \leq (2g - 1)\mathfrak{h}\epsilon(P) + 2g\chi(P).$$

Proof Since w is a shortest element representing some nonidentity conjugacy class in Γ_g , in the language of [Magee and Puder 2023], L_w is a boundary reduced tiled surface. Then the proof of [Magee and Puder 2023, Lemma 5.18] contains the result stated in the lemma. The basic idea of the proof is not complicated and goes back to [Dehn 1912]: if there are too many edges (ie $\epsilon(P)$ is large) then one can find a string of letters in the reduced word of w (eg $aba^{-1}b^{-1}c$) that can be shortened using the relator R (eg $aba^{-1}b^{-1}c = dcd^{-1}$). \square

This inequality plays a crucial role in the next section.

4.5 Proof of Proposition 4.8

Suppose that $g \geq 2$ and $w \in [F_{2g}, F_{2g}]$ is a nonidentity shortest element representing the conjugacy class of $\gamma \in \Gamma_g$. In particular, w is cyclically reduced. We let $R = R_g$. Now fix $k, l \in \mathbb{N}_0$ and suppose $\Sigma^* \in \text{surfaces}(w, k, l)$. The arcs of Σ^* are of three different types:

- (WR) An arc with one endpoint in the w -loop and one endpoint in an R - or R^{-1} -loop.
- (RR) An arc with both endpoints in R - or R^{-1} -loops. By property (P4), the endpoints of such an arc are both in R -loops or both in R^{-1} -loops.
- (WW) An arc with both endpoints in the w -loop.

The boundary of any disc of Σ^* alternates between segments of $\partial\Sigma^*$ and arcs. A disc is a *pre-piece disc* if its boundary contains exactly one segment of the w -loop. A disc is called a *junction disc* if it is not a pre-piece disc. We say that a junction disc is *piece-adjacent* if it meets a WR-arc-side.

To be precise, we view all discs as open discs, and hence not containing any arcs. A disc meets certain arc-sides along its boundary; it is possible for a disc to meet both sides of the same arc and we view this scenario as the disc meeting two separate arc-sides. We say an arc-side has the same type WR/RR/WW as its corresponding arc.

Note that any pre-piece disc cannot meet any WW-arc-side: if it did, the disc could only meet this one arc-side together with one segment of the w -loop and this would contradict the fact that w is cyclically reduced since the arc matches a letter f with a cyclically adjacent letter f^{-1} of w . It is also clear that any pre-piece disc meets exactly two WR-arc-sides: the ones that emanate from the sole segment of the w -loop. So, in light of (P4), a pre-piece disc takes one of the forms shown in Figure 6.

We define a *piece of Σ^** to be a connected component of

$$\{\text{pre-piece discs}\} \cup \{\text{WR-arcs}\}.$$

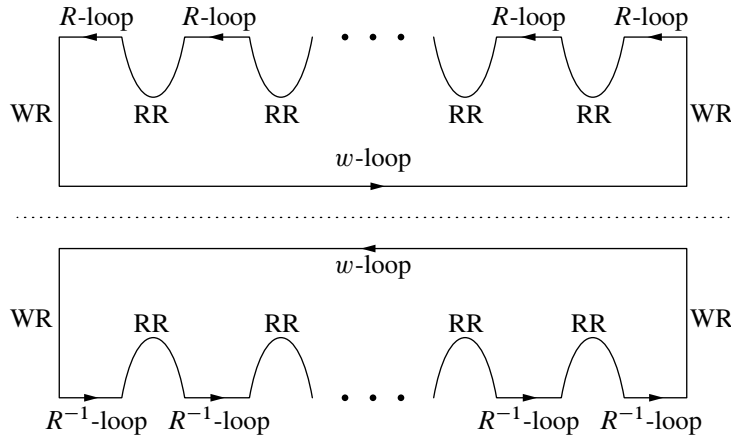


Figure 6: Possible forms of pre-piece discs. The number of R -loop segments or R^{-1} -loop segments is at least 1 and bounded given k and l . The arrows denote the orientations of the boundary loops.

