

UNIFORM EXPONENTIAL MIXING AND RESONANCE FREE REGIONS FOR CONVEX COCOMPACT CONGRUENCE SUBGROUPS OF $SL_2(\mathbb{Z})$

HEE OH AND DALE WINTER

Dedicated to Peter Sarnak on the occasion of his sixty-first birthday.

1. INTRODUCTION

1.1. Uniform exponential mixing. Let $G = SL_2(\mathbb{R})$ and Γ be a non-elementary finitely generated subgroup of $SL_2(\mathbb{Z})$. We will assume that Γ contains the negative identity $-e$ but no other torsion elements. In other words, Γ is the pre-image of a torsion-free subgroup of $PSL_2(\mathbb{Z})$ under the canonical projection $SL_2(\mathbb{Z}) \rightarrow PSL_2(\mathbb{Z})$. For each $q \geq 1$, consider the congruence subgroup of Γ of level q ,

$$\Gamma(q) := \{\gamma \in \Gamma : \gamma \equiv e \pmod{q}\}.$$

For $t \in \mathbb{R}$, let

$$a_t = \begin{pmatrix} e^{t/2} & 0 \\ 0 & e^{-t/2} \end{pmatrix}.$$

As is well known, the right translation action of a_t on $\Gamma \backslash G$ corresponds to the geodesic flow when we identify $\Gamma \backslash G$ with the unit tangent bundle of a hyperbolic surface $\Gamma \backslash \mathbb{H}^2$. We fix a Haar measure dg on G . By abuse of notation, we denote by dg the induced G -invariant measure on $\Gamma(q) \backslash G$. For real-valued functions $\psi_1, \psi_2 \in L^2(\Gamma(q) \backslash G)$, we consider the matrix coefficient

$$\langle a_t \psi_1, \psi_2 \rangle_{\Gamma(q) \backslash G} := \int_{\Gamma(q) \backslash G} \psi_1(ga_t) \psi_2(g) dg.$$

The main aim of this paper is to prove an asymptotic formula (as $t \rightarrow \infty$) with exponential error term for the matrix coefficients $\langle a_t \psi_1, \psi_2 \rangle_{\Gamma(q) \backslash G}$ where the error term is *uniform* for all square free q without small prime divisors.

Denote by $\Lambda(\Gamma)$ the limit set of Γ , that is, the set of all accumulation points of Γ -orbits in the boundary $\partial(\mathbb{H}^2)$ and by $0 < \delta = \delta_\Gamma \leq 1$ the Hausdorff dimension of $\Lambda(\Gamma)$.

The notation C_c^k denotes the space of C^k -functions with compact supports.

Received by the editors April 24, 2015 and, in revised form, August 26, 2015.

2010 *Mathematics Subject Classification*. Primary 37D35, 22E40, 37A25, 37D40, 11F72; Secondary 37F30, 11N45.

The first author was supported in part by NSF Grant 1361673.

Theorem 1.1. *Let $\Gamma < \mathrm{SL}_2(\mathbb{Z})$ be a convex cocompact subgroup; i.e., Γ has no parabolic elements. Then there exist $\eta > 0, C \geq 3$, and $q_0 > 1$ such that for any square free q with $(q, q_0) = 1$ and any $\psi_1, \psi_2 \in C_c^1(\Gamma(q)\backslash G)$, we have*

$$(1.1) \quad e^{(1-\delta)t} \langle a_t \psi_1, \psi_2 \rangle_{\Gamma(q)\backslash G} = \frac{1}{m_q^{\mathrm{BMS}}(\Gamma(q)\backslash G)} m_q^{\mathrm{BR}}(\psi_1) m_q^{\mathrm{BR}^*}(\psi_2) + O(\|\psi_1\|_{C^1} \|\psi_2\|_{C^1} \cdot q^C \cdot e^{-\eta t})$$

as $t \rightarrow +\infty$; here $m_q^{\mathrm{BMS}}, m_q^{\mathrm{BR}}$, and $m_q^{\mathrm{BR}^*}$ denote respectively the Bowen-Margulis-Sullivan measure, the unstable Burger-Roblin measure, and the stable Burger-Roblin measure on $\Gamma(q)\backslash G$ which are chosen compatibly with the choice of dg (see Section 6 for precise definitions).

The implied constant can be chosen uniformly for all C^1 -functions ψ_1, ψ_2 whose supports project to a fixed compact subset of $\Gamma\backslash G$.

If $\Gamma < \mathrm{SL}_2(\mathbb{Z})$ is finitely generated with $\delta > \frac{1}{2}$, then a version of Theorem 1.1 is known by [8] and [12] with a different interpretation of the main term (also see [22], [43], [27]). Therefore the main contribution of Theorem 1.1 lies in the groups Γ with $\delta \leq \frac{1}{2}$; such groups are known to be convex cocompact.

Remark 1.2. (1) Selberg’s celebrated $\frac{3}{16}$ theorem corresponds exactly to this result in the case $\Gamma = \mathrm{SL}_2(\mathbb{Z})$ with the explicit constants $C = 3$ and $\eta(\Gamma) = \frac{1}{4} - \epsilon$. One can therefore regard Theorem 1.1 as yet another generalization of Selberg’s theorem to subgroups of infinite covolume.

- (2) The optimal q^C would be q^3 , which is the growth rate of $[\Gamma : \Gamma(q)]$ in the above error term expressed in C^1 -norms.
- (3) One would expect the results described in this paper to hold without the assumption that q is square free; the missing piece is the ℓ^2 flattening lemma (Lemma 4.7), which is available in the literature only in the case of square free q .
- (4) Theorem 1.1 has an immediate application to counting, equidistribution, and affine sieve; for instance, Theorems 1.7, 1.12, 1.14, 1.16, and 1.17 in [27] are now valid for $\Gamma < \mathrm{SL}_2(\mathbb{Z})$ with $\delta \leq \frac{1}{2}$, with the L^2 -sobolev norms of functions replaced by C^1 -norms; the proofs are verbatim repetition since Theorem 1.1 was the only missing piece in the approach of that paper.

The main term in (1.1) can be related to a Laplace eigenfunction on $\Gamma(q)\backslash \mathbb{H}^2$. Denote by Δ the negative of the Laplacian on \mathbb{H}^2 and $\{\nu_x : x \in \mathbb{H}^2\}$ the Patterson density for Γ . Then $\phi_o(x) := |\nu_x|$ is an eigenfunction of Δ in $C^\infty(\Gamma(q)\backslash \mathbb{H}^2)$ with eigenvalue $\delta(1-\delta)$ [31], and $\phi_o \in L^2(\Gamma(q)\backslash \mathbb{H}^2)$ if and only if $\delta_\Gamma > 1/2$. If we identify $\mathbb{H}^2 = \mathrm{SL}_2(\mathbb{R})/\mathrm{SO}(2)$, and $\psi \in C_c(\Gamma(q)\backslash G)$ is $\mathrm{SO}(2)$ -invariant, then

$$m_q^{\mathrm{BR}}(\psi) = \int_{\Gamma(q)\backslash G} \psi(x) \phi_o(x) dx = m_q^{\mathrm{BR}^*}(\psi).$$

1.2. Uniform resonance free region. When $\delta > \frac{1}{2}$, Bourgain, Gamburd, and Sarnak [8] established a uniform spectral gap for the smallest two Laplace eigenvalues on $L^2(\Gamma(q)\backslash \mathbb{H}^2)$ for all square free $q \in \mathbb{N}$ with no small prime divisors; for some $\epsilon > 0$, there are no eigenvalues between $\delta(1-\delta)$, which is known to be the smallest one, and $\delta(1-\delta) + \epsilon$.

When $\delta \leq \frac{1}{2}$, the L^2 -spectrum of Δ is known to be purely continuous [23], and the relevant spectral quantities are the resonances. The resolvent of the Laplacian

$$R_{\Gamma(q)}(s) := (\Delta - s(1 - s))^{-1} : C_c^\infty(\Gamma(q)\backslash\mathbb{H}^2) \rightarrow C^\infty(\Gamma(q)\backslash\mathbb{H}^2)$$

is holomorphic in the half plane $\Re(s) > \frac{1}{2}$ and has meromorphic continuation to the complex plane \mathbb{C} with poles of finite rank [16] (see also [25]). These poles are called *resonances*. Patterson showed that $s = \delta$ is a resonance of rank 1 and that no other resonances occur in the half-plane $\Re s \geq \delta$ [32]. Naud proved that for some $\epsilon(q) > 0$, the half-plane $\Re s > \delta - \epsilon(q)$ is a resonance free region except at $s = \delta$ [28]. Bourgain, Gamburd, and Sarnak showed that for some $\epsilon > 0$, $\{\Re s > \delta - \epsilon \cdot \min\{1, 1/(\log(1 + |\Im s|))\}\}$ is a resonance free region except for $s = \delta$, for all square free q with no small prime divisors [8]. We will deduce a uniform resonance free half plane from Theorem 1.1 (see Section 6).

Theorem 1.3. *Suppose that $\delta \leq \frac{1}{2}$. There exist $\epsilon > 0$ and $q_0 > 1$ such that for all square free $q \in \mathbb{N}$ with $(q, q_0) = 1$,*

$$\{\Re s > \delta - \epsilon\}$$

is a resonance free region for the resolvent $R_{\Gamma(q)}$ except for a simple pole at $s = \delta$.

Let \mathcal{P}_q denote the set of all primitive closed geodesics in $T^1(\Gamma(q)\backslash\mathbb{H}^2)$ and let $\ell(C)$ denote the length of $C \in \mathcal{P}_q$. The Selberg zeta function given by

$$Z_q(s) := \prod_{k=0}^{\infty} \prod_{C \in \mathcal{P}_q} (1 - e^{-(s+k)\ell(C)})$$

is known to be an entire function when $\Gamma(q)$ is convex cocompact by [21].

Since the resonances of the resolvent of the Laplacian give nontrivial zeros of $Z_q(s)$ by [30], Theorem 1.3 follows from the following.

Theorem 1.4. *There exist $\epsilon > 0$ and $q_0 > 1$ such that for all square free $q \in \mathbb{N}$ with $(q, q_0) = 1$, the Selberg zeta function $Z_q(s)$ is non-vanishing on the set $\{\Re(s) > \delta - \epsilon\}$ except for a simple zero at $s = \delta$.*

1.3. On the proof of Main theorems. Theorem 1.1 is deduced from the following uniform exponential mixing of the Bowen-Margulis-Sullivan measure m_q^{BMS} .

Theorem 1.5. *There exist $\eta > 0, C \geq 3$, and $q_0 > 0$ such that, for all square free $q \in \mathbb{N}$ coprime to q_0 , and for any $\psi_1, \psi_2 \in C_c^1(\Gamma(q)\backslash G)$, we have*

$$\begin{aligned} (1.2) \quad & \int_{\Gamma(q)\backslash G} \psi_1(ga_t)\psi_2(g) dm_q^{\text{BMS}}(g) \\ & = \frac{1}{m_q^{\text{BMS}}(\Gamma(q)\backslash G)} m_q^{\text{BMS}}(\psi_1) \cdot m_q^{\text{BMS}}(\psi_2) + O(\|\psi_1\|_{C^1} \|\psi_2\|_{C^1} \cdot q^C \cdot e^{-\eta t}) \end{aligned}$$

as $t \rightarrow +\infty$, with the implied constant depending only on Γ .

Theorem 1.5 also holds when Γ has a parabolic element by [27]. For a fixed q , Theorem 1.5 was obtained by Stoyanov [39].

We begin by discussing the proof of Theorem 1.5. The first step is to use Markov sections constructed by Ratner [36] and Bowen [5] to build a symbolic model for the a_t -action on the space $\Gamma\backslash G$. The Markov section gives a subshift (Σ, σ) of finite type in an alphabet $\{i_1, \dots, i_k\}$, together with the associated space (Σ^+, σ)

of one sided sequences. Denote by $\tau : \Sigma \rightarrow \mathbb{R}$ the first return time for the flow a_t . The corresponding suspension Σ^τ has a natural flow \mathcal{G}_t , a finite measure μ , and an embedding

$$\zeta : (\Sigma^\tau, \mu, \mathcal{G}_t) \rightarrow (\Gamma \backslash G, m^{\text{BMS}}, a_t)$$

which is an isomorphism of measure theoretic dynamical systems: this is our symbolic model. This framework will be the topic of Section 2.

From the one sided shift we construct, for each $a, b \in \mathbb{R}$, the transfer operator $\mathcal{L}_{ab} : C(\Sigma^+) \rightarrow C(\Sigma^+)$ by

$$(\mathcal{L}_{ab}h)(x) = \sum_{\sigma(y)=x} e^{-(\delta+a-ib)\tau(y)} h(y).$$

Pollicott's observation, later used and refined by many other authors (see [17], [39], [1]), was that the Laplace transform of the correlation function for the system $(\Sigma^\tau, \mu, \mathcal{G}_t)$ can be expressed in terms of transfer operators using the Ruelle-Perron-Frobenius theorem, and that the exponential mixing of $(\Sigma^\tau, \mu, \mathcal{G}_t)$ and hence that of $(\Gamma \backslash G, m^{\text{BMS}}, a_t)$ follows if we prove a uniform spectral bound on \mathcal{L}_{ab} for Hölder observables to be valid on $|a| \leq a_0$ for some $a_0 > 0$.

We write $\text{SL}_2(q)$ for the finite group $\text{SL}_2(\mathbb{Z}/q\mathbb{Z})$. Following this approach, we define congruence transfer operators $\mathcal{M}_{ab,q}$ on the space $C(\Sigma^+, \mathbb{C}^{\text{SL}_2(q)})$ of vector-valued functions for each q satisfying $\text{SL}_2(q) = \Gamma(q) \backslash \Gamma$ (which is the case whenever q does not have small prime divisors): for $x \in \Sigma^+$ and $\gamma \in \text{SL}_2(q)$,

$$(\mathcal{M}_{ab,q}H)(x, \gamma) = \sum_{\sigma(y)=x} e^{-(\delta+a-ib)\tau(y)} H(y, \gamma c^{-1}(y)),$$

where $c : \Sigma^+ \rightarrow \Gamma$ is a cocycle which records the way the a_t -flow moves elements from one fundamental domain to another. The natural extension of Pollicott's idea tells us that *uniform* exponential mixing of $(\Gamma(q) \backslash G, m_q^{\text{BMS}}, a_t)$ will follow if we can establish certain spectral bounds for $\mathcal{M}_{ab,q}$ uniformly for all $|a| \leq a_0$, $b \in \mathbb{R}$, and all q large. This reduction will be carried out in Section 5.

The proof of spectral bounds for transfer operators traditionally falls into two parts. In Section 3 we shall consider the case where $|b|$ is large. The key ideas here are due to Dolgopyat, who gave an ingenious, albeit highly involved, proof of the relevant bounds for \mathcal{L}_{ab} under additional assumptions. We will follow a treatment due to Stoyanov, who carries out the bounds on \mathcal{L}_{ab} for axiom A flows. The bounds follow from an iterative scheme involving Dolgopyat operators, whose construction relies on the highly oscillatory nature of the functions $e^{ib\tau}$ when $|b|$ is large. This oscillation is also sufficient to establish bounds on the congruence transfer operators $\mathcal{M}_{ab,q}$; see Theorem 3.1. Because the oscillation relies only on local non-integrability properties, the bounds we obtain are uniform in q . It is crucial for this argument that the cocycle c is locally constant on an appropriate length scale, so that it does not interfere with the oscillatory argument.

We are left, in Section 4, with the proof of the bounds on $\mathcal{M}_{ab,q}$ for $|a| \leq a_0$ and $|b|$ small. The bounds for \mathcal{L}_{ab} in this region follow immediately from the complex Ruelle-Perron-Frobenius theorem and a compactness argument. Since we require bounds on $\mathcal{M}_{ab,q}$ uniformly in q , however, this compactness argument is not available to us; instead we follow the approach and use the expansion machinery of Bourgain-Gamburd-Sarnak [8].

The expansion approach relies on the idea that $\Gamma(q)\backslash G \sim \mathrm{SL}_2(q) \times \Gamma\backslash G$. Very roughly, the hyperbolic nature of geodesic flow allows us to separate variables and to consider functions that are “independent” of the $\Gamma\backslash G$ component. We are left considering functions on $\mathrm{SL}_2(q)$; for such functions, the right action of the cocycle c , together with the expansion machinery and the ℓ^2 -flattening lemma produce the required decay. One essential estimate in this argument is proved by means of Sullivan’s shadow lemma and the description of the relevant measures in terms of the Patterson-Sullivan density. The memoryless nature of the Markov model for our flow is crucial here, as it allows us to relate these estimates to certain convolutions.

Theorem 1.1 is deduced from Theorem 1.5 by comparing the transverse intersections for the expansion of a horocyclic piece, based on the quasi-product structures of the Haar and the Bowen-Margulis-Sullivan (BMS) measures.

In joint work with Magee [24], we extend a main theorem of [28] uniformly over q as well, which has an application to sieve for orbits of a semigroup as used in the work of Bourgain and Kontorovich on Zaremba’s conjecture [10]. We expect that our methods in this paper generalize to convex cocompact thin subgroups Γ of $\mathrm{SO}(n, 1)$ and moreover to a general rank one group, which we hope to address in a subsequent paper.

Remark. After submission of this paper new arguments have been developed that allow Theorem 1.1 (and hence Theorems 1.3 and 1.4) to be proved without the assumption that q be square free. The key point is to replace the ℓ^2 -flattening lemma with the expansion results of Bourgain and Varjú [14] in the proofs of Propositions 4.18 and 4.20. The new arguments are described in a recent preprint [11] for the setting of [24] and should require only minor modification to apply in our setting.

2. CONGRUENCE TRANSFER OPERATORS

In the whole paper, let $G = \mathrm{SL}_2(\mathbb{R})$ and let $\Gamma < G$ be a non-elementary, convex cocompact subgroup containing the negative identity. We assume that $-e$ is the only torsion element of Γ . If $p : \mathrm{SL}_2(\mathbb{R}) \rightarrow \mathrm{PSL}_2(\mathbb{R})$ is the canonical projection, then $p(\Gamma)$ is a convex cocompact torsion-free subgroup of $\mathrm{PSL}_2(\mathbb{R})$ and we have $\Gamma\backslash \mathrm{SL}_2(\mathbb{R}) = p(\Gamma)\backslash \mathrm{PSL}_2(\mathbb{R})$. Since our results concern the quotient space $\Gamma\backslash G$, we will henceforth abuse notation so that sometimes $G = \mathrm{PSL}_2(\mathbb{R})$ and our Γ is considered as a torsion-free subgroup of $\mathrm{PSL}_2(\mathbb{R})$.

We recall that the limit set $\Lambda(\Gamma)$ is a minimal non-empty closed Γ -invariant subset of the boundary $\partial\mathbb{H}^2$, and its Hausdorff dimension $\delta = \delta_\Gamma$ is equal to the critical exponent of Γ (see [31]).

We denote by $\{\mu_x = \mu_x^{\mathrm{PS}} : x \in \mathbb{H}^2\}$ the Patterson-Sullivan density for Γ ; that is, each μ_x is a finite measure on $\Lambda(\Gamma)$ satisfying

- (1) $\gamma_*\mu_x = \mu_{\gamma x}$ for all $\gamma \in \Gamma$;
- (2) $\frac{d\mu_x}{d\mu_y}(\xi) = e^{\delta\beta_\xi(y,x)}$ for all $x, y \in \mathbb{H}^2$ and $\xi \in \partial(\mathbb{H}^2)$.

Here $\beta_\xi(y, x)$ denotes the Busemann function: $\beta_\xi(y, x) = \lim_{t \rightarrow \infty} d(\xi_t, y) - d(\xi_t, x)$ where ξ_t is a geodesic ray tending to ξ as $t \rightarrow \infty$. Since Γ is convex cocompact, μ_x is simply the δ -dimensional Hausdorff measure on $\Lambda(\Gamma)$ with respect to a spherical metric viewed from x (up to a scaling). See [31] and [41] for references.

Fixing $o \in \mathbb{H}^2$, the map $u \mapsto (u^+, u^-, s = \beta_{u^-}(o, u))$ is a homeomorphism between $\mathrm{T}^1(\mathbb{H}^2)$ and $(\partial(\mathbb{H}^2) \times \partial(\mathbb{H}^2) - \{(\xi, \xi) : \xi \in \partial(\mathbb{H}^2)\}) \times \mathbb{R}$. Using this homeomorphism, and the identification of $\mathrm{PSL}_2(\mathbb{R})$ with $\mathrm{T}^1(\mathbb{H}^2)$, the Bowen-Margulis-Sullivan

measure $\tilde{m}^{\text{BMS}} = \tilde{m}_{\Gamma}^{\text{BMS}}$ on $\text{PSL}_2(\mathbb{R})$ is defined as follows:

$$d\tilde{m}^{\text{BMS}}(u) = e^{\delta\beta_{u^+}(o,u)} e^{\delta\beta_{u^-}(o,u)} d\mu_o^{\text{PS}}(u^+) d\mu_o^{\text{PS}}(u^-) ds.$$

This definition is independent of the choice of $o \in \mathbb{H}^2$, but does depend on Γ .

We denote by m^{BMS} the measure on $\Gamma \backslash G$ induced by \tilde{m}^{BMS} ; it is called the Bowen-Margulis-Sullivan measure on $\Gamma \backslash G$, or the BMS measure for short.

Let $A = \{a_t = \text{diag}(e^{t/2}, e^{-t/2}) : t \in \mathbb{R}\}$. The right translation action of A on $\Gamma \backslash G$ corresponds to the geodesic flow on $\mathbb{T}^1(\Gamma \backslash \mathbb{H}^2)$. It is easy to check that the BMS measure is A -invariant. We choose the left G - and right $\text{SO}_2(\mathbb{R})$ -invariant metric d on G such that $d(e, a_t) = t$.