A piece of Σ^* is therefore either a contiguous collection of pre-piece discs that meet only along WR-arcs, or a single WR-arc. If P is a piece of Σ^* , either $\chi(P) = 1$, or $\chi(P) = 0$, in which case P meets the entire w -loop and is the unique piece.

We now have *two* definitions of pieces: pieces of \widehat{L}_w and pieces of Σ^* . These are, as the names suggest, closely related, and this is the key observation in the proof of Proposition 4.8. Indeed, the reader should carefully consider Figure 7, which leads to the following lemma. In analogy to pieces of \widehat{L}_w , if P is any piece of Σ^* , we write $\epsilon(P)$ for the number of WR-arcs in P , and $\mathfrak{h}\epsilon(P)$ for the number of RR-arc-sides that meet P (this is zero if P is a single WR-arc).

Lemma 4.11 *If w is a shortest element representing the conjugacy class of $\gamma \in \Gamma_g$, $k, l \in \mathbb{N}_0$ and $\Sigma^* \in \text{surfaces}(w, k, l)$, then, for any piece P of Σ^* , we have*

$$\epsilon(P) \leq (2g - 1)\mathfrak{h}\epsilon(P) + 2g\chi(P).$$

Proof Given any piece P of Σ^* , it contains a consecutive (possibly cyclic) series of WR-arcs that correspond to a contiguous collection of edges in the loop L_w . The discs of P correspond to certain vertices of L_w ; each of these vertices has two emanating half-edges belonging to the edges defined by WR-arcs of P . The piece P can either meet only R -loops or meet only R^{-1} -loops.

We define a piece P' of \widehat{L}_w corresponding to P as follows. If P meets R -loops, then P' consists of rectangle sides along the edges of L_w corresponding to the WR-arcs of P together with all hanging half-edges at vertices corresponding to discs of P that are on the *left* of L_w as it is traversed in its assigned orientation (corresponding to reading w along L_w). If P' meets R^{-1} -loops, then P' is defined similarly with the modification that we include instead hanging half-edges on the *right* of L_w . Figure 7 together

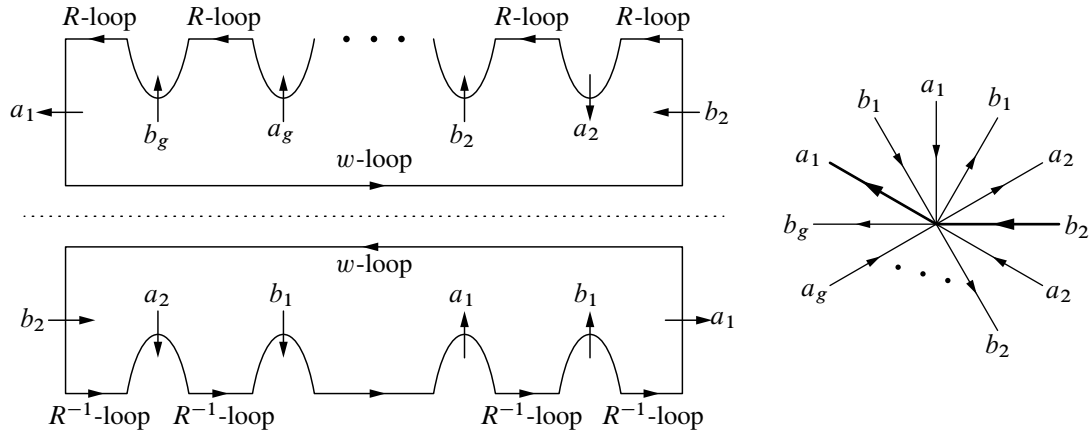


Figure 7: Given a segment of the w -loop corresponding to a juncture between letters $a_1^{-1}b_2^{-1}$ in w , if this segment is part of a pre-piece disc then some possible forms of that disc are shown above. This juncture between letters of w corresponds to a vertex in L_w . The right-hand illustration shows the neighborhood of this vertex in the annulus A_γ , where the bold lines correspond to half-edges of L_w . The right-hand picture actually almost determines the left-hand pictures. Indeed, given the a_1 -arc on the top-left, the next arc has to be a b_g -arc with the given direction, since *only* b_g^{-1} cyclically precedes a_1 in R_g or any power of R_g . Then the next arc a_g with its direction is determined since *only* a_g cyclically precedes b_g in R_g . This continues until an arc labeled by b_2 and with an incoming direction is reached, as in the right arc of the top-left picture. At this point, the boundary of the disc may close up. (This is analogous to what happens in the bottom picture, where an analogous pattern occurs.) The only indeterminacy is that after reaching a b_2 -arc with an incoming direction for the first time, the *entire pattern* shown in the right-hand picture may repeat any number of times, as long as k and l allow it. The upshot of this is that any pre-piece disc has at least as many incident RR-arc-sides as there are hanging half-edges on the corresponding side of L_w , at the corresponding vertex.