Let N^+ and N^- be the expanding and contracting horocyclic subgroups for a_t ,

$$(2.1) \quad N^+ = \{n_s^+ := \begin{pmatrix} 1 & 0 \\ s & 1 \end{pmatrix} : s \in \mathbb{R}\} \quad \text{and} \quad N^- = \{n_s^- := \begin{pmatrix} 1 & s \\ 0 & 1 \end{pmatrix} : s \in \mathbb{R}\}.$$

For $\epsilon > 0$, we will denote by N_ϵ^\pm the intersection of the ϵ ball around the identity, $B_\epsilon(e)$, with N^\pm .

We fix a base point $o \in \mathbb{H}^2$ in the convex hull of the limit set $\Lambda(\Gamma)$, and write Ω for the support of the BMS measure. The geodesic flow $a_t : \Omega \rightarrow \Omega$ is known to be mixing for the BMS measure m^{BMS} by Rudolph [37] (see also [2]). Since Γ is convex cocompact, Ω is compact and there is a uniform positive lower bound for the injectivity radii for points on $\Gamma \backslash G$, which we will simply call the injectivity radius of Γ .

2.1. Markov sections. We refer to [18] for basic facts about Markov sections. Let $\alpha > 0$ be a small number. Consider a finite set z_1, \dots, z_k in Ω and choose small compact neighborhoods U_i and S_i of z_i in $z_i N_\alpha^+ \cap \Omega$ and $z_i N_\alpha^- \cap \Omega$ respectively of diameter at most $\alpha/2$. We write $\text{int}^u(U_i)$ for the interior of U_i in the set $z_i N_\alpha^+ \cap \Omega$ and define $\text{int}^s(S_i)$ similarly. We will assume that U_i (respectively S_i) are proper, that is to say, that $U_i = \overline{\text{int}^u(U_i)}$ (respectively $S_i = \overline{\text{int}^s(S_i)}$). For $x \in U_i$ and $y \in S_i$, we write $[x, y]$ for the unique local intersection of xN^- and yN^+A . We write the rectangles as

$$R_i = [U_i, S_i] := \{[x, y] : x \in U_i, y \in S_i\}$$

and denote their interiors by

$$\text{int}(R_i) = [\text{int}^u(U_i), \text{int}^s(S_i)].$$

Note that $U_i = [U_i, z_i] \subset R_i$. The family $\mathcal{R} = \{R_1, \dots, R_k\}$ is called a complete family of size $\alpha > 0$ if

- (1) $\Omega = \cup_1^k R_i a_{[0, \alpha]}$,
- (2) the diameter of each R_i is at most α , and
- (3) for any $i \neq j$, at least one of the sets $R_i \cap R_j a_{[0, \alpha]}$ or $R_j \cap R_i a_{[0, \alpha]}$ is empty.

Set $R = \coprod_i R_i$. Let $\tau : R \rightarrow \mathbb{R}$ denote the first return time and $\mathcal{P} : R \rightarrow R$ the first return map,

$$\tau(x) := \inf\{t > 0 : xa_t \in R\} \quad \text{and} \quad \mathcal{P}(x) := xa_{\tau(x)}.$$

Definition 2.1 (Markov section). *A complete family $\mathcal{R} := \{R_1 \cdots R_k\}$ of size α is called a Markov section for the flow a_t if the following Markov property is satisfied:*

$$\mathcal{P}([\text{int}^u U_i, x]) \supset [\text{Int}^u U_j, \mathcal{P}(x)] \quad \text{and} \quad \mathcal{P}([x, \text{Int}^s S_i]) \subset [\mathcal{P}(x), \text{Int}^s S_j]$$

whenever $x \in \text{int}(R_i) \cap \mathcal{P}^{-1}(\text{int}(R_j))$.

We consider the $k \times k$ matrix

$$\text{Tr}_{lm} = \begin{cases} 1 & \text{if } \text{int}(R_l) \cap \mathcal{P}^{-1}\text{int}(R_m) \neq \emptyset \\ 0 & \text{otherwise,} \end{cases}$$

which we will refer to as the transition matrix. The transition matrix Tr is called topologically mixing if there exists a positive integer N such that all the entries of Tr^N are positive. Ratner [36] and Bowen [5] established the existence of Markov sections of arbitrarily small size; using an argument of Bowen and Ruelle [3] we may further assume that the associated transition matrix is topologically mixing. We now fix such an $\mathcal{R} = \{R_1 = [U_1, S_1], \dots, R_k = [U_k, S_k]\}$ of size α , where $\alpha > 0$ satisfies

$$\alpha < \frac{1}{1000} \cdot \text{Injectivity radius of } \Gamma \backslash G,$$

and for all $|s| < 4\alpha$,

$$(2.2) \quad d(e, n_s^+) \leq |s| \leq 2d(e, n_s^+).$$

Note that $k \geq 2$ as a consequence of the non-elementary property of Γ .

Write

$$U := \coprod_i U_i \quad \text{and} \quad \text{int}(R) = \coprod_i \text{int}(R_i).$$

The projection map along stable leaves

$$\pi_S : R \rightarrow U, \text{ taking } [x, y] \mapsto x$$

will be very important for us at several stages of the argument. We will write $\hat{\sigma}$ for the map

$$\hat{\sigma} := \pi_S \circ \mathcal{P} : U \rightarrow U.$$

Definition 2.2. We define the cores of R and U by

$$\hat{R} = \{x \in R : \mathcal{P}^m x \in \text{int}(R) \text{ for all } m \in \mathbb{Z}\}, \text{ and}$$

$$\hat{U} = \{u \in U : \hat{\sigma}^m u \in \text{int}^u(U) \text{ for all } m \in \mathbb{Z}_{\geq 0}\}.$$

Note that \hat{R} is \mathcal{P} -invariant and that \hat{U} is $\hat{\sigma}$ -invariant. The cores are residual sets (that is, their complements are countable unions of nowhere dense closed sets).

2.2. Symbolic dynamics. We choose Σ to be the space of bi-infinite sequences $x \in \{1, \dots, k\}^{\mathbb{Z}}$ such that $\text{Tr}_{x_l x_{l+1}} = 1$ for all l . Such sequences will be said to be admissible. We denote by Σ^+ the space of one sided admissible sequences

$$\Sigma^+ = \{(x_i)_{i \geq 0} : \text{Tr}_{x_i x_{i+1}} = 1 \text{ for all } i \geq 0\}.$$

We will write $\sigma : \Sigma \rightarrow \Sigma$ for the shift map $(\sigma x)_i = x_{i+1}$. By abuse of notation we will also allow the shift map to act on Σ^+ .

Definition 2.3. For $\theta \in (0, 1)$, we can give a metric d_θ on Σ (resp. on Σ^+) by choosing

$$d_\theta(x, x') = \theta^{\inf\{|j|: x_j \neq x'_j\}}.$$

For a finite admissible sequence $i = (i_0, \dots, i_m)$, we obtain a cylinder of length m ,

$$(2.3) \quad \mathbb{C}[i] := \{u \in \hat{U}_{i_0} : \hat{\sigma}^j(u) \in \text{int}(U_{i_j}) \text{ for all } 0 \leq j \leq m\}.$$

Note that cylinders of length 0 are precisely U_i 's and that cylinders are open subsets of \hat{U} . By a closed cylinder, we mean the closure of some (open) cylinder. We also take this opportunity to introduce embeddings of the symbolic space into the analytic space.

Definition 2.4 (The map $\zeta : \Sigma \rightarrow \hat{R}$). For $x \in \hat{R}$, we obtain a sequence $\omega = \omega(x) \in \Sigma$ by requiring $\mathcal{P}^k x \in R_{\omega_k}$ for all $k \in \mathbb{Z}$. The set $\hat{\Sigma} := \{\omega(x) : x \in \hat{R}\}$ is a residual set in Σ . Using the fact that any distinct pair of geodesics in \mathbb{H}^2 diverge from one another (in either positive time or negative time), one can show that the map $x \mapsto \omega(x)$ is injective. We now define a continuous function $\zeta : \Sigma \rightarrow \hat{R}$ by choosing $\zeta(\omega(x)) = x$ on $\hat{\Sigma}$ and extending continuously to all of Σ .

The restriction $\zeta : \hat{\Sigma} \rightarrow \hat{R}$ is known to be bijective and to intertwine σ and \mathcal{P} .

Definition 2.5 (The map $\zeta^+ : \hat{\Sigma}^+ \rightarrow \hat{U}$). For $u \in \hat{U}$, we obtain a sequence $\omega'(u) \in \Sigma^+$ by requiring $\mathcal{P}^k x \in R_{\omega'_k}$ for all $k \in \mathbb{Z}_{\geq 0}$. We obtain an embedding $\zeta^+ : \Sigma^+ \rightarrow U$ by sending $\omega'(u) \mapsto u'$ where possible and extending continuously. We write $\hat{\Sigma}^+ := (\zeta^+)^{-1}(\hat{U})$. The restriction $\zeta^+ : \hat{\Sigma}^+ \rightarrow \hat{U}$ is known to be bijective and to intertwine σ and $\hat{\sigma}$.

For θ sufficiently close to 1, the embeddings ζ, ζ^+ are Lipschitz. We fix such a θ once and for all. The space $C_\theta(\Sigma)$ (resp. $C_\theta(\Sigma^+)$) of d_θ -Lipschitz functions on Σ (resp. on Σ^+) is a Banach space with the usual Lipschitz norm

$$\|f\|_{d_\theta} = \sup |f| + \sup_{x \neq y} \frac{|f(x) - f(y)|}{d_\theta(x, y)}.$$

Writing $\tilde{\tau} := \tau \circ \zeta \in C_\theta(\Sigma)$, we form the suspension

$$\Sigma^\tau := \Sigma \times \mathbb{R} / (x, t + \tilde{\tau}x) \sim (\sigma x, t).$$

We write $\hat{\Sigma}^\tau$ for the set $(\hat{\Sigma} \times \mathbb{R} / \sim) \subset \Sigma^\tau$. The suspension embeds into the group quotient via the map

$$\zeta^\tau : \hat{\Sigma}^\tau \rightarrow \Gamma \backslash G, \quad (x, s) \rightarrow \zeta(x)a_s,$$

and has an obvious flow $\mathcal{G}_t : (x, s) \mapsto (x, t + s)$. The restriction $\zeta^\tau : \hat{\Sigma}^\tau \rightarrow \Gamma \backslash G$ intertwines \mathcal{G}_t and a_t .

2.3. Pressure and Gibbs measures.

Definition 2.6. For a real valued function $f \in C_\theta(\Sigma)$, called the potential function, we define the pressure to be the supremum

$$Pr_\sigma(f) := \sup_\mu \left(\int_\Sigma f d\mu + \text{entropy}_\mu(\sigma) \right)$$

over all σ -invariant Borel probability measures μ on Σ ; here $\text{entropy}_\mu(\sigma)$ denotes the measure theoretic entropy of σ with respect to μ .

For a given real valued function $f \in C_\theta(\Sigma)$, there is a unique σ -invariant probability measure on Σ that achieves the supremum above, called the equilibrium state for f . We will denote it ν_f . It satisfies $\nu_f(\hat{\Sigma}) = 1$.

To any σ -invariant measure μ on Σ , we can associate a \mathcal{G}_t -invariant measure μ^τ on Σ^τ ; simply take the local product of μ and the Lebesgue measure on \mathbb{R} . Our interest in these equilibrium states is justified in light of the following fact.

Notation 2.7. We will write ν for the $-\delta(\tau \circ \zeta)$ -equilibrium state on Σ . We remark that the pressure $Pr_\sigma(-\delta(\tau \circ \zeta))$ is known to be zero.

Theorem 2.8. Up to a normalization, the measure m^{BMS} on $\Gamma \backslash G$ coincides with the pushforward $\zeta_*^\tau \nu^\tau$.

Proof. Sullivan [42] proved that m^{BMS} is the unique measure of maximal entropy for the a_t action on $\Gamma \backslash G$. On the other hand, $\zeta_*^\tau \nu^\tau$ is also a measure of maximal entropy on $(\Gamma \backslash G, a_t)$ by [18]. The result follows. \square

In particular, this theorem implies that $(\Sigma^\tau, \mathcal{G}_t, \nu^\tau)$ and $(\Gamma \backslash G, a_t, m^{\text{BMS}})$ are measurably isomorphic as dynamical systems via ζ^τ . One simple consequence is the following corollary.

Corollary 2.9. The measures $(\pi \circ \text{vis}^{-1})_* \mu_o^{\text{PS}}$ and $(\pi_S \circ \zeta)_* \nu$ are mutually absolutely continuous on each U_i with a bounded Radon-Nikodym derivative. Here vis denotes the visual map from a lift \tilde{U}_i to $\partial(\mathbb{H}^2)$, and π is the projection $G \rightarrow \Gamma \backslash G$.

By abuse of notation, we use the notation ν for the measure $(\pi_S \circ \zeta)_* \nu$ on U .

2.4. Transfer operators. The identification of Σ^τ and $\Gamma \backslash G$ above allows the use of symbolic dynamics in the study of the BMS measure. In particular, we will use the theory of transfer operators.

Definition 2.10. For $f \in C_\theta(\Sigma^+)$, we obtain a transfer operator $\mathcal{L}_f : C(\Sigma^+) \rightarrow C(\Sigma^+)$ by taking

$$\mathcal{L}_f(h)(u) := \sum_{\sigma(u')=u} e^{f(u')} h(u').$$

A straightforward calculation shows that \mathcal{L}_f preserves $C_\theta(\Sigma)$. The following is a consequence of the Ruelle-Perron-Frobenius theorem together with the well-known theory of Gibbs measures (see [30], [40]).

Theorem 2.11. For each real valued function $f \in C_\theta(\Sigma^+)$, there exist a positive function $\hat{h} \in C_\theta(\Sigma^+)$, a probability measure $\hat{\nu}$ on Σ^+ , and $\epsilon > 0, c > 0$ such that

- $\mathcal{L}_f(\hat{h}) = e^{Pr_\sigma(f)} \hat{h}$;
- the dual operator satisfies $\mathcal{L}_f^* \hat{\nu} = e^{Pr_\sigma(f)} \hat{\nu}$;
- for all $n \in \mathbb{N}$,

$$|e^{-nPr_\sigma(f)} \mathcal{L}_f^n(\psi)(x) - \hat{\nu}(\psi) \hat{h}(x)| \leq c(1 - \epsilon)^n \|\psi\|_{\text{Lip}(d_\theta)},$$

with \hat{h} normalized so that $\hat{\nu}(\hat{h}) = 1$;

- the measure $\hat{h} \hat{\nu}$ is σ -invariant and is the projection of the f -equilibrium state to Σ^+ .

The constants c, ϵ and the Lipschitz norm of \hat{h} can be bounded in terms of the Lipschitz norm of f ; see [40]

Remark. Using the identification of Σ^+ and \hat{U} by ζ^+ , we can regard the transfer operators defined above as operators on $C(\hat{U})$. We can also regard the metric d_θ as a metric on \hat{U} . We will do both of these freely without further comment.

We also define the normalized transfer operators. For $a \in \mathbb{R}$ with $|a|$ sufficiently small, consider the transfer operator $\mathcal{L}_{-(\delta+a)\tau}$ on the space $C_{d_\theta}(U)$. Let $\lambda_a := e^{Pr_\sigma(-(\delta+a)\tau)}$ be the largest eigenvalue, $\hat{\nu}_a$ the probability measure such that

$\mathcal{L}_{-(\delta+a)\tau}^* \hat{\nu}_a = \lambda_a \hat{\nu}_a$, and let h_a be the associated positive eigenfunction, normalized so that $\int h_a d\hat{\nu}_a = 1$. It is known that $\lambda_0 = 1$, and that λ_a and h_a are Lipschitz in a for $|a|$ small. It is also known that for $|a|$ small, each h_a is Lipschitz in the d -metric [40].

Notation 2.12. For functions $f : \hat{U} \rightarrow \mathbb{R}$ and $h : \Sigma^+ \rightarrow \mathbb{R}$, we will write

$$f_n(u) := \sum_{i=0}^{n-1} f(\hat{\sigma}^i u) \quad \text{and} \quad h_n(\omega) = \sum_{i=0}^{n-1} h(\sigma^i \omega).$$

It follows from the fourth part of Theorem 2.11 that there exist $c_1, c_2 > 0$ such that for all $x \in \Sigma^+$ and for all $n \in \mathbb{N}$,

$$(2.4) \quad c_1 e^{-(\delta+a)\tau_n(x)} \lambda_a^{-n} \leq \hat{\nu}_a(\mathbb{C}[x_0, \dots, x_n]) \leq c_2 e^{-(\delta+a)\tau_n(x)} \lambda_a^{-n};$$

moreover, c_1, c_2 can be taken uniformly for $|a| < a_0$ for a fixed $a_0 > 0$. In particular, $\hat{\nu}_a$ is a Gibbs measure for the potential function $-(\delta + a)\tau$.

We consider

$$(2.5) \quad f^{(a)} := -(\delta + a)\tau + \log h_0 - \log h_0 \circ \sigma - \log \lambda_a,$$

and let $\hat{\mathcal{L}}_{ab} := \mathcal{L}_{f^{(a)}+ib\tau}$, i.e.,

$$\hat{\mathcal{L}}_{ab}(h)(u) := \frac{1}{\lambda_a h_0(u)} \sum_{\sigma(u')=u} e^{(-\delta+a-ib)\tau(u')} (h_0 \cdot h)(u'),$$

be the associated transfer operator. Note that $\hat{\mathcal{L}}_{ab}$ preserves the spaces $C_d(\hat{U})$. We remark that the pressure $Pr_\sigma(f^{(a)})$ is zero, so the leading eigenvalue of $\hat{\mathcal{L}}_{a0}$ is 1, with an eigenfunction h_a/h_0 . Since $f^{(0)}$ is cohomologous to $-\delta\tau$, the corresponding equilibrium states coincide.

2.5. Congruence transfer operators and the cocycle c . Let \mathcal{D} be the intersection of the Dirichlet domain for (Γ, o) in \mathbb{H}^2 and the convex hull of $\Lambda(\Gamma)$. For each $R_j \subset \Gamma \backslash G$, we choose a lift $\tilde{R}_j = [\tilde{U}_j, \tilde{S}_j]$ to G so that the projection of \tilde{R}_j to \mathbb{H}^2 intersects $\bar{\mathcal{D}}$ nontrivially. We write $\tilde{R} := \cup \tilde{R}_j$.

Definition 2.13 (Definition of the cocycle $c : R \rightarrow \Gamma$). For $x \in R$ with (unique) lift $\tilde{x} \in \tilde{R}$, we define the cocycle c by requiring that

$$(2.6) \quad \tilde{x} a_{\tau(x)} \in c(x) \tilde{R}.$$

For $n \in \mathbb{N}$ and $x \in U \subset R$, we write

$$c_n(x) := c(x)c(\hat{\sigma}(x)) \cdots c(\hat{\sigma}^{n-1}x).$$

Lemma 2.14. (1) If $x, x' \in R_j \cap \mathcal{P}^{-1}R_l$, then $c(x) = c(x')$.

(2) If x, x' are both contained in some cylinder of length $n \geq 1$, then $c_n(x) = c_n(x')$.

Proof. Let $x_1, x_2 \in R_j \cap \mathcal{P}^{-1}R_l$. If $\tilde{x}_1, \tilde{x}_2 \in \tilde{R}_j$ with $x_j = \Gamma \backslash \Gamma \tilde{x}_j$, then for $\tilde{y}_i := c(x_i)^{-1} \tilde{x}_i a_{\tau(x_i)} \in \tilde{R}_j$, we have

$$\begin{aligned} d(c(x_1)\tilde{y}_1, c(x_2)\tilde{y}_1) &\leq d(\tilde{x}_1, c(x_2)\tilde{y}_1) + \alpha \\ &\leq d(\tilde{x}_2, c(x_2)\tilde{y}_1) + 2\alpha \\ &\leq d(c(x_2)\tilde{y}_2, c(x_2)\tilde{y}_1) + 3\alpha \\ &\leq 4\alpha, \end{aligned}$$

which is less than the injectivity radius of Γ . Thus $c(x_1) = c(x_2)$ as desired. The second statement is now straightforward from the definition of c_n . \square

Let $\Gamma(q)$ be a normal subgroup of Γ of finite index and denote by F_q the finite group $\Gamma(q)\backslash\Gamma$. We would like a compatible family of Markov sections for the dynamical systems $(\Gamma(q)\backslash G, a_t, m^{\text{BMS}})$. The lifts \tilde{R}_l give a natural choice; for $l \in \{1, \dots, k\}$ and $\gamma \in F_q$, we take

$$R_{l,\gamma}^q = \Gamma(q)\gamma\tilde{R}_l \subset \Gamma(q)\backslash G.$$

The collection

$$\mathcal{R}^q := \{R_{l,\gamma}^q : l \in \{1 \dots k\} \text{ and } \gamma \in F_q\}$$

is a Markov section of size α for $(\Gamma(q)\backslash G, a_t)$ as expected. The first return time τ_q and first return map \mathcal{P}_q associated to \mathcal{R}^q are given rather simply in terms of the cocycle c and the corresponding data for \mathcal{R} .

Remark. Let $\pi_q : \Gamma(q)\backslash G \rightarrow \Gamma\backslash G$ be the natural covering map. If $\tilde{x} \in R_{l,\gamma}^q$ and $\pi_q(\tilde{x}) \in R_l \cap \mathcal{P}^{-1}R_m$, then

$$\tau_q(\tilde{x}) = \tau(\pi_q(\tilde{x}))$$

and $\mathcal{P}_q(x)$ is the lift of $\mathcal{P}(\pi(x))$ to $R_{m,\gamma c(\tilde{x})}^q$.