with its captioned discussion now shows that

$$\mathfrak{h}\epsilon(P') \leq \mathfrak{h}\epsilon(P),$$

and $\epsilon(P) = \epsilon(P')$ by construction. We also have $\chi(P') = \chi(P)$. Therefore Lemma 4.10 applied to P' implies

$$\epsilon(P) = \epsilon(P') \leq (2g - 1)\mathfrak{h}\epsilon(P') + 2g\chi(P') \leq (2g - 1)\mathfrak{h}\epsilon(P) + 2g\chi(P). \quad \square$$

Let N_{RR} be the number of RR-arcs, N_{WR} the number of WR-arcs, and N_{WW} the number of WW-arcs in Σ^* . In the following, we refer to discs of Σ^* simply as discs. Since there are $4g(k + l)$ incidences between arcs and R -loops or R^{-1} -loops, we have

$$(4-3) \quad 2N_{RR} + N_{WR} = 4g(k + l).$$

Let Σ_1 be the surface formed by cutting Σ^* along all RR-arcs. We have

$$\chi(\Sigma_1) = \sum_{\text{discs } D} (1 - \frac{1}{2}d'(D)),$$

where $d'(D)$ is the number of arc-sides meeting D that are *not* of type RR. This formula holds because $d'(D)$ is the degree of the disc D in the dual graph G_1 of Σ_1 , the right-hand side is easily seen to be $\chi(G_1) = V(G_1) - E(G_1)$, and, since Σ_1 deformation retracts to an obvious embedded copy of G_1 , $\chi(G_1) = \chi(\Sigma_1)$. We partition the sum above according to $\chi(\Sigma_1) = S_0 + S_1 + S_2$, where

$$\begin{aligned} S_0 &:= \sum_{\text{pre-piece discs } D} (1 - \frac{1}{2}d'(D)), \\ S_1 &:= \sum_{\text{piece-adjacent junction discs } D} (1 - \frac{1}{2}d'(D)), \\ S_2 &:= \sum_{\text{not piece-adjacent junction discs } D} (1 - \frac{1}{2}d'(D)). \end{aligned}$$

Note first that a pre-piece disc has $d'(D) = 2$ (see Figure 6). Hence $S_0 = 0$. We deal with S_1 next. For a disc D of Σ^* , let $d_{\text{WR}}(D)$ denote the number of WR-arc-sides meeting D . Note that a piece-adjacent junction disc D has $d_{\text{WR}}(D) > 0$ by definition. We rewrite S_1 as

$$\begin{aligned} (4-4) \quad S_1 &= \sum_{\text{piece-adjacent junction discs } D} (1 - \frac{1}{2}d'(D)) \frac{1}{d_{\text{WR}}(D)} \sum_{\text{incidences between } D \text{ and WR-arc-sides}} 1 \\ &= \sum_{\text{pieces } P} \sum_{\text{incidences between } P \text{ and some junction disc } D \text{ along WR-arc}} Q(D), \end{aligned}$$

where, for a piece-adjacent junction disc D ,

$$Q(D) := \frac{1}{d_{\text{WR}}(D)} (1 - \frac{1}{2}d'(D)).$$

Suppose that D is a piece-adjacent junction disc. By parity considerations, $d_{\text{WR}}(D)$ is even. We estimate $Q(D)$ by splitting into two cases. If $d_{\text{WR}}(D) = 2$ then $d'(D) \geq 3$, since, otherwise, D would meet only two WR-arc-sides and other RR-arc-sides, whence be a pre-piece disc and not be a junction disc. In this case,

$$Q(D) = \frac{1}{2} (1 - \frac{1}{2}d'(D)) \leq \frac{1}{2} (1 - \frac{3}{2}) = -\frac{1}{4}.$$

Otherwise, $d_{\text{WR}}(D) \geq 4$ and, since $d'(D) \geq d_{\text{WR}}(D)$, we have

$$Q(D) \leq \frac{1}{d_{\text{WR}}(D)} (1 - \frac{1}{2}d_{\text{WR}}(D)) = \frac{1}{d_{\text{WR}}(D)} - \frac{1}{2} \leq \frac{1}{4} - \frac{1}{2} = -\frac{1}{4}.$$

So we have proved that, for all piece-adjacent junction discs D , $Q(D) \leq -\frac{1}{4}$. Putting this into (4-4) gives

$$\begin{aligned} (4-5) \quad S_1 &\leq -\frac{1}{4} \sum_{\text{pieces } P} \sum_{\text{incidences between } P \text{ and some junction disc } D \text{ along WR-arc}} 1 \\ &= -\frac{1}{4} \sum_{\text{pieces } P} 2\chi(P) = -\frac{1}{2} \sum_{\text{pieces } P} \chi(P). \end{aligned}$$

We now turn to S_2 . Here is the key moment where $w \neq \text{id}$ is used.⁴ Since $w \neq \text{id}$, any disc must meet an arc. Indeed, the only other possibility is that the boundary of the disc is an entire boundary loop that has no emanating arcs. This hypothetical boundary loop cannot be an R - or R^{-1} -loop, so it has to be the w -loop. But this would entail $w = \text{id}$.

Hence any disc contributing to S_2 meets no WR -arc-side, but meets some arc-side. Therefore it meets only WW -arcs or only RR -arcs. Every disc D contributing to S_2 meeting only WW -arcs gives a nonpositive contribution since w is cyclically reduced, and hence $d'(D) \geq 2$. Every disc D contributing to S_2 meeting only RR -arcs, which we will call an RR -disc, has $d'(D) = 0$ and hence contributes 1 to S_2 .

This shows

$$(4-6) \quad S_2 \leq \#\{\text{RR-discs}\}.$$

In total, combining $S_0 = 0$ with (4-5) and (4-6), we get

$$\chi(\Sigma_1) \leq \#\{\text{RR-discs}\} - \frac{1}{2} \sum_{\text{pieces } P \text{ of } \Sigma^*} \chi(P).$$

To obtain Σ^* from Σ_1 we have to glue all cut RR -arcs, of which there are N_{RR} . Each gluing decreases χ by 1, so

$$(4-7) \quad \chi(\Sigma^*) \leq \#\{\text{RR-discs}\} - N_{RR} - \frac{1}{2} \sum_{\text{pieces } P \text{ of } \Sigma^*} \chi(P).$$

Using Lemma 4.11 with the above gives

$$(4-8) \quad \begin{aligned} \chi(\Sigma^*) &\leq \#\{\text{RR-discs}\} - N_{RR} - \frac{1}{2} \sum_{\text{pieces } P \text{ of } \Sigma^*} \chi(P) \\ &\leq \#\{\text{RR-discs}\} - N_{RR} - \frac{1}{4g} \sum_{\text{pieces } P \text{ of } \Sigma^*} \epsilon(P) + \frac{2g-1}{4g} \sum_{\text{pieces } P \text{ of } \Sigma^*} \eta\epsilon(P) \\ &= \#\{\text{RR-discs}\} - N_{RR} - \frac{N_{WR}}{4g} + \frac{2g-1}{4g} \sum_{\text{pieces } P \text{ of } \Sigma^*} \eta\epsilon(P). \end{aligned}$$

Let $\eta\epsilon'(\Sigma^*)$ denote the total number of RR -arc-sides meeting RR -discs. Every RR -disc has to meet at least $4g$ arc-sides; this observation is similar to the reasoning in Figure 7. Therefore

$$(4-9) \quad \eta\epsilon'(\Sigma^*) \geq 4g\#\{\text{RR-discs}\}.$$

Every RR -arc-side either meets a piece P and contributes to $\eta\epsilon(P)$ or a disc meeting only RR -arc-sides and contributes to $\eta\epsilon'(\Sigma^*)$. Hence

$$(4-10) \quad \eta\epsilon'(\Sigma^*) + \sum_{\text{pieces } P \text{ of } \Sigma^*} \eta\epsilon(P) = 2N_{RR}.$$

⁴Although, technically, $w \neq \text{id}$ was used to define L_w and pieces etc, if w is the identity, the proof of Proposition 4.8 could, a priori, circumvent these definitions.