Embedded inside each partition element $R_{l,\gamma}^q$, we have a piece of an unstable leaf. Let \tilde{U}_l be the lift of U_l contained in \tilde{R}_l . Then the subsets

$$U_{l,\gamma}^q := \Gamma(q)\gamma\tilde{U}_l \subset \Gamma(q)\backslash G \text{ and } \hat{U}_{l,\gamma}^q := U_{l,\gamma}^q \cap \pi_q^{-1}(\hat{U})$$

are contained in $R_{l,\gamma}^q$. We then write $\hat{U}^q := \coprod \hat{U}_{l,\gamma}^q$ for the union and $\hat{\sigma}_q : \hat{U}^q \rightarrow \hat{U}^q$ for the natural extension of $\hat{\sigma}$. Just as the partition \mathcal{R} gives rise to a symbolic model of the geodesic flow on $\Gamma\backslash G$, so \mathcal{R}^q provides a model for $\Gamma(q)\backslash G$. In particular, we can identify \hat{U}^q with $\hat{U} \times F_q$ in a natural way; simply send (u, γ) to the image $\gamma\tilde{u}$ where \tilde{u} is the lift of u to \tilde{R} . Note then that $\hat{\sigma}_q$ acts as the map

$$\hat{\sigma}_q(u, \gamma) = (\hat{\sigma}u, \gamma c(u)).$$

For $f_q \in C(\hat{U}^q)$, we may consider the following transfer operators $\mathcal{L}_{f_q,q} : C(\hat{U}^q) \rightarrow C(\hat{U}^q)$ given by

$$\begin{aligned} (\mathcal{L}_{f_q,q}h)(u, \gamma) &:= \sum_{\sigma_q(u', \gamma') = (u, \gamma)} e^{f_q(u', \gamma')} h(u', \gamma') \\ &= \sum_{\sigma(u') = u} e^{f_q(u', \gamma c(u')^{-1})} h(u', \gamma(c(u'))^{-1}). \end{aligned}$$

It will very often be helpful to think of a function $h \in C(\hat{U}^q)$ as a vector valued function $\hat{U} \rightarrow \mathbb{C}^{F_q}$. In the case where $f_q(\omega, \gamma) = f(\omega)$ does not depend on the group element, we can then recover the congruence transfer operator $\mathcal{M}_{f,q} : C(\hat{U}, \mathbb{C}^{F_q}) \rightarrow C(\hat{U}, \mathbb{C}^{F_q})$ given by

$$(\mathcal{M}_{f,q}H)(u) = \sum_{\hat{\sigma}(u') = u} e^{f(u')} H(u')c(u'),$$

where $c(u')$ acts on $H(u') \in \mathbb{C}^{F_q}$ by the right regular action. We will often write $(\mathcal{M}_{f,q}H)(u, \gamma)$ to mean the γ component of $(\mathcal{M}_{f,q}H)(u)$. Most of this paper will be devoted to a study of these congruence transfer operators. The key example

for us will be the normalized congruence transfer operator $\hat{\mathcal{M}}_{ab,q} := \mathcal{M}_{f^{(a)}+ib\tau,q} : C(\hat{U}, \mathbb{C}^{F_q}) \rightarrow C(\hat{U}, \mathbb{C}^{F_q})$,

$$(\hat{\mathcal{M}}_{ab,q}H)(u) = \frac{1}{\lambda_a h_0(u)} \sum_{\hat{\sigma}(u')=u} e^{-(\delta+a-ib)\tau(u)} (h_0 H)(u') \mathbf{c}(u').$$

We then have that for any $n \in \mathbb{N}$,

$$\hat{\mathcal{M}}_{ab,q}^n H(u, \gamma) := \sum_{\hat{\sigma}^n(u')=u} e^{(f_n^{(a)}+ib\tau_n)(u')} H(u', \gamma \mathbf{c}_n^{-1}(u')).$$

The key point will be to establish spectral properties of these congruence transfer operators. To do this we must first establish norms and banach spaces appropriate to the task. We will write $|\cdot|$ for the usual Hermitian norm on \mathbb{C}^{F_q} . For Lipschitz functions $H : \hat{U} \rightarrow \mathbb{C}^{F_q}$, we define the norms

$$(2.7) \quad \|H\|_{1,b} := \sup_{u \in \hat{U}} |H(u)| + \frac{1}{\max(1, |b|)} \sup_{u \neq u'} \frac{|H(u) - H(u')|}{d(u, u')} \text{ and}$$

$$(2.8) \quad \|H\|_2 := \left(\int |H(u)|^2 d\nu(u) \right)^{1/2}.$$

We will sometimes also write $\|\cdot\|_{\text{Lip}(d)} := \|\cdot\|_{1,1}$ for the Lipschitz norm and denote by $C_{\text{Lip}(d)}(\hat{U}, \mathbb{C}^{F_q})$ the space of Lipschitz functions for the norm $\|\cdot\|_{\text{Lip}(d)}$.

Consider the space of functions

$$(2.9) \quad \mathcal{W}(\hat{U}, \mathbb{C}^{F_q}) = \{H \in C_{\text{Lip}(d)}(\hat{U}, \mathbb{C}^{F_q}) : \sum_{\gamma \in F_q} H(u, \gamma) = 0 \text{ for all } u \in \hat{U}\}.$$

We will write $L_0^2(F_q)$ for the space of complex valued functions on F_q that are orthogonal to constants. We can then think of $\mathcal{W}(\hat{U}, \mathbb{C}^{F_q})$ as the space of Lipschitz functions from \hat{U} to $L_0^2(F_q)$.

We are now in a position to state the main technical result of our argument. Suppose that Γ is a (non-elementary) convex cocompact subgroup of $\text{SL}_2(\mathbb{Z})$. We recall the congruence subgroups $\Gamma(q)$ of Γ . Since Γ is Zariski dense in SL_2 , it follows from the strong approximation theorem that there exists $q_0 \geq 1$ such that for all $q \in \mathbb{N}$ with $(q, q_0) = 1$, we have

$$(2.10) \quad \Gamma(q) \backslash \Gamma = G(\mathbb{Z}/q\mathbb{Z}) = \text{SL}_2(q).$$

Theorem 2.15. *There exist $\epsilon > 0, a_0 > 0, C > 0, q'_0 > 1$ such that for all $|a| < a_0, b \in \mathbb{R}$, and for all square free $q \in \mathbb{N}$ with $(q, q_0 q'_0) = 1$, we have*

$$\|\hat{\mathcal{M}}_{ab,q}^m H\|_2 \leq C(1 - \epsilon)^m q^C \|H\|_{1,b}$$

for all $m \in \mathbb{N}$ and all $H \in \mathcal{W}(\hat{U}, \mathbb{C}^{F_q})$.

The next two sections will be focused on the proof of this theorem. In Section 3 we prefer to work with the analytic space \hat{U} and the associated function spaces $C(\hat{U})$, while in Section 4 the symbolic space $\hat{\Sigma}^+$ is preferred. For the most part we can unify these viewpoints through the identification $\zeta : \hat{\Sigma}^+ \rightarrow \hat{U}$; in particular, for those parts of the argument where we consider the transfer operators acting on the $L^2(\nu)$ spaces there is no problem, as the measure theory does not see the precise

geometry of the spaces $\hat{\Sigma}$ and \hat{U} . The one potential difficulty is where we want to use the d -Lipschitz properties of h_a and $f^{(a)}$, which a priori do not follow from the usual statement of the Ruelle-Perron-Frobenius (RPF) Theorem 2.11. This is clarified by [34], which ensures we can proceed as required.

3. DOLGOPYAT OPERATORS AND VECTOR VALUED FUNCTIONS

In this section we aim to prove that Theorem 2.15 holds whenever $|b|$ is sufficiently large.

Theorem 3.1. *There exist $\epsilon > 0, a_0 > 0, b_0 > 0, C > 0$ such that for all $|a| < a_0, |b| > b_0$, and for any normal subgroup $\Gamma(q)$ of Γ of finite index, we have*

$$\|\hat{\mathcal{M}}_{ab,q}^m H\|_2 < C(1 - \epsilon)^m \|H\|_{1,b}$$

for all $m \in \mathbb{N}$ and all $H \in C_{\text{Lip}(d)}(\hat{U}, \mathbb{C}^{F_q})$ for $F_q = \Gamma(q) \backslash \Gamma$.

The strategy here is due to Dolgopyat [17] and uses the construction of so-called Dolgopyat operators. This construction was generalized to axiom A flows by Stoyanov [39], and we will follow his argument. The remaining task is to relate these operators to our vector valued functions. The main reasons we succeed are that (1) the cocycle $c : \hat{U} \rightarrow \Gamma$ is locally constant (Lemma 2.14) and (2) its action on $L^2(F_q)$ is unitary. Both properties are elementary, but they are the critical reasons why our approach works.

Following Stoyanov, we begin by defining a new metric on \hat{U} : for $u, u' \in \hat{U}$, set

$$(3.1) \quad D(u, u') = \inf\{\text{diam}(C) : C \text{ is a cylinder containing } u \text{ and } u'\},$$

where $\text{diam}(C)$ means the diameter of C in the metric d . Note that for all $u, u' \in \hat{U}$,

$$d(u, u') \leq D(u, u').$$

Definition 3.2. *For $E > 0$, we write $K_E(\hat{U})$ for the set of all positive functions $h \in C(\hat{U})$ satisfying*

$$|h(u) - h(u')| \leq Eh(u')D(u, u')$$

for all $u, u' \in \hat{U}$ both contained in \hat{U}_i for some i .

Theorem 3.1 follows from the following technical result as in the works of Dolgopyat and Stoyanov.

Theorem 3.3. *There exist positive constants $N \in \mathbb{N}, E > 1, \epsilon, a_0, b_0$ such that for all a, b with $|a| < a_0, |b| > b_0$ there exist a finite set $\mathcal{J}(b)$ and a family of operators*

$$\mathcal{N}_{J,a} : C(\hat{U}) \rightarrow C(\hat{U}) \text{ for } J \in \mathcal{J}(b)$$

with the properties that

- (1) the operators $\mathcal{N}_{J,a}$ preserve $K_{E|b|}(\hat{U})$;
- (2) we have $\int_{\hat{U}} |\mathcal{N}_{J,a} h|^2 d\nu \leq (1 - \epsilon) \int_{\hat{U}} |h|^2 d\nu$ for all $h \in K_{E|b|}(\hat{U})$;

(3) if $h \in K_{E|b|}(\hat{U})$ and $H \in C(\hat{U}, \mathbb{C}^{F_q})$ satisfy

$$|H(u)| \leq h(u) \text{ and } |H(u) - H(u')| \leq E|b|h(u)D(u, u')$$

for all $u, u' \in \hat{U}$, then there exists $J \in \mathcal{J}(b)$ such that

- $|\hat{\mathcal{M}}_{ab,q}^N H| \leq \mathcal{N}_{J,a} h$;
- for all $u, u' \in \hat{U}$,

$$|\hat{\mathcal{M}}_{ab,q}^N H(u) - \hat{\mathcal{M}}_{ab,q}^N H(u')| \leq E|b|(\mathcal{N}_{J,a} h)(u)D(u, u').$$

The operators $\mathcal{N}_{J,a}$ are called Dolgopyat operators. Before moving on we indicate how to deduce Theorem 3.1 from Theorem 3.3.

Proof that Theorem 3.3 implies Theorem 3.1. Choose $N \in \mathbb{N}$, $\epsilon, |a| < a_0, |b| > b_0, E$, and H as in Theorem 3.3 and set h_0 to be the constant function $\|H\|_{1,b}$. Theorem 3.3 allows us to inductively construct sequences $J_l \in \mathcal{J}(b)$ and $h_l \in K_{E|b|}(\hat{U})$ such that

- (1) $h_{l+1} = \mathcal{N}_{J_l, a} h_l$,
- (2) $|\hat{\mathcal{M}}_{ab,q}^{lN} H(u)| \leq h_l(u)$ pointwise, and
- (3) $\|\hat{\mathcal{M}}_{ab,q}^{lN} H\|_2 \leq \|h_l\|_2 \leq (1 - \epsilon)^l \|H\|_{1,b}$.

Now choose $\epsilon' > 0$ such that $(1 - \epsilon')^N = (1 - \epsilon)$. There is a uniform upper bound, say $R_0 > 1$, on the $L^2(\nu)$ operator norm of $\hat{\mathcal{M}}_{ab,q}$, valid for all b and all $|a| < a_0$. For any $m = lN + r$, with $r < N$, we have

$$\begin{aligned} \|\hat{\mathcal{M}}_{ab,q}^m H\|_2 &= \left(\int_{\hat{U}} |\hat{\mathcal{M}}_{ab,q}^r \hat{\mathcal{M}}_{ab,q}^{lN} H(u)|^2 d\nu \right)^{1/2} \\ &\leq R_0^r \left(\int_{\hat{U}} |\hat{\mathcal{M}}_{ab,q}^{lN} H|^2 d\nu \right)^{1/2} \\ &\leq R_0^r (1 - \epsilon)^l \|H\|_{1,b} \\ &\leq R_0^r (1 - \epsilon')^{lN} \|H\|_{1,b} \\ &\leq R_0^N (1 - \epsilon')^{m-N} \|H\|_{1,b}. \end{aligned}$$

This proves the claim. □

3.1. Notation and constants. We fix notations and constants that will be needed later on. From hyperbolicity properties of the map $\hat{\sigma}$, we obtain constants $c_0 \in (0, 1), \kappa_1 > \kappa > 1$, such that for all $n \in \mathbb{N}$,

$$(3.2) \quad c_0 \kappa^n d(u, u') \leq d(\hat{\sigma}^n u, \hat{\sigma}^n u') \leq c_0^{-1} \kappa_1^n d(u, u')$$

for all $u, u' \in \hat{U}_i$ both contained in some cylinder of length n . Note that this implies a similar estimate for D ,

$$(3.3) \quad c_0 \kappa^n D(u, u') \leq D(\hat{\sigma}^n u, \hat{\sigma}^n u') \leq c_0^{-1} \kappa_1^n D(u, u')$$

for all $u, u' \in \hat{U}_i$ both contained in some cylinder of length n . Fix $0 < a'_0 < 0.1$. The functions τ and h_0 , and hence $f^{(a)}$, are not d -Lipschitz globally, but they are *essentially d -Lipschitz* in the following sense: there exists $0 < T_0 < \infty$ such that

$$(3.4) \quad T_0 \geq \max_{|a| \leq a'_0} \left\{ \|f^{(a)}\|_\infty \right\} + \|\tau\|_\infty,$$

and

$$(3.5) \quad T_0 \geq \frac{|f^{(a)}(u) - f^{(a)}(u')| + |\tau(u) - \tau(u')|}{d(u, u')}$$

for all $|a| < a'_0$ and all u, u' both contained in the same cylinder of length 1. The following lemma follows from the Markov property.

Lemma 3.4. *Suppose that $C[i_0, \dots, i_N]$ is a non-empty cylinder. The map $\hat{\sigma}^n : C[i_0, \dots, i_N] \rightarrow C[i_n, \dots, i_N]$ is a bi-Lipschitz homeomorphism. Moreover, any section v of $\hat{\sigma}^n$ whose image contains $C[i_0, \dots, i_N]$ restricts to a bi-Lipschitz homeomorphism $C[i_n, \dots, i_N] \rightarrow C[i_0, \dots, i_N]$.*

The proof is omitted for brevity. We choose a small $r_0 > 0$ and $z_i \in \hat{U}_i$ such that $2r_0 < \min_i(\text{diam}(U_i))$ and $z_i N_{r_0}^+ \cap \Omega \subset \text{int}^u(U_i)$ for each i (here again Ω denotes the support of the BMS measure). We fix $C_1 > 0$ and $\rho_1 > 0$ to satisfy the following lemma.

Lemma 3.5 ([39, Lemma 3.2]). *There exist $C_1 > 0$ and $\rho_1 > 0$ such that, for any cylinder $C[i]$ of length m , we have*

$$c_0 r_0 \kappa_1^{-m} \leq \text{diam}(C[i]) \leq C_1 \rho_1^m.$$

We also fix $p_0 \in \mathbb{N}$ and $\rho \in (0, 1)$ to satisfy the following proposition.

Proposition 3.6 ([39, Proposition 3.3]). *There exist $p_0 \in \mathbb{N}$ and $\rho \in (0, 1)$ such that, for any n , any cylinder $C[i]$ of length n , and any sub-cylinders $C[i']$, $C[i'']$ of length $(n + 1)$ and $(n + p_0)$, respectively, we have*

$$\text{diam}(C[i'']) \leq \rho \text{diam}(C[i]) \leq \text{diam}(C[i']).$$

Choose also $p_1 > 1$ such that

$$(3.6) \quad 1/4 \leq 1/2 - 2\rho^{p_1-1}.$$

Fact 3.7. *It follows from a property of an equilibrium state and the fact that $Pr_\sigma(-\delta\tau) = 0$ that there is a constant $0 < c_1 < 1$ such that for any $m \in \mathbb{N}$,*

$$c_1 e^{-\delta\tau_m(y)} \leq \nu(C[i]) \leq c_1^{-1} e^{-\delta\tau_m(y)}$$

for any cylinder $C[i]$ of length m and any $y \in C[i]$.

Now we need to recall some consequences of non-joint-integrability of the N^+ , N^- foliations.

Lemma 3.8 (Main Lemma of [39]). *There exist $n_1 \in \mathbb{N}$, $\delta_0 \in (0, 1)$, a non-empty subset $U_0 \subset U_1$ which is a finite union of cylinders of length $n_1 \geq 1$, and $z_0 \in U_0$ such that, setting $\mathcal{U} = \sigma^{n_1}(U_0)$, \mathcal{U} is dense in U and that for any $N > n_1$,*

- (1) *there exist Lipschitz sections $v_1, v_2 : U \rightarrow U$ such that $\sigma^N(v_i(x)) = x$ for all $x \in \mathcal{U}$, and $v_i(\mathcal{U})$ is a finite union of open cylinders of length N ;*
- (2) *$v_1(U) \cap v_2(U) = \emptyset$;*
- (3) *for all $s \in \mathbb{R}$ such that $z_0 n_s^+ \in U_0$, all $0 < |t| < \delta_0$ with $z_0 n_{s+t}^+ \in U_0 \cap \Omega$, we have*

$$\frac{1}{t} |(\tau_N \circ v_2 \circ \hat{\sigma}^{n_1} - \tau_N \circ v_1 \circ \hat{\sigma}^{n_1})(z_0 n_{t+s}^+) - (\tau_N \circ v_2 \circ \hat{\sigma}^{n_1} - \tau_N \circ v_1 \circ \hat{\sigma}^{n_1})(z_0 n_s^+)| \geq \frac{\delta_0}{2}$$

(see (2.1) for other notation).

The next step is to establish certain a priori bounds on the transfer operators. Fix notation as in the previous subsection and choose

$$(3.7) \quad A_0 > 2c_0^{-1} e^{\frac{T_0}{c_0(\kappa-1)}} \max \left\{ 1, \frac{T_0}{\kappa-1} \right\}.$$

Lemma 3.9. *For all $a \in \mathbb{R}$ with $|a| < a'_0$ as in (3.4) and all $|b| > 1$, the following hold:*

- if $h \in K_B(\hat{U})$ for some $B > 0$, then

$$\left| \frac{\hat{\mathcal{L}}_{a_0}^m h(u) - \hat{\mathcal{L}}_{a_0}^m h(u')}{\hat{\mathcal{L}}_{a_0}^m h(u')} \right| \leq A_0 \left[\frac{B}{\kappa^m} + \frac{T_0}{\kappa-1} \right] D(u, u')$$

for all $m \geq 0$ and for all $u, u' \in \hat{U}_i$ for some i ;

- if the functions $0 < h \in C(\hat{U}), H \in C(\hat{U}, \mathbb{C}^{F_q})$ and the constant $B > 0$ are such that

$$|H(v) - H(v')| \leq Bh(v')D(v, v')$$

whenever $v, v' \in \hat{U}_i$ for some i , then for any $m \in \mathbb{N}$ and any $|b| > 1$,

$$|\hat{\mathcal{M}}_{ab,q}^m H(u) - \hat{\mathcal{M}}_{ab,q}^m H(u')| \leq A_0 \left[\frac{B}{\kappa^m} \hat{\mathcal{L}}_{a_0}^m h(u') + |b|(\hat{\mathcal{L}}_{a_0}^m |H|)(u') \right] D(u, u')$$

whenever $u, u' \in \hat{U}_i$ for some i .

Proof. The first part is essentially proved in [39]. We concentrate on the second claim. Let $u, u' \in \hat{U}_i$ for some i and let $m > 0$ be an integer. Given $v \in \hat{U}$ with $\hat{\sigma}^m v = u$, let $C[i_0, \dots, i_m]$ be the cylinder of length m containing v . Note that $i_m = i$ and that $\hat{\sigma}^m C[i_0, \dots, i_m] = \hat{U}_i$ by the Markov property. Moreover, we know that $\hat{\sigma}^m : C[i_0, \dots, i_m] \rightarrow \hat{U}_i$ is a homeomorphism, so there exists $v' = v'(v)$ with $\hat{\sigma}^m v' = u'$. We therefore have

$$d(\hat{\sigma}^j v', \hat{\sigma}^j v) \leq \frac{1}{c_0 \kappa^{m-j}} d(u, u')$$

and so

$$\begin{aligned} |f_m^{(a)}(v) - f_m^{(a)}(v')| &\leq \sum_{j=0}^{m-1} |f^{(a)}(\hat{\sigma}^j v) - f^{(a)}(\hat{\sigma}^j v')| \\ &\leq \sum_{j=0}^{m-1} \|f^{(a)}\|_{\text{Lip}(d)} \frac{D(u, u')}{c_0 \kappa^{m-j}} \\ &\leq \frac{T_0}{c_0(\kappa-1)} D(u, u'). \end{aligned}$$

A similar estimate holds for $|\tau_m(v'(v)) - \tau_m(v)|$ by a similar calculation. In particular,

$$(3.8) \quad e^{f_m^{(a)}(v)} \leq c_0 A_0 e^{f_m^{(a)}(v'(v))},$$

and

$$\begin{aligned} &|e^{(f_m^{(a)} + ib\tau_m)(v) - (f_m^{(a)} + ib\tau_m)(v'(v))} - 1| \\ &\leq e^{|f_m^{(a)}(v) - f_m^{(a)}(v')|} |(f_m^{(a)} + ib\tau_m)(v) - (f_m^{(a)} + ib\tau_m)(v'(v))| \\ (3.9) \quad &\leq |b| A_0 D(u, u'). \end{aligned}$$

Remark. This type of estimate will be used repeatedly for the rest of the paper, often with little comment.

Recall that $c_m(v'(v)) = c_m(v)$ by Lemma 2.14. Using the fact that the diameter of \hat{U}_i is bounded above by 1, we now compute

$$\begin{aligned} & |\hat{\mathcal{M}}_{ab,q}^m H(u) - \hat{\mathcal{M}}_{ab,q}^m H(u')| \\ & \leq \sum_{\hat{\sigma}^m v=u} \left| e^{(f_m^{(a)} - ib\tau_m)(v)} H(v) - e^{(f_m^{(a)} - ib\tau_m)(v'(v))} H(v'(v)) \right| \\ & \leq \sum_{\hat{\sigma}^m v=u} e^{f_m^{(a)}(v)} |H(v) - H(v'(v))| \\ & + \sum_{\hat{\sigma}^m v=u} \left| e^{(f_m^{(a)} - ib\tau_m)(v)} - e^{(f_m^{(a)} - ib\tau_m)(v'(v))} \right| \cdot |H(v'(v))| \\ & \leq \sum_{\hat{\sigma}^m v=u} e^{f_m^{(a)}(v)} Bh(v'(v))D(v, v'(v)) \\ & + \sum_{\hat{\sigma}^m v=u} e^{f_m^{(a)}(v'(v))} \left| e^{(f_m^{(a)} + ib\tau_m)(v) - (f_m^{(a)} + ib\tau_m)(v'(v))} - 1 \right| \cdot |H(v'(v))| \\ & \leq c_0 A_0 B D(v, v'(v)) \sum_{\hat{\sigma}^m v=u} e^{f_m^{(a)}(v'(v))} h(v'(v)) \\ & + |b| A_0 D(u, u') \sum_{\hat{\sigma}^m v=u} e^{f_m^{(a)}(v'(v))} |H(v'(v))| \end{aligned}$$

by (3.8) and (3.9). By definitions and (3.3) this then yields

$$\begin{aligned} & |\hat{\mathcal{M}}_{ab,q}^m H(u) - \hat{\mathcal{M}}_{ab,q}^m H(u')| \\ & \leq \frac{A_0 B D(u, u')}{\kappa^m} \hat{\mathcal{L}}_{a0} h(u') + |b| A_0 D(u, u') \hat{\mathcal{L}}_{a0} |H|(u') \\ & \leq A_0 \left(\frac{B}{\kappa^m} \hat{\mathcal{L}}_{a0} h(u') + |b| \hat{\mathcal{L}}_{a0} |H|(u') \right) D(u, u') \end{aligned}$$

as expected. □

3.2. Construction of Dolgopyat operators. We now recall the construction of Dolgopyat operators. Their definitions rely on a number of constants, which we now fix. The meanings of these constants will become clear throughout the rest of the section. Choose

$$(3.10) \quad E > \max \left\{ \frac{2A_0 T_0}{\kappa - 1}, 4A_0, 1 \right\};$$

$$(3.11) \quad N > n_1 \text{ such that } \kappa^N > \max \left\{ \frac{E}{4c_0}, 6A_0, \frac{512\kappa_1^{n_1} E}{c_0^2 \delta_0 \rho}, \frac{200\kappa_1^{n_1} A_0}{c_0^2} \right\};$$

$$(3.12) \quad \epsilon_1 < \min \left\{ \frac{c_0^2 (\kappa - 1)}{16T_0 \kappa_1^{n_1}}, \frac{c_0 r_0}{\kappa_1^{n_1}}, \frac{\delta_0}{2} \right\};$$

$$(3.13) \quad \mu < \min \left(\frac{1}{4}, \frac{c_0^2 \rho^{p_0, p_1 + 2} \epsilon_1}{4\kappa_1^N}, \frac{c_2^2 \epsilon_1^2}{256} \right);$$

where A_0 is given in (3.7), and other constants are as in Section 3.1. Moreover, set

$$b_0 = 1.$$

For the rest of this subsection, we fix $|b| > b_0$. Let

$$\{C_m := C_m(b)\}$$

be the family of maximal closed cylinders contained in $\overline{U_0}$ (see Lemma 3.8) with $\text{diam}(C_m) \leq \epsilon_1/|b|$. As a consequence of (3.12) and Lemma 3.5 we have the following lemma.

Lemma 3.10. *Each of the cylinders C_m has length at least $n_1 + 1$.*

Corollary 3.11. *Let v_1, v_2 be the sections for $\hat{\sigma}^N$ constructed by Lemma 3.8. If $u, u' \in C_m \cap \hat{U}$, then $c_N(v_i(\hat{\sigma}^{n_1}u)) = c_N(v_i(\hat{\sigma}^{n_1}u'))$ for $i = 1, 2$.*

Proof. Choose $u, u' \in C_m \cap \hat{U}$. They are both contained in some cylinder of length $n_1 + 1$. Thus $\hat{\sigma}^{n_1}u, \hat{\sigma}^{n_1}u'$ are both contained in some cylinder of length 1. But then $v_i(\hat{\sigma}^{n_1}u), v_i(\hat{\sigma}^{n_1}u')$ are both contained in the same cylinder of length N by Lemma 3.8. The result then follows by Lemma 2.14. \square

Notation 3.12. *We set $c_i^{(m)} = c_N(v_i(\hat{\sigma}^{n_1}u)) \in \Gamma$ for any $u \in C_m$; this is well defined by Corollary 3.11.*

Let $\{D_j := D_j(b) : j = 1, \dots, p\}$ be the collection of sub-cylinders of the C_m of length $\text{length}(C_m) + p_0 p_1$. We will say that $D_j, D_{j'}$ are adjacent if they are both contained in the same C_m . We set

$$\begin{aligned} \Xi(b) &:= \{1, 2\} \times \{1, \dots, p(b)\}, \\ \hat{D}_j &:= D_j \cap \hat{U}, \quad Z_j := \overline{\sigma^{n_1}(\hat{D}_j)}, \quad \hat{Z}_j := Z_j \cap \hat{U}, \end{aligned}$$

and

$$X_{i,j} := \overline{v_i(\hat{Z}_j)}, \quad \hat{X}_{i,j} := X_{i,j} \cap \hat{U}$$

for each $i \in \{1, 2\}$ and $j \in \{1, \dots, p\}$. For $J \subset \Xi(b)$, we define $\beta_J : C(\hat{U}) \rightarrow \mathbb{R}$ by

$$\beta_J = 1 - \mu \sum_{(i,j) \in J} w_{i,j},$$

where $w_{i,j}$ is the indicator function of $X_{i,j}$. We recall a number of consequences of the constructions above:

- (1) Each cylinder C_m is contained in some U_n and has diameter at least $\rho\epsilon_1/|b|$; apply Lemma 3.10 and Proposition 3.6.
- (2) $\rho^{p_0 p_1 + 1} \frac{\epsilon_1}{|b|} \leq \text{diam}(D_j) \leq \rho^{p_1} \frac{\epsilon_1}{|b|}$; this follows from the definition of D_j and Proposition 3.6.
- (3) The sections v_i are d -Lipschitz on each \hat{U}_i , with Lipschitz constant no larger than $\frac{1}{c_0 \kappa^N}$; this follows from (3.2).
- (4) The sets $\hat{X}_{i,j}$ are pairwise disjoint cylinders with diameters

$$(3.14) \quad \frac{c_0^2 \epsilon_1 \rho^{p_0 p_1 + 1} \kappa^{n_1}}{\kappa_1^N |b|} \leq \text{diam}(\hat{X}_{i,j}) \leq \frac{\epsilon_1 \rho^{p_1} \kappa_1^{n_1}}{c_0^2 \kappa^N |b|};$$

apply the previous two comments and (3.2).

- (5) The function β_J is D -Lipschitz on \hat{U} with Lipschitz constant

$$(3.15) \quad \frac{\mu \kappa_1^N |b|}{c_0^2 \epsilon_1 \rho^{p_0 p_1 + 1} \kappa^{n_1}};$$

this follows from the previous comment and the definition (3.1) of the metric D .

(6) If $u, u' \in \hat{\sigma}^{n_1}(C_m) \cap \hat{U}$ for some m , then

$$(3.16) \quad D(v_i(u), v_i(u')) \leq \frac{\epsilon_1 \kappa_1^{n_1}}{c_0^2 |b| \kappa^N} \text{ for all } i \in \{1, 2\};$$

see the definition of C_m and (3.2).

(7) If $u', u'' \in \hat{\sigma}^{n_1}(C_m) \cap \hat{U}$, then

$$(3.17) \quad |b| \cdot |(\tau_N(v_2(u')) - \tau_N(v_1(u'))) - (\tau_N(v_2(u'')) - \tau_N(v_1(u'')))| \leq \frac{1}{8};$$

this follows from the definition of C_m , the choice (3.12) of ϵ_1 , and (3.2).

Our next lemma, a simple special case of [39, Lemma 5.9] encapsulates the essential output of non-integrability for our argument. It is deduced from Lemma 3.8.

Lemma 3.13. *For any C_m , there exist $D_{j'}, D_{j''} \subset C_m$ such that*

$$(3.18) \quad |b| \cdot |(\tau_N(v_2(u')) - \tau_N(v_1(u'))) - (\tau_N(v_2(u'')) - \tau_N(v_1(u'')))| \geq \frac{\epsilon_1 \delta_0 \rho}{16}$$

for all $u' \in \hat{Z}_{j'}$ and $u'' \in \hat{Z}_{j''}$.

Proof. Fix m and choose v'_0, v''_0, j', j'' such that $v'_0 \in D_{j'} \subset C_m$, and $v''_0 \in D_{j''} \subset C_m$ with $d(v'_0, v''_0) > \frac{1}{2} \text{diam}(C_m)$. For any $v' \in D_{j'}$ and $v'' \in D_{j''}$, we have

$$d(v', v'') \geq d(v'_0, v''_0) - \text{diam}(D_{j'}) - \text{diam}(D_{j''}) \geq \frac{\epsilon_1 \rho}{|b|} \left(\frac{1}{2} - 2\rho^{p_1-1} \right) \geq \frac{\epsilon_1 \rho}{4}$$

by (3.6). Now we recall z_0 as in Lemma 3.8 and choose $s_1, s_2 \in (-\alpha, \alpha)$ such that $v' = z_0 n_{s_1}^+$ and $v'' = z_0 n_{s_2}^+$. Thus $|s_1 - s_2| \geq \frac{\epsilon_1 \rho}{8}$ by (2.2). On the other hand,

$$\frac{\delta_0}{2} \geq \text{diam}(C_m) \geq d(v', v'') \geq \frac{1}{2} |s_1 - s_2|$$

by (3.12) and (2.2); the result follows by Lemma 3.8 part 3. □

We are finally in a position to give the definition of our Dolgopyat operators. For $|b| > b_0, |a| < a'_0$ and for each $J \subset \Xi(b)$, we define an operator $\mathcal{N}_{J,a} : C(\hat{U}) \rightarrow C(\hat{U})$ by

$$\mathcal{N}_{J,a}(h) := \hat{\mathcal{L}}_{a_0}^N(\beta_J h).$$

3.3. Vector valued transfer operators and Dolgopyat operators. We will now check that appropriate operators $\mathcal{N}_{J,a}$ satisfy the conditions of Theorem 3.3. First choose the subsets $J \in \Xi$ that will be of interest.

Definition 3.14. *A subset $J \subset \Xi(b)$ will be called dense if for every C_m , there exists $(i, j) \in J$ with $D_j \subset C_m$. We write $\mathcal{J}(b)$ for the collection of all dense subsets of $\Xi(b)$.*

The following proves parts 1 and 2 of Theorem 3.3.

Lemma 3.15. *There exist $a_0 \in (0, a'_0)$ and $\epsilon > 0$ such that for any $|a| < a_0$ and $|b| > b_0$, the family of operators $\{\mathcal{N}_{J,a} : J \subset \mathcal{J}(b)\}$ satisfies the following:*

- (1) $\mathcal{N}_{J,a} h \in K_{E|b|}(\hat{U})$ whenever $h \in K_{E|b|}(\hat{U})$;
- (2) $\int_U |\mathcal{N}_{J,a} h|^2 d\nu \leq (1 - \epsilon) \int_U |h|^2 d\nu$ for all $h \in K_{E|b|}(\hat{U})$;

(3) if $H \in C_D(\hat{U}, \mathbb{C}^{F_a})$ and $h \in K_{E|b|}(\hat{U})$ are such that $|H| \leq h$ and

$$|H(v) - H(v')| \leq E|b|h(v')D(v, v'),$$

then

$$|\hat{\mathcal{M}}_{ab,q}^N H(v) - \hat{\mathcal{M}}_{ab,q}^N H(v')| \leq E|b|(\mathcal{N}_{J,a}h)(v')D(v, v'),$$

where N is given as in (3.11).

Proof. The second part is Lemma 5.8 of [39]; although that paper uses a differently normalized transfer operator, the error is at most a factor $\sup_{|a| \leq a_0} \frac{\sup h_a}{\inf h_a}$, which can be absorbed into the decay term for a_0 sufficiently small. The other parts are contained in the same paper for complex valued functions; we include the argument for completeness. Suppose that $h \in K_{E|b|}(\hat{U})$, and that $u, u' \in \hat{U}$. We compute

$$|h\beta_J(u) - h\beta_J(u')| \leq |h(u) - h(u')| + h(u')|\beta_J(u) - \beta_J(u')|.$$

Thus, recalling (3.15),

$$\begin{aligned} |h(u) - h(u')| + h(u')|\beta_J(u) - \beta_J(u')| &\leq |bD(u, u')h(u')(E + \frac{\mu\kappa^N}{c_0\rho^{p_0 p_1 + 1}}) \\ &\leq |bD(u, u')h(u')(E + \frac{\rho}{4}) \end{aligned}$$

by (3.13). It follows that

$$h\beta_J \in K_{(E+\frac{\rho}{4})|b|/(1-\mu)}(\hat{U}).$$

We may now apply Lemma 3.9 above to give

$$\begin{aligned} |\mathcal{N}_{J,a}h(u) - \mathcal{N}_{J,a}h(u')| &= |\mathcal{L}_{a_0}^N(h\beta_J)(u) - \mathcal{L}_{a_0}^N(h\beta_J)(u')| \\ &\leq A_0 \left(\frac{(E + \rho/4)|b|}{\kappa^N(1-\mu)} + \frac{T_0}{\kappa - 1} \right) D(u, u')\mathcal{L}_{a_0}^N(h\beta_J)(u') \\ &\leq A_0 \left(2\frac{E|b|}{\kappa^N} + \frac{T_0}{\kappa - 1} \right) D(u, u')\mathcal{N}_{J,a}h(u') \\ &\leq E|bD(u, u')\mathcal{N}_{J,a}h(u') \end{aligned}$$

as required. The final part also follows as a direct calculation using Lemma 3.9. \square

Our final task for this section is to prove the following key proposition, which completes the proof of Theorem 3.3 by addressing part 3.

Proposition 3.16. *There exists $a_0 > 0$ with the following property. For any h, H as in Theorem 3.3, any $|a| < a_0$, and any $|b| > b_0$ there exists $J \in \mathcal{J}(b)$ such that for all $v \in \hat{U}$,*

$$|\hat{\mathcal{M}}_{ab,q}^N H(v)| \leq \mathcal{N}_{J,a}h(v),$$

where N is given as in (3.11).

We proceed via a series of lemmas.

Lemma 3.17. *For $|b| > b_0$, functions h, H as in Theorem 3.3, we have, for any $(i, j) \in \Xi(b)$,*

(1)

$$\frac{1}{2} \leq \frac{h(v_i(u'))}{h(v_i(u))} \leq 2 \quad \text{for all } u, u' \in \hat{Z}_j;$$

(2) either $|H(v_i(u))| \leq \frac{3}{4}h(v_i(u))$ for all $u \in \hat{Z}_j$ or $|H(v_i(u))| \geq \frac{1}{4}h(v_i(u))$ for all $u \in \hat{Z}_j$.

Proof. For $h \in K_{E|b|}(\hat{U})$ and $u, u' \in \hat{Z}_j$ with $D_j \subset C_m$, we simply calculate

$$\begin{aligned} h(v_i(u')) &\leq h(v_i(u)) + E|b|D(v_i(u'), v_i(u))h(v_i(u)) \\ &\leq h(v_i(u)) \left(1 + E|b|\text{diam}(\hat{X}_{i,j})\right) \\ &\leq 2h(v_i(u)) \end{aligned}$$

by (3.14) and (3.11). The other bound follows by symmetry. The second part of the observation follows by similar calculations, which we shall omit. \square

Definition 3.18. Let $a \in (0, a'_0)$ and choose $|b| > b_0, h, H$ as in Theorem 3.3. For each fixed $C_m = C_m(b)$, recall that

$$c_i^{(m)} = c_N(v_i(\hat{\sigma}^{n_1}(u))) \quad \text{for all } u \in C_m.$$

Define the functions

$$\chi^{(1)}[H, h](u) := \frac{|e^{(f_N^{(a)} + ib\tau_N)(v_1(u))} H(v_1(u))c_1^{(m)} + e^{(f_N^{(a)} + ib\tau_N)(v_2(u))} H(v_2(u))c_2^{(m)}|}{(1 - \mu)e^{f_N^{(a)}(v_1(u))}h(v_1(u)) + e^{f_N^{(a)}(v_2(u))}h(v_2(u))}$$

and

$$\chi^{(2)}[H, h](u) := \frac{|e^{(f_N^{(a)} + ib\tau_N)(v_1(u))} H(v_1(u))c_1^{(m)} + e^{(f_N^{(a)} + ib\tau_N)(v_2(u))} H(v_2(u))c_2^{(m)}|}{e^{f_N^{(a)}(v_1(u))}h(v_1(u)) + (1 - \mu)e^{f_N^{(a)}(v_2(u))}h(v_2(u))}.$$

We claim the following lemma.

Lemma 3.19. For every C_m , there exist $i \in \{1, 2\}$ and $j \in \{1 \cdots p\}$ such that $D_j \subset C_m$ and $\chi^{(i)}[H, h](u) \leq 1$ for all $u \in \hat{Z}_j$.

Proof. Fix m and choose j', j'' as in Lemma 3.13. Consider $\hat{Z}_{j'}$ and $\hat{Z}_{j''}$. If there exist $t \in \{j', j''\}$ and $i \in \{1, 2\}$ such that the first alternative of Lemma 3.17 (2) holds for \hat{Z}_t , then $\chi^{(i)}[H, h](u) \leq 1$ for all $u \in \hat{Z}_t$. So from now on in this proof we assume the converse, i.e., for each $i, |H(v_i(u))| \geq \frac{1}{4}h(v_i(u))$ for all $u \in \hat{Z}_{j'} \cup \hat{Z}_{j''}$.

Consider now $u' \in \hat{Z}_{j'}$ and $u'' \in \hat{Z}_{j''}$. Then the properties of h and H imply

$$\begin{aligned} \frac{|H(v_i(u')) - H(v_i(u''))|}{\min\{|H(v_i(u'))|, |H(v_i(u''))|\}} &\leq \frac{E|b|h(v_i(u'))D(v_i(u'), v_i(u''))}{\min\{|H(v_i(u'))|, |H(v_i(u''))|\}} \\ &\leq 4E|b|D(v_i(u'), v_i(u'')) \\ &< \frac{\epsilon_1 \delta_0 \rho}{128} \text{ by (3.16),} \end{aligned}$$

where we have assumed $|H(v_i(u'))| \leq |H(v_i(u''))|$ without loss of generality.

In particular, this is less than $\frac{1}{2}$. We write $c_2 = \frac{\delta_0 \rho}{16}$. The sine of the angle θ_i between $H(v_i(u'))$ and $H(v_i(u''))$ is therefore at most $\sin \theta_i \leq \frac{c_2 \epsilon_1}{8}$, so

$$(3.19) \quad \theta_i \leq \frac{c_2 \epsilon_1}{4}.$$

We need to use this to show that at least one of the angles

$$\theta(e^{ib\tau_N(v_1(u'))} H(v_1(u'))c_1^{(m)}, e^{ib\tau_N(v_2(u'))} H(v_2(u'))c_2^{(m)})$$

or

$$\theta(e^{ib\tau_N(v_1(u''))}H(v_1(u''))c_1^{(m)}, e^{ib\tau_N(v_2(u''))}H(v_2(u''))c_2^{(m)})$$

is greater than $c_2\epsilon_1/4$. Supposing that the first term is less than $c_2\epsilon_1/4$, we will show that the second term is bigger than $c_2\epsilon_1/4$. Write

$$\phi(w) := b \cdot (\tau_N(v_2(w)) - \tau_N(v_1(w)))$$

and note that

$$c_2\epsilon_1 \leq |\phi(u') - \phi(u'')| \leq \frac{1}{8}$$

for all $u' \in \hat{Z}_j'$ and all $u'' \in \hat{Z}_j''$ by (3.17) and Lemma 3.13.