Combining (4-3), (4-9) and (4-10) with (4-8) gives

$$\begin{aligned}
 \chi(\Sigma^*) &\leq \frac{\mathfrak{h}\epsilon'(\Sigma^*)}{4g} - N_{\text{RR}} - \frac{N_{\text{WR}}}{4g} + \frac{2g-1}{4g} \sum_{\text{pieces } P \text{ of } \Sigma^*} \mathfrak{h}\epsilon(P) && \text{(by (4-9))} \\
 &= \frac{\mathfrak{h}\epsilon'(\Sigma^*)}{4g} - N_{\text{RR}} - \frac{N_{\text{WR}}}{4g} + \frac{(2g-1)N_{\text{RR}}}{2g} - \frac{2g-1}{4g} \mathfrak{h}\epsilon'(\Sigma^*) && \text{(by (4-10))} \\
 &= -\frac{1}{4g}(2N_{\text{RR}} + N_{\text{WR}}) - \frac{2g-2}{4g} \mathfrak{h}\epsilon'(\Sigma^*) \\
 &\leq -\frac{1}{4g}(2N_{\text{RR}} + N_{\text{WR}}) = -\frac{4g(k+l)}{4g} && \text{(by (4-3))} \\
 &= -(k+l).
 \end{aligned}$$

This completes the proof of Proposition 4.8. □

5 Proof of the main theorem

Proof of Theorem 1.2 Assume $\gamma \in [\Gamma_g, \Gamma_g]$ is not the identity and that $w \in [F_{2g}, F_{2g}]$ is a shortest element representing the conjugacy class of γ , hence also not the identity. By Corollary 2.10, we have

$$\mathbb{E}_{g,n}[\text{Tr}_\gamma] = \zeta(2g-2; n)^{-1} \sum_{(\mu, \nu) \in \tilde{\Omega}} D_{\mu, \nu}(n) \mathcal{F}_n(w, \nu, \mu) + O_{w,g}\left(\frac{1}{n}\right),$$

where $\tilde{\Omega}$ is a finite collection of pairs of Young diagrams. We know $\lim_{n \rightarrow \infty} \zeta(2g-2; n) = 1$ from (2-11) and, for each fixed (μ, ν) , $D_{\mu, \nu}(n) \mathcal{F}_n(w, \nu, \mu) = D_{\nu, \mu}(n) \mathcal{F}_n(w, \nu, \mu) = O_{w, \mu, \nu}(1)$ by Theorem 3.1. Hence $\mathbb{E}_{g,n}[\text{Tr}_\gamma] = O_\gamma(1)$ as $n \rightarrow \infty$, as required. □

References

- [Atiyah and Bott 1983] **MF Atiyah, R Bott**, *The Yang–Mills equations over Riemann surfaces*, Philos. Trans. Roy. Soc. A 308 (1983) 523–615 MR Zbl
- [Benkart et al. 1994] **G Benkart, M Chakrabarti, T Halverson, R Leduc, C Lee, J Stroomeer**, *Tensor product representations of general linear groups and their connections with Brauer algebras*, J. Algebra 166 (1994) 529–567 MR Zbl
- [Birman and Series 1987] **JS Birman, C Series**, *Dehn’s algorithm revisited, with applications to simple curves on surfaces*, from “Combinatorial group theory and topology” (SM Gersten, JR Stallings, editors), Ann. of Math. Stud. 111, Princeton Univ. Press (1987) 451–478 MR Zbl
- [Collins 2003] **B Collins**, *Moments and cumulants of polynomial random variables on unitary groups, the Itzykson–Zuber integral, and free probability*, Int. Math. Res. Not. 2003 (2003) 953–982 MR Zbl
- [Collins and Male 2014] **B Collins, C Male**, *The strong asymptotic freeness of Haar and deterministic matrices*, Ann. Sci. École Norm. Sup. 47 (2014) 147–163 MR Zbl

- [Collins and Śniady 2006] **B Collins, P Śniady**, *Integration with respect to the Haar measure on unitary, orthogonal and symplectic group*, *Comm. Math. Phys.* 264 (2006) 773–795 MR Zbl
- [Culler 1981] **M Culler**, *Using surfaces to solve equations in free groups*, *Topology* 20 (1981) 133–145 MR Zbl
- [Dahlqvist and Lemoine 2023] **A Dahlqvist, T Lemoine**, *Large N limit of Yang–Mills partition function and Wilson loops on compact surfaces*, *Probab. Math. Phys.* 4 (2023) 849–890 MR Zbl
- [Dehn 1912] **M Dehn**, *Transformation der Kurven auf zweiseitigen Flächen*, *Math. Ann.* 72 (1912) 413–421 MR Zbl
- [Enomoto and Izumi 2016] **T Enomoto, M Izumi**, *Indecomposable characters of infinite dimensional groups associated with operator algebras*, *J. Math. Soc. Japan* 68 (2016) 1231–1270 MR Zbl
- [Fulton and Harris 1991] **W Fulton, J Harris**, *Representation theory: a first course*, *Graduate Texts in Math.* 129, Springer (1991) MR Zbl
- [Goldman 1984] **W M Goldman**, *The symplectic nature of fundamental groups of surfaces*, *Adv. Math.* 54 (1984) 200–225 MR Zbl
- [Gromov and Milman 1983] **M Gromov, V D Milman**, *A topological application of the isoperimetric inequality*, *Amer. J. Math.* 105 (1983) 843–854 MR Zbl
- [Guralnick et al. 2012] **R Guralnick, M Larsen, C Manack**, *Low degree representations of simple Lie groups*, *Proc. Amer. Math. Soc.* 140 (2012) 1823–1834 MR Zbl
- [Haagerup and Thorbjørnsen 2005] **U Haagerup, S Thorbjørnsen**, *A new application of random matrices: $\text{Ext}(C_{\text{red}}^*(F_2))$ is not a group*, *Ann. of Math.* 162 (2005) 711–775 MR Zbl
- [Häsä and Stasinski 2019] **J Häsä, A Stasinski**, *Representation growth of compact linear groups*, *Trans. Amer. Math. Soc.* 372 (2019) 925–980 MR Zbl
- [Koike 1989] **K Koike**, *On the decomposition of tensor products of the representations of the classical groups: by means of the universal characters*, *Adv. Math.* 74 (1989) 57–86 MR Zbl
- [Larsen and Lubotzky 2008] **M Larsen, A Lubotzky**, *Representation growth of linear groups*, *J. Eur. Math. Soc.* 10 (2008) 351–390 MR Zbl
- [Lemoine 2022] **T Lemoine**, *Large N behaviour of the two-dimensional Yang–Mills partition function*, *Combin. Probab. Comput.* 31 (2022) 144–165 MR Zbl
- [Magee 2022] **M Magee**, *Random unitary representations of surface groups, I: Asymptotic expansions*, *Comm. Math. Phys.* 391 (2022) 119–171 MR Zbl
- [Magee and Puder 2015] **M Magee, D Puder**, *Word measures on unitary groups*, preprint (2015) arXiv 1509.07374
- [Magee and Puder 2019] **M Magee, D Puder**, *Matrix group integrals, surfaces, and mapping class groups, I: $\mathcal{U}(n)$* , *Invent. Math.* 218 (2019) 341–411 MR Zbl
- [Magee and Puder 2023] **M Magee, D Puder**, *The asymptotic statistics of random covering surfaces*, *Forum Math. Pi* 11 (2023) art. id. e15 MR Zbl
- [Mirzakhani 2007] **M Mirzakhani**, *Simple geodesics and Weil–Peterson volumes of moduli spaces of bordered Riemann surfaces*, *Invent. Math.* 167 (2007) 179–222 MR Zbl
- [Nica 1994] **A Nica**, *On the number of cycles of given length of a free word in several random permutations*, *Random Structures Algorithms* 5 (1994) 703–730 MR Zbl