We compute

$$\begin{aligned} & \theta(e^{ib\tau_N(v_1(u''))}H(v_1(u''))c_1^{(m)}, e^{ib\tau_N(v_2(u''))}H(v_2(u''))c_2^{(m)}) \\ &= \theta(e^{-i\phi(u'')}H(v_1(u''))c_1^{(m)}, H(v_2(u''))c_2^{(m)}) \\ &\geq \theta(e^{-i\phi(u'')}H(v_1(u''))c_1^{(m)}, e^{-i\phi(u')}H(v_1(u''))c_1^{(m)}) \\ &\quad - \theta(e^{-i\phi(u')}H(v_1(u''))c_1^{(m)}, H(v_2(u''))c_2^{(m)}) \\ &\geq |\phi(u') - \phi(u'')| - \theta(e^{-i\phi(u')}H(v_1(u''))c_1^{(m)}, H(v_2(u''))c_2^{(m)}) \\ &\geq c_2\epsilon_1 - c_2\epsilon_1/2 - \theta(e^{-i\phi(u')}H(v_1(u'))c_1^{(m)}, H(v_2(u'))c_2^{(m)}) \\ &\geq c_2\epsilon_1/4 \end{aligned}$$

by (3.19) and the assumption. Write

$$\begin{aligned} v &= e^{(f_N^{(a)} + ib\tau_N(v_1(u'')))H(v_1(u''))c_1^{(m)}} \text{ and} \\ w &= e^{(f_N^{(a)} + ib\tau_N(v_2(u'')))H(v_2(u''))c_2^{(m)}} \end{aligned}$$

so that $|v + w|$ is the numerator of $\chi^{(i)}[H, h](u'')$. Without loss of generality, we assume that $|v| \leq |w|$. We now claim that $\chi^{(1)}[H, h](u'') \leq 1$ for all $u'' \in \hat{Z}_j''$. This now follows from rather simple trigonometry. Since the angle $\tilde{\theta}$ between v and w is at least $c_2\epsilon_1/4$, we have

$$1 + 2\cos\tilde{\theta} \leq 2 + \cos\tilde{\theta} \leq 3 - \frac{c_2^2\epsilon_1^2}{16} \leq 3(1 - \mu)^2 \leq (1 - \mu)^2 + 2(1 - \mu).$$

Thus

$$|v| + 2|v|\cos\tilde{\theta} \leq (1 - \mu)^2|v| + 2(1 - \mu)|v|.$$

Now

$$|v| + 2|w|\cos\tilde{\theta} \leq (1 - \mu)^2|v| + 2(1 - \mu)|w|,$$

and so $(1 - \mu)|v| + |w| \geq |v + w|$ and $\chi^{(1)}[H, h] \leq 1$ on \hat{Z}_j'' as expected. \square

Proof of Proposition 3.16. Choose $h, H, |b| > b_0$ as in the hypotheses of Theorem 3.3 and choose $a_0 \in (0, a'_0)$ to satisfy Lemma 3.15. We choose a subset $J \in \mathcal{J}(b)$ as follows. First include in J all $(1, j) \in \Xi$ such that $\chi^{(1)}[H, h] \leq 1$ on \hat{Z}_j . Then for any $j \in \{1 \cdots p\}$, include $(2, j)$ in J if $(1, j)$ is not already in J and $\chi^{(2)}[H, h] \leq 1$ on \hat{Z}_j . By Lemma 3.19, this subset J is dense (in the sense of Definition 3.14), so that $J \in \mathcal{J}(b)$. We will show that for all $u \in \hat{U}$

$$|\hat{\mathcal{M}}_{ab,q}^N H(u)| \leq \mathcal{N}_{J,a} h(u).$$

Let $u \in \hat{U}$. Suppose first that $u \notin \hat{Z}_j$ for any $(i, j) \in J$; then $\beta_J(v) = 1$ whenever $\hat{\sigma}^N(v) = u$, and the bound follows. Suppose instead that $u \in \hat{Z}_j \subset C_m$ with

$(1, j) \in J$. Then $(2, j) \notin J$ and so $\beta_J(v_1(u)) \geq 1 - \mu$ and $\beta_J(v_2(u)) = 1$. We therefore have $\chi^{(1)}[H, h] \leq 1$ on \hat{Z}_j , so

$$\begin{aligned} |\hat{\mathcal{M}}_{ab,q}^N H(u)| &\leq \sum_{\hat{\sigma}^N v=u, v \neq v_1(u), v_2(u)} e^{f_N^{(a)}(v)} |H(v)| \\ &+ \left| e^{(f_N^{(a)} + ib\tau_N)(v_1(u))} H(v_1(u)) \mathbf{c}_1^{(m)} + e^{(f_N^{(a)} + ib\tau_N)(v_2(u))} H(v_2(u)) \mathbf{c}_2^{(m)} \right| \\ &\leq \sum_{\hat{\sigma}^N v=u, v \neq v_1(u), v_2(u)} e^{f_N^{(a)}(v)} |h(v)| \\ &+ (1 - \mu) e^{f_N^{(a)}(v_1(u))} h(v_1(u)) + e^{f_N^{(a)}(v_2(u))} h(v_2(u)) \\ &\leq \mathcal{N}_{J,a} h(u). \end{aligned}$$

The case $u \in \hat{Z}_j$ with $(2, j) \in J$ is similar. This finishes the proof. □

Together with Lemma 3.15, this completes the proof of Theorem 3.3.

4. THE EXPANSION MACHINERY

4.1. Some reductions. In this section we assume that Γ is a convex cocompact subgroup in $\text{SL}_2(\mathbb{Z})$ and that q_0 is as in (2.10). Let $b_0 > 0$ be as in Theorem 3.1. The main aim of this section is to prove the following theorem.

Theorem 4.1. *There exist $\epsilon \in (0, 1), a_0 > 0, C > 1, q'_0 > 1$ such that for all $|a| < a_0, |b| \leq b_0$, and all square free $q \geq 1$ with $(q, q_0 q'_0) = 1$, we have*

$$\|\hat{\mathcal{M}}_{ab,q}^m H\|_2 < C(1 - \epsilon)^m q^C \|H\|_{\text{Lip}(d)}$$

for all $m \in \mathbb{N}$ and all $H \in \mathcal{W}(\hat{U}, \mathbb{C}^{\text{SL}_2(q)})$; see (2.9) for notation.

Since the $\text{Lip}(d)$ norm and the $\|\cdot\|_{1,b}$ norm are equivalent for all $|b| \leq b_0$, this theorem and Theorem 3.1 imply Theorem 2.15.

The key ingredient of the proof of Theorem 4.1 is the expander technology, introduced in this context by Bourgain, Gamburd, and Sarnak [8], from which we draw heavily throughout this section. The idea of the expansion machinery is that random walks on the Cayley graphs of $\text{SL}_2(q)$ have good spectral properties. We do not have a random walk in the usual sense, but the randomness inherent in the Gibbs measure provides the same effect.

We recall the sequence spaces Σ^+, Σ , the shift map σ , and the embedding $\zeta : \Sigma^+ \rightarrow \hat{U}$.

Notation 4.2. *For any function $H \in C(\hat{U}, \mathbb{C}^{\text{SL}_2(q)})$, we will denote $\tilde{H} = H \circ \zeta : \Sigma^+ \rightarrow \mathbb{C}^{\text{SL}_2(q)}$. Similarly $\tilde{\tau}$ will denote $\tau \circ \zeta$.*

We recall the constant $\theta \in (0, 1)$ chosen sufficiently close to one (see Section 2.2) and the metric d_θ on Σ (resp. on Σ^+). Write

$$\|\tilde{H}\|_\infty := \sup_{\omega \in \hat{\Sigma}^+} |\tilde{H}(\omega)|$$

and

$$\text{Lip}_{d_\theta}(\tilde{H}) := \sup_{\omega \neq \omega' \in \hat{\Sigma}^+} \frac{|\tilde{H}(\omega) - \tilde{H}(\omega')|}{d_\theta(\omega, \omega')}$$

which is the minimal Lipschitz constant of \tilde{H} . We also write

$$\|\tilde{H}\|_{d_\theta} := \|\tilde{H}\|_\infty + \text{Lip}_{d_\theta}(\tilde{H}).$$

We fix the following constant for later convenience:

$$(4.1) \quad \eta_\theta := \frac{\text{Lip}_{d_\theta}(\tau) + \sup_{|a| < 1} \text{Lip}_{d_\theta}(f^{(a)})}{1 - \theta}.$$

Rather than proving Theorem 4.1 directly, we will instead start by describing some reductions to a simpler form. For $q'|q$, we define $\hat{E}_{q'}^q \subset L_0^2(\text{SL}_2(q))$ to be the space of functions invariant under the left action of $\Gamma(q')$. We may then write

$$E_{q'}^q := \hat{E}_{q'}^q \cap \left(\bigoplus_{q' \neq q''|q} \hat{E}_{q''}^q \right)^\perp.$$

We think of $E_{q'}^q$ as the space of *new* functions at the level q' . We can then define $\tilde{E}_{q'}^q$ as the subspace of functions H in $\mathcal{W}(\hat{U}, \mathbb{C}^{\text{SL}_2(q)})$ with $H(u, \cdot) \in E_{q'}^q$ for all u . We recall the orthogonal decomposition

$$L_0^2(\text{SL}_2(q)) = \bigoplus_{1 \neq q'|q} E_{q'}^q$$

and the induced direct sum decomposition

$$\mathcal{W}(\hat{U}, \mathbb{C}^{\text{SL}_2(q)}) = \bigoplus_{1 \neq q'|q} \tilde{E}_{q'}^q.$$

Write

$$e_{q,q'} : \mathcal{W}(\hat{U}, \mathbb{C}^{\text{SL}_2(q)}) \rightarrow \tilde{E}_{q'}^q$$

for the projection operator, and note that $e_{q,q'}$ is norm decreasing for both the $\|\cdot\|_{\text{Lip}(d)}$ norm and the $\|\cdot\|_2$ norm.

Remark. The projection operators commute with the congruence transfer operators: we have

$$e_{q,q'} \circ \hat{\mathcal{M}}_{ab,q} = \hat{\mathcal{M}}_{ab,q} \circ e_{q,q'}$$

for any $q'|q$.

The first reduction is that we only need to consider functions in \tilde{E}_q^q .

Theorem 4.3. *There exist $\epsilon \in (0, 1)$, $a_0 > 0$, $C > 1$, $q_1 > 1$ such that for all $|a| < a_0$, $|b| < b_0$, and $q \geq q_1$ square free with $(q, q_0) = 1$, we have*

$$(4.2) \quad \|\hat{\mathcal{M}}_{ab,q}^m H\|_2 < C(1 - \epsilon)^m q^C \|H\|_{\text{Lip}(d)}$$

for all $m \in \mathbb{N}$ and all $H \in \tilde{E}_q^q$.

Proof that Theorem 4.3 implies Theorem 4.1. Set q'_0 to be the product of all primes less than or equal to q_1 . We will first explain how to deduce Theorem 4.1 from Theorem 4.3. Fix $\epsilon, a_0, b_0, C, q_1, q_0$ as in Theorem 4.3. Fix also q square free such that $(q, q_0 q'_0) = 1$. For $q'|q$, we consider the projection maps

$$\text{proj}_{q,q'} : E_{q'}^q \rightarrow E_{q'}^{q'}$$

by choosing $(\text{proj}_{q,q'} F)(\gamma) = F(\tilde{\gamma})$, where $\tilde{\gamma}$ is any pre-image of γ under the natural projection map $\text{SL}_2(q) \rightarrow \text{SL}_2(q')$. By abuse of notation we will also write $\text{proj}_{q,q'}$ for the induced maps $\tilde{E}_{q'}^q \rightarrow \tilde{E}_{q'}^{q'}$. Write $\spadesuit_{q,q'} := \frac{\#\text{SL}_2(q')}{\#\text{SL}_2(q)}$. We note that

$$\|(\text{proj}_{q,q'} H)\|_{\text{Lip}(d)} = \sqrt{\spadesuit_{q,q'}} \|H\|_{\text{Lip}(d)},$$

that

$$\|(\text{proj}_{q,q'} H)\|_2 = \sqrt{\spadesuit_{q,q'}} \|H\|_2,$$

and that

$$\hat{\mathcal{M}}_{ab,q'} \circ \text{proj}_{q,q'} = \text{proj}_{q,q'} \circ \hat{\mathcal{M}}_{ab,q}.$$

Now consider $H \in \mathcal{W}(\hat{U}, \mathbb{C}^{\text{SL}_2(q)})$. We calculate, for $|a| < \min(a_0, a'_0)$ and $|b| \leq b_0$,

$$\begin{aligned} \|\hat{\mathcal{M}}_{ab,q}^m H\|_2^2 &= \sum_{1 \neq q' | q} \|e_{q,q'} \hat{\mathcal{M}}_{ab,q}^m H\|_2^2 \\ &= \sum_{1 \neq q' | q} \|\hat{\mathcal{M}}_{ab,q}^m(e_{q,q'} H)\|_2^2 \\ &= \sum_{1 \neq q' | q} \spadesuit_{q,q'} \|\text{proj}_{q,q'}(\hat{\mathcal{M}}_{ab,q}^m(e_{q,q'} H))\|_2^2 \\ &= \sum_{1 \neq q' | q} \spadesuit_{q,q'} \|\hat{\mathcal{M}}_{ab,q'}^m(\text{proj}_{q,q'}(e_{q,q'} H))\|_2^2. \end{aligned}$$

Applying Theorem 4.3, we obtain

$$\begin{aligned} \|\hat{\mathcal{M}}_{ab,q}^m H\|_2^2 &\leq C^2(1 - \epsilon)^{2m} (q')^{2C} \sum_{1 \neq q' | q} \spadesuit_{q,q'} \|\text{proj}_{q,q'}(e_{q,q'} H)\|_{\text{Lip}(d)}^2 \\ &\leq C^2(1 - \epsilon)^{2m} \sum_{1 \neq q' | q, q' \geq q_1} (q')^{2C} \|e_{q,q'} H\|_{\text{Lip}(d)}^2 \\ &\leq (C'')^2(1 - \epsilon'')^{2m} q^{2C''+1} \|H\|_{\text{Lip}(d)}^2 \end{aligned}$$

as expected. □

The most convenient formulation to prove will be the following theorem.

Theorem 4.4. *There exist $\kappa > 0, a_0 > 0, q_1 > 0$ such that*

$$\|\hat{\mathcal{M}}_{ab,q}^{ln_q} \tilde{H}\|_2 \leq q^{-l\kappa} \|\tilde{H}\|_{d_\theta}$$

for all $|a| < a_0, |b| \leq b_0, l \in \mathbb{N}$, all $q > q_1$ square free and coprime to q_0 , and all $H \in \tilde{E}_q^q$; here n_q denotes the integer part of $\log q$.

Proof that Theorem 4.4 implies Theorem 4.3. Choose $a_0 > 0$ small enough that $|\log \lambda_a| \leq \epsilon \leq \min(\kappa/2, 1)$ for all $|a| < a_0$. Set

$$C := \max \left\{ \log \left(\sup_{|a| \leq a_0, b \in \mathbb{R}} \|\hat{\mathcal{M}}_{ab,q}\|_2 \right), 0 \right\}.$$

Then for all $0 \leq r < n_q$, we have $\|\hat{\mathcal{M}}_{ab,q}\|_2^r \leq q^C$.

For any $m \in \mathbb{N}$, we write $m = ln_q + r$, with $0 \leq r < n_q$. Thus Theorem 4.4 yields

$$\begin{aligned} \|\hat{\mathcal{M}}_{ab,q}^m \tilde{H}\|_2 &\leq \|\hat{\mathcal{M}}_{ab,q}\|_2^r \cdot \|\hat{\mathcal{M}}_{ab,q}^{ln_q} \tilde{H}\|_2 \\ &\leq q^C q^{-l\kappa} \|\tilde{H}\|_{d_\theta} \\ &\leq q^C e^{-ln_q \epsilon} \|\tilde{H}\|_{d_\theta} \\ &\leq q^{C+1} e^{-\epsilon m} \|\tilde{H}\|_{d_\theta}, \end{aligned}$$

as desired. □

4.2. The ℓ^2 -flattening lemma. The rest of this section is devoted to a proof of Theorem 4.4. The key ingredient is a version of the ℓ^2 -flattening lemma 4.7 of Bourgain-Gamburd-Sarnak [8, Lemma 7.2]. For the rest of this section we will assume that

$$q \text{ is square free and coprime to } q_0 \text{ (as in (2.10)).}$$

Definition 4.5. For a complex valued measure μ on $SL_2(q)$ and $q'|q$, we define $|||\pi_{q'}(\mu)|||_\infty$ to be the maximum weight of $|\mu|$ over all cosets of subgroups of $SL_2(q')$ that have proper projection in each divisor of q' .

Notation 4.6. For a function ϕ and a measure μ on $SL_2(q)$, we denote the convolution by

$$\mu * \phi(g) = \sum_{\gamma \in SL_2(q)} \mu(\gamma)\phi(g\gamma^{-1}).$$

Lemma 4.7 ([7], [8]). Given $\kappa > 0$ there exist $\kappa' > 0$ and $C > 0$ such that if μ satisfies $||\mu||_1 \leq B$ and

$$|||\pi_{q'}(\mu)|||_\infty < q^{-\kappa} B \text{ for all } q'|q, q' > q^{1/10} \text{ for some } B > 0,$$

then for each q and $\phi \in E_q^g$,

$$||\mu * \phi||_2 \leq Cq^{-\kappa'} B ||\phi||_2.$$

4.3. Measure estimates on cylinders. Before we can apply the expansion machinery we must first establish certain a priori measure estimates on cylinders; that will be the topic of this subsection. We define the following notation.

Notation 4.8. • For $x = (x_1, x_2, \dots) \in \Sigma^+$ and a sequence i_1, \dots, i_n of symbols, we denote the concatenation by

$$(i_n, \dots, i_1, x) = (i_n, \dots, i_1, x_1, x_2, \dots);$$

• For a function f on Σ^+ and $x \in \Sigma^+$, we set

$$f_n(x) := f(x) + f(\sigma(x)) + \dots + f(\sigma^{n-1}(x));$$

• For $x = (x_i) \in \Sigma^+$, put $c(x) = c(\zeta(x))$, and

$$c_n(x) := c(\zeta(x))c(\zeta(\sigma x)) \dots c(\zeta(\sigma^{n-1}x)) \in \Gamma.$$

Lemma 4.9. For sequences $x, y \in \hat{\Sigma}^+$ with $x_i = y_i$ for $i = 0 \dots k$ for some $k \geq 1$, we have $c_k(x) = c_k(y)$.

Proof. This is a straightforward consequence of Lemma 2.14. □

We may therefore write $c(x) = c(x_0, x_1)$, and more generally, for $n \geq 2$, $c_n(x)$ is the product $c(x_0, x_1) \dots c(x_{n-1}, x_n)$.

Notation 4.10. In the rest of the section, the notation $\sum_{i_1, \dots, i_\ell}$ means the sum taken over all sequences (i_1, \dots, i_ℓ) such that any concatenation following the sum sign is admissible.

Lemma 4.11. *There exist $0 < a_0 < 1$ and $c > 1$ such that for all $|a| < a_0$, $x \in \hat{\Sigma}^+$ and for all $n \in \mathbb{N}$,*

$$\sum_{i_n, \dots, i_1} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \leq c.$$

Proof. This follows easily from (2.4). □

We recall that for $x, y \in \mathbb{H}^2$ and $r > 0$, the shadow $O_r(x, y)$ is defined to be the set of all points $\xi \in \partial(\mathbb{H}^2)$ such that the geodesic ray from x to ξ intersects the ball $B_r(y)$ non-trivially. We need the following: recall the Patterson-Sullivan density $\{\mu_x^{\text{PS}} : x \in \mathbb{H}^2\}$ for Γ .

Lemma 4.12 (Sullivan’s shadow lemma [41]). *Let $x \in \mathbb{H}^2$. There exists $r_0 = r_0(x) > 1$ such that for all $r > r_0$, there exists $c > 1$ such that for all $\gamma \in \Gamma$,*

$$c^{-1}e^{-\delta d(x, \gamma x)} \leq \mu_x^{\text{PS}}(O_r(x, \gamma x)) \leq ce^{-\delta d(x, \gamma x)}.$$

Lemma 4.13. *There exists $c' > 1$ such that for any $x \in \hat{\Sigma}^+$, $\gamma \in \Gamma$, any $m \in \mathbb{N}$, and any fixed i_{m+1} , we have*

$$\sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, i_m, \dots, i_1, x) = \gamma} e^{f_m^{(0)}(i_m, \dots, i_1, x)} \leq c' \cdot e^{-\delta m \inf(\tau)}.$$

Proof. Recalling the definition (2.5) of $f^{(0)}$ in terms of τ , we calculate

$$\begin{aligned} & \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, i_m, \dots, i_1, x) = \gamma} e^{f_m^{(0)}(i_m, \dots, i_1, x)} \\ & \leq \frac{\sup(h_0)}{\inf(h_0)} \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, i_m, \dots, i_1, x) = \gamma} e^{-\delta \tau_m(i_m, \dots, i_1, x)} \\ & \leq e^{\delta \sup(\tau)} \frac{\sup(h_0)}{\inf(h_0)} \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, i_m, \dots, i_1, x) = \gamma} e^{-\delta \tau_{m+1}(i_{m+1}, \dots, i_1, x)} \\ & \leq c'_1 \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, i_m, \dots, i_1, x) = \gamma} \nu(\mathbb{C}[i_{m+1}, i_m, \dots, i_1]) \text{ (see Fact 3.7),} \end{aligned}$$

where $c'_1 = c_1 e^{\delta \sup(\tau)} \frac{\sup(h_0)}{\inf(h_0)}$.