- [Voiculescu 1985] **D Voiculescu**, *Symmetries of some reduced free product C^* -algebras*, from “Operator algebras and their connections with topology and ergodic theory” (H Araki, C C Moore, Ş Strătilă, D Voiculescu, editors), Lecture Notes in Math. 1132, Springer (1985) 556–588 MR Zbl
- [Voiculescu 1986] **D Voiculescu**, *Addition of certain noncommuting random variables*, J. Funct. Anal. 66 (1986) 323–346 MR Zbl
- [Voiculescu 1987] **D Voiculescu**, *Multiplication of certain noncommuting random variables*, J. Operator Theory 18 (1987) 223–235 MR Zbl
- [Voiculescu 1990] **D Voiculescu**, *Noncommutative random variables and spectral problems in free product C^* -algebras*, Rocky Mountain J. Math. 20 (1990) 263–283 MR Zbl
- [Voiculescu 1991] **D Voiculescu**, *Limit laws for random matrices and free products*, Invent. Math. 104 (1991) 201–220 MR Zbl
- [Voiculescu et al. 1992] **D V Voiculescu, K J Dykema, A Nica**, *Free random variables*, CRM Monogr. Ser. 1, Amer. Math. Soc., Providence, RI (1992) MR Zbl
- [Weingarten 1978] **D Weingarten**, *Asymptotic behavior of group integrals in the limit of infinite rank*, J. Math. Phys. 19 (1978) 999–1001 MR Zbl
- [Witten 1991] **E Witten**, *On quantum gauge theories in two dimensions*, Comm. Math. Phys. 141 (1991) 153–209 MR Zbl
- [Xu 1997] **F Xu**, *A random matrix model from two-dimensional Yang–Mills theory*, Comm. Math. Phys. 190 (1997) 287–307 MR Zbl
- [Zagier 1994] **D Zagier**, *Values of zeta functions and their applications*, from “First European Congress of Mathematics, II” (A Joseph, F Mignot, F Murat, B Prum, R Rentschler, editors), Progr. Math. 120, Birkhäuser, Basel (1994) 497–512 MR Zbl