Recall that \mathcal{D} denotes the intersection of the Dirichlet domain for (Γ, o) with the convex core of Γ , and the lifts \tilde{U}_i of U_i chosen to intersect \mathcal{D} . Recall also the projection map π from G to $\Gamma \backslash G$. It is a consequence of the definition of c (2.6) that

$$(4.3) \quad \bigcup_{i_m, \dots, i_1: c_{m+1}(i_{m+1}, \dots, i_1, x) = \gamma} \mathbb{C}[i_{m+1}, i_m, \dots, i_1] \subset \pi \left(\{ \tilde{u} \in \tilde{U}_{i_{m+1}} : d(\tilde{u}a_{\tau_{m+1}(\tilde{u})}, \gamma o) < R_1 \} \right),$$

where R_1 denotes thrice the size of the Markov section plus twice the diameter of \mathcal{D} plus the constant $r_0(o)$ defined in Lemma 4.12.

Case 1. *If $d(\tilde{u}, \gamma \tilde{u}) < (m + 1) \inf(\tau) - R_1$, then*

$$\{ \tilde{u} \in \tilde{U}_{i_{m+1}} : ua_{\tau_{m+1}(\tilde{u})} \in B_{R_1}(\gamma \tilde{u}) \} = \emptyset,$$

and the claim follows.

Case 2. We now assume that $d(\tilde{u}, \gamma\tilde{u}) \geq (m + 1)\inf(\tau) - R_1$. Then $d(o, \gamma o) \geq (m + 1)\inf(\tau) - 2R_1$. A straightforward argument in hyperbolic geometry yields

$$\begin{aligned} \{\tilde{u} \in \tilde{U}_{i_{m+1}} : d(\tilde{u}a_{\tau_{m+1}(\tilde{u})}, \gamma o) < R_1\} &\subset \{\tilde{u} \in \tilde{U}_{i_{m+1}} : \text{vis}(\tilde{u}) \in O_{R_1}(\tilde{u}, \gamma o)\} \\ &\subset \{\tilde{u} \in \tilde{U}_{i_{m+1}} : \text{vis}(\tilde{u}) \in O_{2R_1}(o, \gamma o)\}. \end{aligned}$$

Applying Corollary 2.9 and Lemma 4.12, we obtain

$$\begin{aligned} &\sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, i_m, \dots, i_1, x) = \gamma} e^{-\delta\tau_m(i_m, \dots, i_1, x)} \\ &\leq c'_1 \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, i_m, \dots, i_1, x) = \gamma} \nu(\mathbb{C}[i_{m+1}, i_m, \dots, i_1]) \\ &\leq c'_1 \nu(\{\pi\tilde{u} : \tilde{u} \in \tilde{U}_{i_{m+1}} \text{ and } \text{vis}(\tilde{u}) \in O_{2R_1}(o, \gamma o)\}) \\ &\leq c'_2 \mu_o^{\text{PS}}(O_{2R_1}(o, \gamma o)) \\ &\leq c'_3 e^{-\delta\inf(\tau)m} \end{aligned}$$

as required. □

4.4. Decay estimates for convolutions. We will now use the cylinder estimates and the expansion machinery to provide technical estimates on the L^2 norm of certain convolutions. This is the last preparatory step before we begin the proof of Theorem 4.4 in earnest.

We observe that the set

$$\mathcal{S} := \{\pm c(x), \pm c(x)^{-1} \in \Gamma : x \in \Sigma^+\}$$

is a finite symmetric subset of Γ .

Lemma 4.14. *The set \mathcal{S} generates Γ .*

Proof. By our assumption, the projection $p(\Gamma)$ is a torsion-free convex cocompact subgroup of $\text{PSL}_2(\mathbb{R})$, and hence it is a classical Schottky group by [13]. Therefore the Dirichlet domain D for $(p(\Gamma), o)$ is the common exterior of finitely many disks D_i , $i = 1, \dots, 2\ell$, which meets $\partial(\mathbb{H}^2)$ perpendicularly and whose closures are pairwise disjoint. It is now clear from the definition of the cocycle that $\{c(x)\}$ contains all $\gamma \in p(\Gamma)$ such that $\overline{D} \cap \gamma(\overline{D})$ is non-empty, and hence contains a generating set for $p(\Gamma)$. □

Notation 4.15. *For $m \in \mathbb{N}$, we write $\mathcal{B}_m(e) \subset \Gamma$ for the ball of radius m around the identity e in the word metric defined by \mathcal{S} .*

Notation 4.16. *We write n_q for the integer part of $\log q$. There exists $d_0 > 3$ such that for any $m_q \leq \frac{1}{d_0} \log q$, the ball $\mathcal{B}_{m_q}(e)$ injects to $\text{SL}_2(q')$ whenever $q' | q, q' > q^{1/10}$. For each q we fix a choice*

$$\frac{n_q}{2d_0} < m_q < \frac{n_q}{d_0}$$

and denote $r_q = n_q - m_q$, so that

$$\frac{(d_0 - 1)n_q}{d_0} < r_q < \frac{(d_0 - 1/2)n_q}{d_0}.$$

Notation 4.17. For each element i of the alphabet defining Σ , we choose an element $\omega(i) \in \hat{\Sigma}^+$ such that the concatenation $(i, \omega(i))$ is admissible.

For $\gamma \in \text{SL}_2(q)$, we write δ_γ for the dirac measure at γ . Given real numbers a, b , an element $x \in \hat{\Sigma}^+$, and an admissible sequence $(i_{n_q}, \dots, i_{m_q+1})$, we define a complex valued measure $\mu_{x, (i_{n_q}, \dots, i_{m_q+1})}^{a, b}$ on $\text{SL}_2(q)$ by

$$(4.4) \quad \mu_{x, (i_{n_q}, \dots, i_{m_q+1})}^{a, b} := \sum_{i_1, \dots, i_{m_q}} e^{(f_{n_q}^{(a)} + ib\tau_{n_q})(i_{n_q}, \dots, i_1, x)} \delta_{c_{m_q+1}(i_{m_q+1}, \dots, i_1, x)}.$$

That is, for $\gamma \in \text{SL}_2(q)$, the value $\mu_{x, (i_{n_q}, \dots, i_{m_q+1})}^{a, b}(\gamma)$ is given by the sum

$$\sum e^{(f_{n_q}^{(a)} + ib\tau_{n_q})(i_{n_q}, \dots, i_1, x)}$$

over all indices (i_1, \dots, i_{m_q}) satisfying $c(i_{m_q+1}, i_{m_q}) \cdots c(i_2, i_1)c(i_1, x_0) = \gamma$.

Our first goal in this subsection is to prove the following proposition, which is essential to prove bounds on the supremum norm $\|\hat{\mathcal{M}}_{ab, q} \tilde{H}\|_\infty$.

Proposition 4.18. Let

$$\mu := \mu_{x, (i_{n_q}, \dots, i_{m_q+1})}^{a, b}$$

and

$$B = B_{i_{n_q}, \dots, i_{m_q+1}}^a := ce^{f_{r_q}^{(a)}(i_{n_q}, \dots, i_{m_q+1}, \omega)} e^{\eta_\theta},$$

with $c > 1$ the constant from Lemma 4.11, $\omega = \omega(i_{m_q+1})$, and η_θ as in (4.1). There exist constants $\kappa > 0, a_0 > 0, q_1 > 1, C > 1$ such that for any $x \in \hat{\Sigma}^+$, for any square free $q > q_1$, $|a| < a_0, |b| \leq b_0$, and for all $\phi \in E_q^a$,

$$\|\mu * \phi\|_2 \leq BCq^{-\kappa} \|\phi\|_2.$$

The constants C, κ may be chosen independent of q, x, a, b , and $i_{n_q}, \dots, i_{m_q+1}$.

Proof. The idea is to apply the ℓ^2 flattening lemma 4.7. For ease of notation, we will fix $n = n_q$ and $r = r_q$ throughout this proof. We will assume that a_0 is small enough so that we may apply Lemma 4.11.

Claim 1. We have the following bound:

$$(4.5) \quad \|\mu\|_1 \leq B.$$

We first observe that

$$f_n^{(a)}(i_n, \dots, i_1, x) = f_r^{(a)}(i_n, \dots, i_1, x) + f_m^{(a)}(i_m, \dots, i_1, x)$$

and

$$\begin{aligned} & |f_r^{(a)}(i_n, \dots, i_1, x) - f_r^{(a)}(i_n, \dots, i_{m+1}, \omega(i_{m+1}))| \\ & \leq \sum_{j=0}^{r-1} |f^{(a)}(i_{n-j}, \dots, i_{m+1}, \omega(i_{m+1})) - f^{(a)}(i_{n-j}, \dots, i_1, x)| \\ & \leq \sum_{j=0}^{r-1} \theta^{r-1-j} \text{Lip}_{d_\theta}(f^{(a)}) \leq \eta_\theta. \end{aligned}$$

Using the triangle inequality, we deduce

$$\begin{aligned} & f_r^{(a)}(i_n, \dots, i_1, x) \\ & \leq f_r^{(a)}(i_n, \dots, i_{m+1}, \omega(i_{m+1})) + |f_r^{(a)}(i_n, \dots, i_1, x) - f_r^{(a)}(i_n, \dots, i_{m+1}, \omega(i_{m+1}))| \\ & \leq f_r^{(a)}(i_n, \dots, i_{m+1}, \omega(i_{m+1})) + \eta\theta. \end{aligned}$$

We therefore have

$$\begin{aligned} \|\mu\|_1 & \leq \sum_{i_1, \dots, i_m} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \\ & \leq \sum_{i_1, \dots, i_m} e^{f_r^{(a)}(i_n, \dots, i_{m+1}, \omega(i_{m+1}))} e^{f_m^{(a)}(i_m, \dots, i_1, x)} e^{\eta\theta} \\ & \leq B. \end{aligned}$$

Claim 2. For some $\kappa_1 > 0$,

$$(4.6) \quad \|\mu\|_\infty \leq q^{-\kappa_1} B.$$

Using the bound on $m < n/d_0$, it suffices to bound

$$\begin{aligned} & \left| \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, \dots, i_1, x) = \gamma} e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} \right| \\ & \leq B \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, \dots, i_1, x) = \gamma} e^{f_m^{(a)}(i_m, \dots, i_1, x)} \\ & \leq B e^{m(|a| \sup \tau + |\log \lambda_a|)} \sum_{i_1, \dots, i_m, c_{m+1}(i_{m+1}, \dots, i_1, x) = \gamma} e^{f_m^{(0)}(i_m, \dots, i_1, x)} \\ & \leq B c' e^{m(|a| \sup \tau + |\log \lambda_a|)} e^{-\delta m \inf(\tau)} \text{ (see Lemma 4.13)} \\ & \leq B q^{-\kappa_1} \end{aligned}$$

as long as we choose a_0 so small that

$$\max_{|a| < a_0} e^{m(|\tau|_\infty |a| + |\log \lambda_a|)} \leq e^{\frac{1}{3} \delta \inf(\tau)},$$

and $\kappa_1 > 0$ is chosen such that (recalling $m \asymp \log q/d_0$)

$$q^{\kappa_1} \leq e^{\frac{1}{3} \delta m \inf(\tau)}$$

and $q > q_1 > (c')^{1/3\kappa_1}$.

Claim 3. We have

$$(4.7) \quad \|\pi_{q'}(\mu)\|_\infty < B(q')^{-\kappa}$$

for all $q' | q$ with $q' > q^{1/10}$ for some $\kappa > 0$.

Choose such a q' , and let $\Gamma_0 < \Gamma$ be a subgroup such that the projection $\pi_p(\Gamma_0)$ is a proper subgroup of $SL_2(p)$ for each divisor $p|q'$. As in [8] (see also Lemma 5.5 of [24] for more details), we know that $\#(a\Gamma_0 \cap B_m(e))$ grows sub-exponentially in m , and in particular, we have

$$\#(a\Gamma_0 \cap \mathcal{B}_m(e)) = O(q^{\kappa_1/2}),$$

with κ_1 as in Claim 2. Claim 3 now follows from Claim 2.

By Claims 1 and 3, we have now verified the conditions of the flattening lemma (Lemma 4.7). We therefore apply it to obtain

$$\|\mu * \phi\|_2 \leq BCq^{-\kappa'} \|\phi\|_2$$

for q large. □

The following bound, which will be useful in a number of places, follows from direct calculation.

Lemma 4.19. *There is $\tilde{c} > 0$ such that, for any $x, y \in \Sigma^+$ with $d_\theta(x, y) < 1$, any admissible sequence (i_{n_q}, \dots, i_1, x) , and any $|a| < 1, |b| < b_0$, we have*

$$\left| 1 - e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, y) - (f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} \right| \leq \tilde{c} \cdot d_\theta(x, y).$$

Remark. It is in this type of bound that large values of $|b|$ cause problems. This characterizes the difference between Dolgopyat’s approach and that of Bourgain-Gamburd-Sarnak.

Before moving on we will establish another proposition; this one will furnish Lipschitz bounds on $\hat{\mathcal{M}}_{ab,q}\tilde{H}$. For real numbers a, b , for $x, y \in \hat{\Sigma}^+$ with $d_\theta(x, y) \leq \theta$, and a sequence $i_{n_q}, \dots, i_{m_q+1}$, we define

$$\begin{aligned} (4.8) \quad \mu' &= \mu'_{x,y,i_{m_q+1}, \dots, i_{n_q}}^{(a,b)} \\ &= \sum_{i_1, \dots, i_{m_q}} \left(e^{(f_{n_q}^{(a)} + ib\tau_{n_q})(i_{n_q}, \dots, i_1, x)} - e^{(f_{n_q}^{(a)} + ib\tau_{n_q})(i_{n_q}, \dots, i_1, y)} \right) \delta_{\mathbf{c}_{m_q+1}(i_{m_q+1}, \dots, i_1, x)}. \end{aligned}$$

Proposition 4.20. *As usual, we write n, m, r for n_q, m_q, r_q . Let*

$$(4.9) \quad \mu' = \mu'_{x,y,i_{m+1}, \dots, i_n}^{(a,b)}$$

and

$$\begin{aligned} B' &= B'_{a,i_{m+1}, \dots, i_n} \\ &:= c\tilde{c}e^{\eta\theta} e^{f_r^{(a)}(i_n, \dots, i_{m+1}, \omega(i_{m+1}))} \end{aligned}$$

with $c > 1$ the constant from Lemma 4.11. There exist constants $\kappa > 0, a_0 > 0, q_1 > 1, C > 1$ such that for any $x, y \in \hat{\Sigma}^+$ with $d_\theta(x, y) < 1$ and any square free $q > q_1, |a| < a_0, |b| \leq b_0$

$$\|\mu' * \phi\|_2 \leq B'C'q^{-\kappa} \|\phi\|_2 d_\theta(x, y)$$

for all $\phi \in E_q^q$. The constants C, κ may be chosen independent of q, x, a, b , and i_n, \dots, i_{m+1} .

Proof.

Claim 1. *A calculation similar to that for (4.5) yields*

$$\|\mu'\|_1 \leq B'd_\theta(x, y).$$

Claim 2. *There exists $\kappa > 0$ with $\|\mu'\|_\infty < B'q^{-\kappa}$. For any $\gamma \in \text{SL}_2(q)$, we estimate*

$$\begin{aligned} |\mu'(\gamma)| &= \left| \sum_{i_1, \dots, i_m: \mathfrak{c}_{m+1}(i_{m+1}, \dots, i_1, x) = \gamma} e^{(f_n^{(a)} + ib\tau_n)_n(i_n, \dots, i_1, x)} - e^{(f_n^{(a)} + ib\tau_n)_n(i_n, \dots, i_1, y)} \right| \\ &\leq \sum_{i_1, \dots, i_m: \mathfrak{c}_{m+1}(i_{m+1}, \dots, i_1, x) = \gamma} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \left| 1 - e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, y) - (f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} \right| \\ &\leq \tilde{c}d_\theta(x, y) \sum_{i_1, \dots, i_m: \mathfrak{c}_{m+1}(i_{m+1}, \dots, i_1, x) = \gamma} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \text{ by Lemma 4.19.} \end{aligned}$$

The same argument as used in claim 2 of Proposition 4.18 now yields

$$\|\mu'\|_\infty < B'q^{-\kappa}$$

as expected.

Claim 3. *An argument similar to the one leading to Claim 3 in the proof of Proposition 4.18 gives that*

$$\|\|\pi_{q'}(\mu')\|\|_\infty \leq B'd_\theta(x, y)(q')^{-\kappa} \text{ for } q'|q, q' > q^{1/10}.$$

The proposition now follows from the ℓ^2 flattening lemma as in the proof of Proposition 4.18. \square

4.5. Supremum bounds and Lipschitz bounds. The purpose of all the estimates in the last two subsections is to provide bounds on the congruence transfer operators. We will need to bound both the supremum norms and the Lipschitz constants. Start with the supremum norm. We observe that $(\hat{\mathcal{M}}_{ab,q}^n \tilde{H})(x, g)$ has a good approximation by an appropriate sum of the convolutions, and we use this fact, together with the convolution estimates in Propositions 4.18 and 4.20 to estimate the supremum and Lipschitz norms of $(\hat{\mathcal{M}}_{ab,q}^n \tilde{H})$.

For any q , for $\tilde{H} \in \tilde{E}_q^q$, and a sequence $i_{n_q}, \dots, i_{m_q+1}$, define the function ϕ on $\text{SL}_2(q)$ by

$$(4.10) \quad \phi(g) = \phi_{\tilde{H}, (i_{n_q}, \dots, i_{m_q+1})}(g)$$

$$(4.11) \quad := \tilde{H}(i_{n_q}, \dots, i_{m_q+1}, \omega(i_{m_q+1})), g\mathfrak{c}_{r_q-1}^{-1}(i_{n_q}, \dots, i_{m_q+1}, \omega(i_{m_q+1})).$$

Note that

$$|\phi| \leq \|\tilde{H}\|_\infty \quad \text{and} \quad \phi \in E_q^q.$$

Lemma 4.21. *There exist $\tilde{C} > 1$ and $a_0 > 0$ such that the following holds for any $q, |a| < a_0, |b| \leq b_0, x \in \tilde{\Sigma}^+$, and any $\tilde{H} \in \tilde{E}_q^q$:*

$$\begin{aligned} &\left| (\hat{\mathcal{M}}_{ab,q}^{n_q} \tilde{H})(x, \cdot) - \sum_{i_{m_q+1}, \dots, i_{n_q}} \mu_{x, (i_{n_q}, \dots, i_{m_q+1})}^{a,b} * \phi_{\tilde{H}, (i_{n_q}, \dots, i_{m_q+1})}(\cdot) \right| \\ &\leq \tilde{C} \text{Lip}_{d_\theta}(\tilde{H})\theta^r. \end{aligned}$$

Proof. Fix q, x , and \tilde{H} and write $n = n_q, r = r_q, m = m_q$. Choose $|a| < a_0$, the constant from Lemma 4.11. For a sequence i_n, \dots, i_{m+1} , set

$$(4.12) \quad \phi(g) = \phi_{\tilde{H}, (i_n, \dots, i_{m+1})}(g).$$

Choose $|a| < a_0$, the constant from Lemma 4.11. We observe, as a consequence of the definitions and of Lemma 4.9, that

$$\begin{aligned} & \sum_{i_{m+1}, \dots, i_n} \mu_{x, (i_n, \dots, i_{m+1})}^{a, b} * \phi_{(i_n, \dots, i_{m+1})}(g) \\ = & \sum_{i_1, \dots, i_n} e^{f_n^{(a)} + ib\tau_n}(i_n, \dots, i_1, x) \tilde{H}((i_n, \dots, i_{m+1}, \omega(i_{m+1})), g c_n^{-1}(i_n, \dots, i_1, x)). \end{aligned}$$

We may therefore compute

$$\begin{aligned} & \left| \left(\hat{\mathcal{M}}_{ab, q}^n \tilde{H} \right)(x, \cdot) - \sum_{i_{m+1}, \dots, i_n} \mu_{x, (i_n, \dots, i_{m+1})} * \phi_{\tilde{H}, (i_n, \dots, i_{m+1})}(\cdot) \right| \\ \leq & \sum_{i_1, \dots, i_n} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \\ & \left| \left(\tilde{H}((i_n, \dots, i_1, x), \cdot) - \tilde{H}((i_n, \dots, i_{m+1}, \omega(i_{m+1})), \cdot) \right) \right| \\ \leq & \sum_{i_1, \dots, i_n} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \text{Lip}_{d_\theta}(\tilde{H}) \theta^{r-1} \\ \leq & c \text{Lip}_{d_\theta}(\tilde{H}) \theta^r, \end{aligned}$$

where $c > 1$ is the constant from Lemma 4.11. \square

The next lemma provides bounds on the supremum norm for $\hat{\mathcal{M}}_{ab, q}^{n_q} \tilde{H}$ using the description in terms of the convolutions we just proved together with the convolution estimate, Proposition 4.18.

Lemma 4.22. *There exist constants $a_0 > 0$, $q_1 > 0$, $\kappa' > 0$ such that the following holds for any $|a| < a_0$, $|b| \leq b_0$, $q > q_1$, and $\tilde{H} \in \tilde{E}_q^q$:*

$$\|\hat{\mathcal{M}}_{ab, q}^{n_q} \tilde{H}\|_\infty \leq \frac{1}{2} q^{-\kappa'} \left(\|\tilde{H}\|_\infty + \text{Lip}_{d_\theta}(\tilde{H}) \theta^{n_q/2} \right).$$

Proof. We choose $a_0 > 0$ small and q_1 large as on Lemma 4.11 and Proposition 4.18. Consider $q > q_1$, $|a| < a_0$, $|b| \leq b_0$, and a function $\tilde{H} \in \tilde{E}_q^q$. We write n for n_q and r for r_q . We recall the function ϕ on $\text{SL}_2(q)$ as in (4.11). Summing over all admissible sequences i_{m+1}, \dots, i_n and applying Lemma 4.21 and Proposition 4.18 we obtain, for $x \in \tilde{\Sigma}^+$,

$$\begin{aligned} & \left| \hat{\mathcal{M}}_{ab, q}^n \tilde{H}(x) \right| \\ \leq & \left| \sum_{i_{m+1}, \dots, i_n} \mu_{x, (i_n, \dots, i_{m+1})}^{a, b} * \phi_{\tilde{H}, (i_n, \dots, i_{m+1})}(\cdot) \right| + \tilde{C} \text{Lip}_{d_\theta}(\tilde{H}) \theta^r \\ \leq & q^{-\kappa} C \sum_{i_{m+1}, \dots, i_n} B_{i_n, \dots, i_{m+1}}^a |\phi_{\tilde{H}, (i_n, \dots, i_{m+1})}| + \tilde{C} \text{Lip}_{d_\theta}(\tilde{H}) \theta^r \\ \leq & q^{-\kappa} c C e^{\eta_\theta} \|\tilde{H}\|_\infty \sum_{i_{m+1}, \dots, i_n} e^{f_r^{(a)}(i_n, \dots, i_{m+1}, \omega(i_{m+1}))} + \tilde{C} \text{Lip}_{d_\theta}(\tilde{H}) \theta^r \\ \leq & c^2 C q^{-\kappa} e^{\eta_\theta} \|\tilde{H}\|_\infty + \tilde{C} \text{Lip}_{d_\theta}(\tilde{H}) \theta^r \end{aligned}$$

by Lemma 4.11. We may therefore choose $\kappa' > 0$ and $q_1 > 1$ so that $q^{-\kappa'} > 2c^2 C e^{\eta\theta} q^{-\kappa}$ and $2\theta^{r_q - n_q/2} \tilde{C} < q^{-\kappa'}$ for all $q > q_1$ and hence obtain

$$\|\hat{\mathcal{M}}_{ab,q}^n \tilde{H}\|_\infty \leq \frac{1}{2} q^{-\kappa'} \left(\|\tilde{H}\|_\infty + \text{Lip}_{d_\theta}(\tilde{H}) \theta^{n/2} \right)$$

so long as $q > q_1$. □

We'd like to iterate this argument, but before we can do that we need to estimate $\text{Lip}_{d_\theta}(\hat{\mathcal{M}}_{ab,q}^n \tilde{H})$. The proof of the next lemma is similar to the proof of the last one, though slightly longer.

Lemma 4.23. *There exist $C > 0, q_1 > 0, \kappa' > 0, a_0 > 0$ such that for all $|a| < a_0, |b| \leq b_0, q > q_1$, and $\tilde{H} \in \tilde{E}_q^q$, we have*

$$(4.13) \quad \text{Lip}_{d_\theta}(\hat{\mathcal{M}}_{ab,q}^n \tilde{H}) \leq \frac{1}{2} q^{-\kappa'} \left(\|\tilde{H}\|_\infty + \text{Lip}_{d_\theta}(\tilde{H}) \theta^{n_q/2} \right).$$

Proof. Again, we choose $a_0 > 0$ small and q_1 large as on Lemma 4.11 and Propositions 4.18 and 4.20. Consider $q > q_1, |a| < a_0, |b| \leq b_0$, and a function $\tilde{H} \in \tilde{E}_q^q$. We write n for n_q and r for r_q . For $x, y \in \hat{\Sigma}^+$ with $x_i = y_i$ for all $i \leq l$ (that is, with $d_\theta(x, y) \leq \theta^l < 1$) we have

$$\begin{aligned} & \left| \hat{\mathcal{M}}_{ab,q}^n \tilde{H}(x, g) - \hat{\mathcal{M}}_{ab,q}^n \tilde{H}(y, g) \right| \\ \leq & \sum_{i_1, \dots, i_n} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \\ & \left| \tilde{H}((i_n, \dots, i_1, x), g\mathbf{c}_n^{-1}(i_n, \dots, i_1, x)) - \tilde{H}((i_n, \dots, i_1, y), g\mathbf{c}_n^{-1}(i_n, \dots, i_1, y)) \right| \\ & + \left| \left(\sum_{i_1, \dots, i_n} e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} - e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, y)} \right) \right. \\ & \left. \tilde{H}((i_n, \dots, i_1, y), g\mathbf{c}_n^{-1}(i_n, \dots, i_1, y)) \right| \\ := & W + V. \end{aligned}$$

The first term W is bounded as

$$(4.14) \quad W \leq c \text{Lip}_{d_\theta}(\tilde{H}) \theta^n d_\theta(x, y)$$

by Lemma 4.11. We estimate the other term as

$$\begin{aligned} V & \leq \left| \sum_{i_1, \dots, i_n} \left(e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} - e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, y)} \right) \right. \\ & \quad \left. \times \tilde{H}((i_n, \dots, i_{m+1}, \omega(i_{m+1})), g\mathbf{c}_n^{-1}(i_n, \dots, i_1, x)) \right| \\ & \quad + \theta^{r-1} \text{Lip}_{d_\theta}(\tilde{H}) \sum_{i_1, \dots, i_n} \left| e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} - e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, y)} \right| \\ := & L + K. \end{aligned}$$

Next address K ,

$$\begin{aligned} K &= \theta^{r-1} \text{Lip}_{d_\theta}(\tilde{H}) \sum_{i_1, \dots, i_n} \left| e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} - e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, y)} \right| \\ &\leq \theta^{r-1} \text{Lip}_{d_\theta}(\tilde{H}) \sum_{i_1, \dots, i_n} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \\ &\quad \times \left| 1 - e^{(f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, y) - (f_n^{(a)} + ib\tau_n)(i_n, \dots, i_1, x)} \right| \\ &\leq \tilde{c}\theta^{r-1} d_\theta(x, y) \text{Lip}_{d_\theta}(\tilde{H}) \sum_{i_1, \dots, i_n} e^{f_n^{(a)}(i_n, \dots, i_1, x)} \end{aligned}$$

by Lemma 4.19. A final application of Lemma 4.11 then gives

$$(4.15) \quad K \leq c\tilde{c}\theta^{r-1} \text{Lip}_{d_\theta}(\tilde{H})d_\theta(x, y).$$

The bound on L uses the ℓ^2 flattening lemma once again. We observe that

$$L = \left| \sum_{i_n, \dots, i_{m+1}} \mu_{x, y, i_1, \dots, i_m}^{I_{a, b}} * \phi_{\tilde{H}}(i_n, \dots, i_{m+1}) \right|$$

for μ', ϕ as in (4.9), (4.11), respectively. Proposition 4.20 then gives

$$|\mu' * \phi| \leq C'q^{-\kappa'} B'_{a, i_{m+1}, \dots, i_n} \|\tilde{H}\|_\infty d_\theta(x, y),$$

and summation over i_n, \dots, i_{m+1} yields

$$(4.16) \quad L \leq C''q^{-\kappa'} d_\theta(x, y) \|\tilde{H}\|_\infty$$

for an appropriately chosen constant C'' (more precisely, $C'' = c^2\tilde{c}C'e^{\eta\theta}$ will do). Putting together Equations (4.14), (4.15), (4.16), we see that

$$(4.17) \quad \text{Lip}_{d_\theta}(\hat{\mathcal{M}}_{ab, q}^n \tilde{H}) \leq \tilde{C}'q^{-\kappa'} \left(\|\tilde{H}\|_\infty + \text{Lip}_{d_\theta}(\tilde{H})\theta^{n/2} \right)$$

for an appropriate constant \tilde{C}' . Now choose $\kappa'' = \kappa'/2$ and q_1 large enough that $C''q^{-\kappa''} < \frac{1}{2}$ for all $q > q_1$. Then

$$(4.18) \quad \text{Lip}_{d_\theta}(\hat{\mathcal{M}}_{ab, q}^n \tilde{H}) \leq \frac{1}{2}q^{-\kappa''} \left(\|\tilde{H}\|_\infty + \text{Lip}_{d_\theta}(\tilde{H})\theta^{n/2} \right). \quad \square$$

Proof of Theorem 4.4. Combining Lemmas 4.22 and 4.23, we obtain that for some $\kappa' > 0$,

$$\|\hat{\mathcal{M}}_{ab, q}^n \tilde{H}\|_\infty + \theta^{n/2} \text{Lip}_{d_\theta}(\hat{\mathcal{M}}_{ab, q}^n \tilde{H}) \leq q^{-\kappa'} (\|\tilde{H}\|_\infty + \theta^{n/2} \text{Lip}_{d_\theta}(\tilde{H})),$$

where $n = n_q$. Iterating, we obtain that for any $l \in \mathbb{N}$,

$$\begin{aligned} &\|\hat{\mathcal{M}}_{ab, q}^{ln} \tilde{H}\|_\infty + \theta^{n/2} \text{Lip}_{d_\theta}(\hat{\mathcal{M}}_{ab, q}^{ln} \tilde{H}) \\ &\leq q^{-l\kappa'} (\|\tilde{H}\|_\infty + \theta^{n/2} \text{Lip}_{d_\theta}(\tilde{H})). \end{aligned}$$

It follows that

$$\|\hat{\mathcal{M}}_{ab, q}^{ln} \tilde{H}\|_2 \leq \|\hat{\mathcal{M}}_{ab, q}^{ln} \tilde{H}\|_\infty \leq q^{-l\kappa'} \|\tilde{H}\|_{d_\theta}$$

as desired. □

5. UNIFORM MIXING OF THE BMS MEASURE AND THE HAAR MEASURE

We assume that $\Gamma < \text{SL}_2(\mathbb{Z})$ is convex cocompact. For each $q \in \mathbb{N}$, we denote by m_q^{BMS} the measure on $\Gamma(q) \backslash G$ induced by \tilde{m}^{BMS} and normalized so that its total mass is $\# \text{SL}_2(q)$.

5.1. Uniform exponential mixing. Our aim in this subsection is to prove Theorem 1.5 using Theorem 2.15 on spectral bounds for the transfer operators. Although this argument is similar to that contained in [17] and [1], we shall include it in order to understand the dependence of the implied constants on the level q . First we establish some more notation. We fix q such that $\Gamma(q) \backslash \Gamma = \text{SL}_2(q)$. We recall the equivalence relation $(u, t) \sim (\sigma u, t - \tau(u))$ on $\Sigma \times \mathbb{R}$ and the suspension space

$$\Sigma^\tau := \Sigma \times \mathbb{R} / \sim .$$

Definition. Similarly, we write

$$(5.1) \quad \hat{U}^{q,\tau} := \hat{U} \times \text{SL}_2(q) \times \mathbb{R}^+ / (u, \gamma, t + \tau(u)) \sim (\hat{\sigma}(u), \gamma c(u), t).$$

For a function $\phi : \hat{U}^{q,\tau} \rightarrow \mathbb{C}$, we say $\phi \in \mathcal{B}_0$ if $\|\phi\|_{\mathcal{B}_0} < \infty$ where

$$\|\phi\|_{\mathcal{B}_0} := \|\phi\|_\infty + \sup \left\{ \frac{|\phi(u, \gamma, s) - \phi(u', \gamma, s')|}{d(u, u') + |s - s'|} : u \neq u', \gamma \in \text{SL}_2(q), s \in [0, \tau(u)), s' \in [0, \tau(u')) \right\}.$$

We also say $\phi \in \mathcal{B}_1$ if $\|\phi\|_{\mathcal{B}_1} < \infty$ where

$$\|\phi\|_{\mathcal{B}_1} := \|\phi\|_\infty + \sup \{ \text{Var}_{0, \tau(u)}(t \mapsto \phi(u, \gamma, t)) : u \in \hat{U}, \gamma \in \text{SL}_2(q) \}.$$

For a bounded measurable function $\phi : \hat{U}^{q,\tau} \rightarrow \mathbb{C}$, we define the function $\hat{\phi}_\xi$ on $\hat{U} \times \text{SL}_2(q)$ by

$$\hat{\phi}_\xi(u, \gamma) := \int_0^{\tau(u)} \phi(u, \gamma, t) e^{-\xi t} dt;$$

we will sometimes regard this as a vector valued function on \hat{U} . The following lemma can be easily checked.

Lemma 5.1. *If $\psi \in \mathcal{B}_0$ with $\sum_{\gamma \in \Gamma} \psi(u, \gamma, s) = 0$ for all $(u, s) \in \hat{U}^\tau$, then $\hat{\psi}_\xi \in \mathcal{W}(\hat{U}, \mathbb{C}^{\text{SL}_2(q)})$ when considered as a vector valued function.*

For functions $\phi \in \mathcal{B}_1$ and $\psi \in \mathcal{B}_0$, we define the correlation function,

$$(5.2) \quad \tilde{\rho}_{\phi, \psi}(t) := \sum_{\gamma \in \text{SL}_2(q)} \int_{\hat{U}} \int_0^{\tau(u)} \phi(u, \gamma, s+t) \psi(u, \gamma, s) ds d\nu(u).$$

In order to establish an exponential decay for $\tilde{\rho}_{\phi, \psi}(t)$ for a suitable class of functions ϕ, ψ , we consider its Laplace transform and relate it with the transfer operators. We decompose $\tilde{\rho}_{\phi, \psi}(t)$ as

$$\begin{aligned} \tilde{\rho}_{\phi, \psi}(t) &= \sum_{\gamma \in \text{SL}_2(q)} \int_{\hat{U}} \int_{\max(0, \tau(u)-t)}^{\tau(u)} \phi(u, \gamma, s+t) \psi(u, \gamma, s) ds d\nu(u) \\ &\quad + \sum_{\gamma \in \text{SL}_2(q)} \int_{\hat{U}} \int_0^{\max(0, \tau(u)-t)} \phi(u, \gamma, s+t) \psi(u, \gamma, s) ds d\nu(u) \\ &:= \rho_{\phi, \psi}(t) + \bar{\rho}_{\phi, \psi}(t). \end{aligned}$$

The reason for this decomposition is that the Laplace transform of $\rho_{\phi,\psi}(t)$ can be expressed neatly in terms of transfer operators (see Lemma 5.2 below). More importantly, the Laplace transform of $\rho_{\phi,\psi}$ has better decay properties than the Laplace transform of $\tilde{\rho}_{\phi,\psi}$; this is needed when we apply the inverse Laplace transform at the end of the argument. Moreover, since $\tilde{\rho}_{\phi,\psi}(t) = \rho_{\phi,\psi}(t)$ for all $t \geq \sup \tau$, the exponential decay of $\tilde{\rho}_{\phi,\psi}(t)$ follows from that of $\rho_{\phi,\psi}(t)$.

So, consider the Laplace transform $\hat{\rho}$ of ρ : for $\xi \in \mathbb{C}$,

$$\hat{\rho}_{\phi,\psi}(\xi) = \int_0^\infty e^{-\xi t} \rho_{\phi,\psi}(t) dt.$$

For the rest of the section, we shall use the notation

$$\xi = a - ib.$$

The first task is to write $\hat{\rho}(\xi)$ in terms of the transfer operators:

Lemma 5.2. *For $\phi \in \mathcal{B}_1$, $\psi \in \mathcal{B}_0$, and $\Re(\xi) > 0$, we have*

$$\hat{\rho}_{\phi,\psi}(\xi) = \sum_{k=1}^\infty \lambda_a^k \int_{\hat{U}} \hat{\phi}_\xi(u) \cdot \hat{\mathcal{M}}_{ab,q}^k \hat{\psi}_{-\xi}(u) d\nu(u),$$

where λ_a is the lead eigenvalue of $\mathcal{L}_{-(\delta+a)\tau}$ as in Section 2. The right hand side should be understood as an inner product between two vectors in $\mathbb{C}^{\text{SL}_2(q)}$.

Proof. We calculate

$$\begin{aligned} \hat{\rho}_{\phi,\psi}(\xi) &= \sum_{\gamma \in \text{SL}_2(q)} \int_{\hat{U}} \int_{s=0}^{\tau(u)} \int_{\tau(u)-s}^\infty e^{-\xi t} \phi(u, \gamma, s+t) \psi(u, \gamma, a) dt ds d\nu(u) \\ &= \sum_{\gamma \in \text{SL}_2(q)} \int_{\hat{U}} \int_0^{\tau(u)} \int_{\tau(u)}^\infty e^{-\xi(t-s)} \phi(u, \gamma, t) \psi(u, \gamma, s) dt ds d\nu(u) \\ &= \sum_{\gamma \in \text{SL}_2(q)} \sum_{k=1}^\infty \int_{\hat{U}} \int_0^{\tau(u)} \int_{\tau_k(u)}^{\tau_{k+1}(u)} e^{-\xi(t-s)} \phi(u, \gamma, t) \psi(u, \gamma, s) dt ds d\nu(u) \\ &= \sum_{\gamma \in \text{SL}_2(q)} \sum_{k=1}^\infty \int_{\hat{U}} \int_0^{\tau(u)} \\ &\quad \times \int_0^{\tau(\hat{\sigma}^k u)} e^{-\xi(t+\tau_k(u)-s)} \phi(\hat{\sigma}^k(u), \gamma c_k(u), t) \psi(u, \gamma, s) dt ds d\nu(u) \\ &= \sum_{\gamma \in \text{SL}_2(q)} \sum_{k=1}^\infty \int_{\hat{U}} e^{-\xi \tau_k(u)} \hat{\phi}_\xi(\hat{\sigma}^k(u), \gamma c_k(u)) \hat{\psi}_{-\xi}(u, \gamma) d\nu(u) \\ &= \sum_{k=1}^\infty \lambda_a^k \int_{\hat{U}} \hat{\phi}_\xi(u) \cdot \hat{\mathcal{M}}_{ab,q}^k \hat{\psi}_{-\xi}(u) d\nu(u) \end{aligned}$$

using the fact that $\hat{\mathcal{L}}_{00}^*(\nu) = \nu$. □

Lemma 5.3. *If $\phi \in \mathcal{B}_1$, then $\|\hat{\phi}_\xi\|_2 \leq \|\hat{\phi}_\xi\|_\infty \leq \frac{2\sqrt{\#\text{SL}_2(q)}e^{|\alpha| \sup(\tau)} \|\phi\|_{\mathcal{B}_1}}{\max(1,|b|)}$.*

Proof. This follows from integration by parts in the flow direction. □

Lemma 5.4. *If $\psi \in \mathcal{B}_0$, then $\|\hat{\psi}_\xi\|_{1,b} \leq \frac{\sqrt{\#\text{SL}_2(q)}e^{|\alpha| \sup(\tau)} (3 \sup(\tau) + \text{Lip}_a(\tau)) \|\psi\|_{\mathcal{B}_0}}{\max(1,|b|)}$.*

Proof. The trivial bound $\text{var}_{[0, \tau(u)]} \psi(u, \gamma, \cdot) \leq \|\psi\|_{\mathcal{B}_0} \sup(\tau)$ provides

$$(5.3) \quad \|\hat{\psi}_\xi\|_\infty \leq \frac{2\sqrt{\#\text{SL}_2(q)} e^{|a| \sup(\tau)} \sup(\tau) \|\psi\|_{\mathcal{B}_0}}{\max(1, |b|)}.$$

On the other hand, consider any $u, u' \in \hat{U}$, $\gamma \in \text{SL}_2(q)$, and suppose, without loss of generality, that $\tau(u') \geq \tau(u)$. Then

$$\begin{aligned} & |\hat{\psi}_\xi(u, \gamma) - \hat{\psi}_\xi(u', \gamma)| \\ & \leq \int_0^{\tau(u)} |\psi(u, \gamma, t) - \psi(u', \gamma, t)| e^{|a|t} dt + \int_{\tau(u)}^{\tau(u')} |\psi(u', \gamma, t)| e^{|a|t} dt \\ & \leq d(u, u') e^{|a| \sup(\tau)} (\sup(\tau) \|\psi\|_{\mathcal{B}_0} + \text{Lip}_d(\tau) \|\psi\|_\infty). \end{aligned}$$

Together with (5.3), this proves the claim. □

We will now use the spectral bounds (Theorem 2.15) to prove a rate of decay for the correlation functions.

Proposition 5.5. *Let a_0, q_0, q'_0 be as in Theorem 2.15. There exist $C > 0, \eta > 0$ such that for all square free q with $(q, q_0 q'_0) = 1$, we have*

$$|\tilde{\rho}_{\phi, \psi}(t)| \leq C q^C \|\phi\|_{\mathcal{B}_1} \|\psi\|_{\mathcal{B}_0} e^{-\eta t}$$

for all $\phi \in \mathcal{B}_1$ and $\psi \in \mathcal{B}_0$ satisfying $\sum_{\gamma \in \text{SL}_2(q)} \psi(u, \gamma, s) = 0$.