*Department of Mathematical Sciences, Durham University
Durham, United Kingdom*

michael.r.magee@durham.ac.uk

<https://www.mmagee.net/>

Proposed: David Fisher

Seconded: John Lott, Leonid Polterovich

Received: 2 December 2021

Revised: 3 January 2023

GEOMETRY & TOPOLOGY

msp.org/gt

MANAGING EDITORS

Robert Lipshitz University of Oregon
lipshitz@uoregon.edu

András I Stipsicz Alfréd Rényi Institute of Mathematics
stipsicz@renyi.hu

BOARD OF EDITORS

Mohammed Abouzaid	Stanford University abouzaid@stanford.edu	Mark Gross	University of Cambridge mgross@dpmms.cam.ac.uk
Dan Abramovich	Brown University dan_abramovich@brown.edu	Rob Kirby	University of California, Berkeley kirby@math.berkeley.edu
Ian Agol	University of California, Berkeley ianagol@math.berkeley.edu	Bruce Kleiner	NYU, Courant Institute bkleiner@cims.nyu.edu
Arend Bayer	University of Edinburgh arend.bayer@ed.ac.uk	Sándor Kovács	University of Washington skovacs@uw.edu
Mark Behrens	University of Notre Dame mbehren1@nd.edu	Urs Lang	ETH Zürich urs.lang@math.ethz.ch
Mladen Bestvina	University of Utah bestvina@math.utah.edu	Marc Levine	Universität Duisburg-Essen marc.levine@uni-due.de
Martin R Bridson	University of Oxford bridson@maths.ox.ac.uk	Ciprian Manolescu	University of California, Los Angeles cm@math.ucla.edu
Jim Bryan	University of British Columbia jbryan@math.ubc.ca	Haynes Miller	Massachusetts Institute of Technology hrm@math.mit.edu
Dmitri Burago	Pennsylvania State University burago@math.psu.edu	Tomasz Mrowka	Massachusetts Institute of Technology mrowka@math.mit.edu
Tobias H Colding	Massachusetts Institute of Technology colding@math.mit.edu	Aaron Naber	Northwestern University anaber@math.northwestern.edu
Simon Donaldson	Imperial College, London s.donaldson@ic.ac.uk	Peter Ozsváth	Princeton University petero@math.princeton.edu
Yasha Eliashberg	Stanford University eliash-gt@math.stanford.edu	Leonid Polterovich	Tel Aviv University polterov@post.tau.ac.il
Benson Farb	University of Chicago farb@math.uchicago.edu	Colin Rourke	University of Warwick gt@maths.warwick.ac.uk
David M Fisher	Rice University davidfisher@rice.edu	Roman Sauer	Karlsruhe Institute of Technology roman.sauer@kit.edu
Mike Freedman	Microsoft Research michaelf@microsoft.com	Stefan Schwede	Universität Bonn schwede@math.uni-bonn.de
David Gabai	Princeton University gabai@princeton.edu	Natasa Sesum	Rutgers University natasas@math.rutgers.edu
Stavros Garoufalidis	Southern U. of Sci. and Tech., China stavros@mpim-bonn.mpg.de	Gang Tian	Massachusetts Institute of Technology tian@math.mit.edu
Cameron Gordon	University of Texas gordon@math.utexas.edu	Ulrike Tillmann	Oxford University tillmann@maths.ox.ac.uk
Jesper Grodal	University of Copenhagen jg@math.ku.dk	Nathalie Wahl	University of Copenhagen wahl@math.ku.dk
Misha Gromov	IHÉS and NYU, Courant Institute gromov@ihes.fr	Anna Wienhard	Universität Heidelberg wienhard@mathi.uni-heidelberg.de

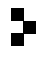
See inside back cover or msp.org/gt for submission instructions.

The subscription price for 2025 is US \$865/year for the electronic version, and \$1210/year (+\$75, if shipping outside the US) for print and electronic. Subscriptions, requests for back issues and changes of subscriber address should be sent to MSP. Geometry & Topology is indexed by Mathematical Reviews, Zentralblatt MATH, Current Mathematical Publications and the Science Citation Index.

Geometry & Topology (ISSN 1465-3060 printed, 1364-0380 electronic) is published 9 times per year and continuously online, by Mathematical Sciences Publishers, c/o Department of Mathematics, University of California, 798 Evans Hall #3840, Berkeley, CA 94720-3840. Periodical rate postage paid at Oakland, CA 94615-9651, and additional mailing offices. POSTMASTER: send address changes to Mathematical Sciences Publishers, c/o Department of Mathematics, University of California, 798 Evans Hall #3840, Berkeley, CA 94720-3840.

GT peer review and production are managed by EditFlow[®] from MSP.

PUBLISHED BY

 **mathematical sciences publishers**
nonprofit scientific publishing

<http://msp.org/>

© 2025 Mathematical Sciences Publishers

GEOMETRY & TOPOLOGY

Volume 29 Issue 3 (pages 1115–1691) 2025

A cubical model for (∞, n) -categories	1115
TIM CAMPION, KRZYSZTOF KAPULKIN and YUKI MAEHARA	
Rank-one Hilbert geometries	1171
MITUL ISLAM	
Random unitary representations of surface groups, II: The large n limit	1237
MICHAEL MAGEE	
Partial Okounkov bodies and Duistermaat–Heckman measures of non-Archimedean metrics	1283
MINGCHEN XIA	
Global homotopy theory via partially lax limits	1345
SIL LINSKENS, DENIS NARDIN and LUCA POL	
An h-principle for complements of discriminants	1441
ALEXIS AUMONIER	
The motivic lambda algebra and motivic Hopf invariant one problem	1489
WILLIAM BALDERRAMA, DOMINIC LEON CULVER and J D QUIGLEY	
Exotic Dehn twists on sums of two contact 3-manifolds	1571
EDUARDO FERNÁNDEZ and JUAN MUÑOZ-ECHÁNIZ	
On boundedness and moduli spaces of K-stable Calabi–Yau fibrations over curves	1619
KENTA HASHIZUME and MASAFUMI HATTORI	