Proof. We will establish that the Laplace transform $\hat{\rho}_{\phi, \psi}$ extends to an appropriate half plane and then apply the inversion formula. Lemma 5.2 gives

$$\hat{\rho}_{\phi, \psi}(\xi) = \sum_{k=0}^\infty \lambda_a^k \int_{\hat{U}} \hat{\phi}_\xi(u) \cdot \hat{\mathcal{M}}_{ab, q}^k \hat{\psi}_{-\xi}(u) d\nu(u)$$

for $\Re(\xi) > 0$. We claim an analytic continuation of $\hat{\rho}_{\phi, \psi}(\xi)$ to $\Re(\xi) > -a_0$ for some $a_0 > 0$. Each term of the above infinite sum is analytic, so it suffices to check that the sum is absolutely convergent. For $|a| \leq \min(1, a_0)$, Theorem 2.15, together with Lemma 5.1, gives that for some $\epsilon > 0$,

$$\begin{aligned} \lambda_a^k \int_{\hat{U}} \hat{\phi}_\xi(u) \cdot \hat{\mathcal{M}}_{ab, q}^k \hat{\psi}_{-\xi}(u) d\nu & \leq \lambda_a^k \|\hat{\mathcal{M}}_{ab, q}^k \hat{\psi}_{-\xi}\|_2 \|\phi_\xi\|_2 \\ & \leq \lambda_a^k C q^C e^{-\epsilon k} \|\hat{\psi}_{-\xi}\|_{1, b} \|\phi_\xi\|_2 \\ & \leq \lambda_a^k \frac{C' q^{C'} e^{-\epsilon k}}{\max(1, |b|)^2} \|\psi\|_{\mathcal{B}_0} \|\phi\|_{\mathcal{B}_1}, \end{aligned}$$

where C' is given by Lemmas 5.3 and 5.4; this is clearly summable so long as we choose a_0 small enough that

$$\max_{|a| \leq a_0} \lambda_a \leq e^{\epsilon/2}.$$

This computation also gives that for some absolute constant $C_1 > 0$,

$$|\hat{\rho}_{\phi, \psi}(\xi)| \leq \frac{C_1 q^{C'}}{1 + |b|^2} \|\psi\|_{\mathcal{B}_0} \|\phi\|_{\mathcal{B}_1}$$

for all ξ with $|\Re(\xi)| < a_0$. Now $\rho_{\phi, \psi}(t)$ is Lipschitz, so we may apply the inverse Laplace transform formula [44, Chapter II, Theorem 7.3] and obtain for all $t > 0$,

$$(5.4) \quad \rho_{\phi, \psi}(t) = e^{-\frac{a_0}{2} t} \lim_{T \rightarrow \infty} \int_{-\frac{a_0}{2} - iT}^{-\frac{a_0}{2} + iT} \hat{\rho}_{\phi, \psi}\left(-\frac{a_0}{2} - ib\right) e^{-ibt} db.$$

Since $\int_{-\frac{a_0}{2}-iT}^{-\frac{a_0}{2}+iT} |\hat{\rho}_{\phi,\psi}(-\frac{a_0}{2} - ib)| db \ll q^{C'} \int_0^T \frac{1}{1+b^2} db < \infty$, the limit in the right hand side of (5.4) is $O(q^{C'})$ with the implied constant independent of t , yielding the result for a uniform constant $C > 0$ with $\rho_{\phi,\psi}$ in place of $\tilde{\rho}_{\phi,\psi}$. Since those two functions agree on $t > \sup(\tau)$, and since $\tilde{\rho}_{\phi,\psi}$ is bounded as $q^C \|\psi\|_{\mathcal{B}_0} \|\phi\|_{\mathcal{B}_1}$, the result follows. \square

We can convert a function ϕ on $\Gamma(q)\backslash G$ to give a function ϕ_t on $\hat{U}^{q,\tau}$ as follows: for $t > 0$, $u \in \hat{U}_i, 0 \leq s \leq \tau(u)$, and $\gamma \in \Gamma(q)\backslash \Gamma$, we set \tilde{u} to be the lift of u to \tilde{U} , and

$$(5.5) \quad \phi_t(u, \gamma, s) := \int_{\tilde{S}_i} \phi(\gamma[\tilde{u}, \tilde{y}]a_{t+s}) d\nu_u(\tilde{y}),$$

where ν_u is the probability measure on \tilde{S}_i conditioned from the measure ν at u . For a general $s > 0$, we define

$$(5.6) \quad \phi_t(u, \gamma, s) := \phi_t(\hat{\sigma}^k(u), \gamma c_k(u), s - \tau_k(u)),$$

where $k \in \mathbb{N}$ is such that $0 \leq s - \tau_k(u) \leq \tau(\hat{\sigma}^k(u))$. By the equivalence relation (5.1), this defines ϕ_t on all of $\hat{U} \times \text{SL}_2(q) \times \mathbb{R}_{\geq 0}$.

Lemma 5.6. *There exists $C > 0$ such that, for any $\tilde{y} \in \tilde{S}_i$, we have*

$$|\phi(\gamma[\tilde{u}, \tilde{y}]a_{2t+s}) - \phi_t(u, \gamma, s + t)| \leq C e^{-t} \|\phi\|_{C^1}.$$

Proof. Let $u, \tilde{u}, \tilde{y} \in \tilde{S}_i$ be as above. Choose $k \in \mathbb{N}$ such that $0 \leq t + s - \tau_k(u) \leq \tau(\hat{\sigma}^k u)$, and write $u' = \hat{\sigma}^k u \in U_j$. Set \tilde{u}' to be the lift of u' to \tilde{U}_j . If $y' \in S_j$ with lift $\tilde{y}' \in \tilde{S}_j$, then the definition of the cocycle c tells us that both $\gamma[\tilde{u}, \tilde{y}]a_{\tau_k(u)}$ and $\gamma c_k(u)[\tilde{u}', \tilde{y}']$ lie in the stable leaf of $\gamma c_k(u)\tilde{R}_j \subset G$. It follows that for some $c_1 > 0$,

$$d(\gamma[\tilde{u}, \tilde{y}]a_{2t+s}, \gamma c_k(u)[\tilde{u}', \tilde{y}']a_{2t+s-\tau_k(u)}) \leq c_1 e^{-(2t+s-\tau_k(u))} \leq c_1 e^{-t},$$

and so that

$$|\phi(\gamma[\tilde{u}, \tilde{y}]a_{2t+s}) - \phi(\gamma c_k(u)[\tilde{u}', \tilde{y}']a_{2t+s-\tau_k(u)})| \leq c_1 \|\phi\|_{C^1} e^{-t}.$$

Integrating this inequality over $\tilde{y}' \in S_j$ and using (5.6) which gives

$$\int_{\tilde{S}_j} \phi(\gamma c_k(u)[\tilde{u}', \tilde{y}']a_{2t+s-\tau_k(u)}) d\nu_u = \phi_t(u, \gamma, s + t),$$

this gives the required result. \square

We therefore have the following lemma (cf. [1, Lemma 8.2]).

Lemma 5.7. *There are constants $\eta > 0, C > 0$ independent of ψ, ϕ, q such that*

$$\left| \int_{\Gamma(q)\backslash G} (\phi \circ a_{2t}) \cdot \psi \, dm_q^{\text{BMS}} - \frac{\tilde{\rho}_{\phi_t, \psi_0}(t)}{\nu(\tau)} \right| < C \cdot \#\text{SL}_2(q) \cdot \|\phi\|_{C^1} \|\psi\|_{\infty} e^{-\eta t}$$

for all $\phi, \psi \in C^1(\Gamma(q)\backslash G)$.

Proof of Theorem 1.5. We assume that $m^{\text{BMS}}(\Gamma\backslash G) = 1$ without loss of generality, so that the total mass of m_q^{BMS} is equal to $\#\text{SL}_2(q)$. Fix q with $(q, q_0 q'_0) = 1$ and compactly supported functions $\psi, \phi \in C^1(\Gamma(q)\backslash G)$.

We write

$$\psi = \psi' + \psi'',$$

where ψ' is (left) Γ invariant, and ψ'' satisfies $\sum_{\gamma \in \Gamma(q) \setminus \Gamma} \psi''(\gamma x) = 0$ for all $x \in \Gamma(q) \setminus G$. Exponential mixing of ψ' (with constant independent of q) follows from the bounds established in Section 3 together with the complex RPF theorem, as was carried out in the work of Dolgopyat and Stoyanov [39]. So we can and shall assume that $\psi = \psi''$, so that $\sum_{\gamma \in \Gamma(q) \setminus \Gamma} \psi(\gamma x) = 0$ for all $x \in \Gamma(q) \setminus G$.

We consider the functions ϕ_t, ψ_0 as defined in (5.5); note that ψ_0 satisfies $\sum_{\gamma \in \Gamma} \psi_0(u, \gamma, s) = 0$ and that $\|\psi_0\|_{\mathcal{B}_0} \ll \|\psi\|_{C^1}$. We also need to bound $\|\phi_t\|_{\mathcal{B}_1}$. It is clear that $\sup |\phi_t| \leq \|\phi\|_{C^1}$. On the other hand, we know that, for fixed u , and s such that (u, s) is not of the form $(u', 0)$, $\phi_t(u, s)$ is differentiable in the flow direction with derivative bounded by $\|\phi\|_{C^1}$. On the other hand, there are at most $\frac{\sup \tau}{\inf \tau} + 1$ values of s such that $(u, s) \sim (u', 0)$. Each of these may be a discontinuity, but each jump is at most $2\|\phi\|_{C^1}$. We can therefore bound the variation as

$$\text{var}_{[0, \tau(u))}(s \rightarrow \phi_t(u, \gamma, s)) \leq \left(\tau(u) + 2 \left(\frac{\sup \tau}{\inf \tau} + 1 \right) \right) \|\phi\|_{C^1}.$$

In other words,

$$\|\phi_t\|_{\mathcal{B}_1} \ll \|\phi\|_{C^1}.$$

Now calculate: for any $t > 0$

$$\begin{aligned} & \left| \int_{\Gamma(q) \setminus G} \phi(ga_{2t})\psi(g) dm^{\text{BMS}} \right| \\ & \leq \left| \frac{\tilde{\rho}_{\phi_t, \psi_0}(t)}{\nu(\tau)} \right| + C(\#\text{SL}_2(q)) \|\phi\|_{C^1} \|\psi\|_{C^0} e^{-\eta t} \\ & \leq C' q^{C'} (\|\psi_0\|_{\mathcal{B}_0} \|\phi_t\|_{\mathcal{B}_1} + \|\phi\|_{C^1} \|\psi\|_{C^0}) e^{-\eta' t} \\ & \leq C'' q^{C''} \|\phi\|_{C^1} \|\psi\|_{C^0} e^{-\eta' t} \end{aligned}$$

for some $C', C'', \eta' > 0$, by Lemma 5.7 and Proposition 5.5. □

5.2. Exponential decay of the matrix coefficients. Let Γ be a geometrically finite subgroup of $\text{PSL}_2(\mathbb{R})$. We begin by recalling the definitions of measures m^{BR} , $m^{\text{BR}*}$, and m^{Haar} . Similar to the definition of the BMS measure

$$d\tilde{m}^{\text{BMS}}(u) = e^{\delta\beta_{u^+}(o, u)} e^{\delta\beta_{u^-}(o, u)} d\mu_o^{\text{PS}}(u^+) d\mu_o^{\text{PS}}(u^-) ds$$

given in Section 2, the measures $\tilde{m}^{\text{BR}} = \tilde{m}_{\Gamma}^{\text{BR}}$, $\tilde{m}^{\text{BR}*} = \tilde{m}_{\Gamma}^{\text{BR}*}$, and \tilde{m}^{Haar} on $\text{PSL}_2(\mathbb{R})$ are defined as follows:

$$\begin{aligned} d\tilde{m}^{\text{BR}}(u) &= e^{\beta_{u^+}(o, u)} e^{\delta\beta_{u^-}(o, u)} dm_o(u^+) d\mu_o^{\text{PS}}(u^-) ds; \\ d\tilde{m}^{\text{BR}*}(u) &= e^{\delta\beta_{u^+}(o, u)} e^{\beta_{u^-}(o, u)} dm_o(u^-) d\mu_o^{\text{PS}}(u^+) ds; \\ d\tilde{m}^{\text{Haar}}(u) &= e^{\beta_{u^+}(o, u)} e^{\beta_{u^-}(o, u)} dm_o(u^+) dm_o(u^-) ds, \end{aligned}$$

where m_o is the unique probability measure on $\partial(\mathbb{H}^2)$ which is invariant under the stabilizer of o .

These measures are all left Γ -invariant and induce measures on $\Gamma \setminus G$, which we will denote by $m^{\text{BR}}, m^{\text{BR}*}, m^{\text{Haar}}$, respectively.

Let

$$N = \{n_s := \begin{pmatrix} 1 & 0 \\ s & 1 \end{pmatrix} : s \in \mathbb{R}\} \quad \text{and} \quad H = \{h_s := \begin{pmatrix} 1 & s \\ 0 & 1 \end{pmatrix} : s \in \mathbb{R}\}.$$

For $g \in G$, denote by g^\pm the forward and backward end points of the geodesic determined by g and set

$$\alpha(g, \Lambda(\Gamma)) := \inf\{|s| : (gn_s)^+ \in \Lambda(\Gamma)\} + \inf\{|s| : (gh_s)^- \in \Lambda(\Gamma)\} + 1.$$

It follows from the continuity of the visual map that for any compact subset $\mathcal{Q} \subset G$,

$$\alpha(\mathcal{Q}, \Lambda(\Gamma)) := \sup \alpha(g, \Lambda(\Gamma)) < \infty.$$

If Γ' is a normal subgroup of Γ of finite index, then $\Lambda(\Gamma) = \Lambda(\Gamma')$, and hence $\alpha(\mathcal{Q}, \Lambda(\Gamma)) = \alpha(\mathcal{Q}, \Lambda(\Gamma'))$. Therefore the following theorem implies that Theorem 1.1 can be deduced from Theorem 1.5; note that even though we need the following theorem only for Γ convex cocompact in this paper, we record it for a general geometrically finite group Γ of G for future reference. Let $\pi : G \rightarrow \Gamma \backslash G$ be the canonical projection.

Theorem 5.8. *Let $\mathcal{Q} \subset G$ be a compact subset. Suppose that there exist constants $c_\Gamma > 0$ and $\eta_\Gamma > 0$ such that for any $\Psi, \Phi \in C^1(\Gamma \backslash G)$ supported on $\pi(\mathcal{Q})$,*

$$(5.7) \quad \int_{\Gamma \backslash G} \Psi(ga_t)\Phi(g)dm^{\text{BMS}} = \frac{m^{\text{BMS}}(\Psi) \cdot m^{\text{BMS}}(\Phi)}{m^{\text{BMS}}(\Gamma \backslash G)} + O(c_\Gamma \cdot \|\Psi\|_{C^1} \|\Phi\|_{C^1} \cdot e^{-\eta_\Gamma t}),$$

where the implied constant depends only on \mathcal{Q} . Then for any $\Psi, \Phi \in C^1(\Gamma \backslash G)$ supported on $\pi(\mathcal{Q})$, as $t \rightarrow +\infty$,

$$(5.8) \quad e^{(1-\delta)t} \int_{\Gamma \backslash G} \Psi(ga_t)\Phi(g)dm^{\text{Haar}} \\ = \frac{m^{\text{BR}}(\Psi) \cdot m^{\text{BR}^*}(\Phi)}{m^{\text{BMS}}(\Gamma \backslash G)} + O(c_\Gamma \cdot \|\Psi\|_{C^1} \|\Phi\|_{C^1} \cdot e^{-\eta'_\Gamma t}),$$

where $\eta'_\Gamma = \frac{\eta_\Gamma}{8+2\eta_\Gamma}$ and the implied constant depends only on \mathcal{Q} and $\alpha(\mathcal{Q}, \Lambda(\Gamma))$.

The rest of this section is devoted to the proof of this theorem. The proof involves effectivizing the original argument of Roblin [35], extended in [38], [29], [27], while making the dependence of the implied constant on the relevant functions precise.

For $\epsilon > 0$ and a subset S of G , S_ϵ denotes the set $\{s \in S : d(s, e) \leq \epsilon\}$.

Let

$$P := HA.$$

Then the sets $B_\epsilon := P_\epsilon N_\epsilon$, $\epsilon > 0$ form a basis of neighborhoods of e in G .

For $g \in \text{PSL}_2(\mathbb{R})$, we define measures on gN ,

$$d\tilde{\mu}_{gN}^{\text{Leb}}(gn) = e^{\beta_{(gn)^+}(o, gn)} dm_o(gn^+); \\ d\tilde{\mu}_{gN}^{\text{PS}}(gn) = e^{\delta\beta_{(gn)^+}(o, gn)} d\mu_o^{\text{PS}}(gn^+).$$

If $x = [g] \in \Gamma \backslash G$, for a compact subset N_0 of N such that gN_0 injects to $\Gamma \backslash G$, and for a function ψ on xN_0 , we write $d\mu_{xN}^{\text{Leb}}(\psi)$ and $d\mu_{xN}^{\text{PS}}(\psi)$ for the pushforward of the above measures to xN_0 via the isomorphism gN_0 with xN_0 . The measure $d\tilde{\mu}_{gN}^{\text{Leb}}(gn)$ is simply the Haar measure on N , and hence we write $dn = d\tilde{\mu}_{gN}^{\text{Leb}}(gn)$.

The quasi-product structure of \tilde{m}^{BMS} is a key ingredient in the arguments below: for $\Psi \in C_c(G)$ supported on gB_ϵ for all $\epsilon > 0$ small,

$$\tilde{m}^{\text{BMS}}(\Psi) = \int_{gP_\epsilon} \int_{gpN_\epsilon} \Psi(gpn) d\tilde{\mu}_{gpN}^{\text{PS}}(gpn) dv_{gP}(gp),$$

where $dv_{gP}(gp) = e^{\delta\beta_{(gp)^-}(o, gp)} d\mu_o^{\text{PS}}(gp^-) ds$ for $s = \beta_{gp^-}(o, gp)$.

In the rest of this section, we fix a compact subset \mathcal{Q} of G , and assume that the hypotheses of Theorem 5.8 are satisfied for functions supported in $\pi(\mathcal{Q})$. Let $2\epsilon_0 > 0$ be the injectivity radius of $\pi(\mathcal{Q})$. Fix $x = [g] \in \pi(\mathcal{Q})$ and functions $\Psi, \Phi \in C^1(\Gamma \backslash G)$ which are supported in $xB_{\epsilon_0/2}$.

Proposition 5.9. *Fix $y \in xP_{\epsilon_0}$ and put $\phi := \Phi|_{yN_{\epsilon_0}} \in C^1(yN_{\epsilon_0})$. Then for $t > 1$,*

$$\int_{yN_{\epsilon_0}} \Psi(yna_t)\phi(y)n d\mu_{yN}^{\text{PS}}(yn) = \frac{\mu_{yN}^{\text{PS}}(\phi)}{|m^{\text{BMS}}|} m^{\text{BMS}}(\Psi) + O(c_\Gamma \|\Psi\|_{C^1} \|\phi\|_{C^1} e^{-\eta_1 t}),$$

where $\eta_1 = \eta_\Gamma / (4 + \eta_\Gamma)$ and the implied constant depends only on \mathcal{Q} and $\alpha(\mathcal{Q}, \Lambda(\Gamma))$.

Proof. Set $R_0 := \alpha(y, \Lambda(\Gamma)) + 2$. For a sufficiently small $\epsilon \in (0, 1)$, if we set $t_0 := \log(R_0 \epsilon^{-1})$, $y_0 = ya_{t_0}$, then $\nu(y_0 P_\epsilon) > 0$. Hence we may choose a smooth positive function ρ_ϵ supported on $y_0 P_\epsilon$ such that $\nu(\rho_\epsilon) = 1$ and that $\|\rho_\epsilon\|_{C^1} \ll \epsilon^{-3}$. Define a C^1 -function Φ^\dagger supported on $y_0 P_\epsilon N_{\epsilon^{-1}R_0}$ as follows:

$$\Phi^\dagger(y_0 pn) := e^{-\delta t_0 - \delta \beta_{n_p^+}(n_p, pn)} \phi(y_0 n_p a_{-t_0}) \rho_\epsilon(y_0 p),$$

where $n_p \in N$ is the unique element such that $p^{-1}n_p \in nP$. This is well-defined if $R_0 < \epsilon_0$. When R_0 is large, we should define Φ^\dagger as the projection to $\Gamma \backslash G$ of a suitably defined function on G . We have $m^{\text{BMS}}(\Phi^\dagger) = \mu_{yN}^{\text{PS}}(\phi)$. Now by the hypothesis of Theorem 5.8, we have

$$\begin{aligned} \int_{yN_{\epsilon_0}} \Psi(yna_t)\phi(y)n d\mu_{yN}^{\text{PS}}(yn) &= (1 + O(\epsilon)) \langle a_{t-t_0} \Psi, \Phi^\dagger \rangle_{m^{\text{BMS}}} \\ &= (1 + O(\epsilon)) \left(\frac{\mu_{yN}^{\text{PS}}(\phi)}{|m^{\text{BMS}}|} m^{\text{BMS}}(\Psi) + O(c_\Gamma \epsilon^{-3} e^{-\eta(t-t_0)}) \right) \\ &= \frac{\mu_{yN}^{\text{PS}}(\phi)}{|m^{\text{BMS}}|} m^{\text{BMS}}(\Psi) + O(\epsilon + c_\Gamma R_0^\eta \epsilon^{-\eta-3} e^{-\eta t}), \end{aligned}$$

where the implied constant depends only on the C^1 -norms of Ψ , ϕ , and \mathcal{Q} . By taking $\epsilon = e^{-\eta t / (4 + \eta)}$ and by setting $\eta_1 := \eta / (4 + \eta)$, we obtain

$$\int_{yN_{\epsilon_0}} \Psi(yna_t)\phi(y)n d\mu_{yN}^{\text{PS}}(yn) = \frac{\mu_{yN}^{\text{PS}}(\phi)}{|m^{\text{BMS}}|} m^{\text{BMS}}(\Psi) + O(c_\Gamma R_0^\eta e^{-\eta_1 t}).$$

Since R_0 is bounded above by $\alpha(\mathcal{Q}, \Lambda(\Gamma))$, this proves the claim. □

Proposition 5.10. *Keeping the same notation as in Proposition 5.9, we have*

$$e^{(1-\delta)t} \int_{yN_{\epsilon_0}} \Psi(yna_t)\phi(y)n dn = \frac{\mu_{yN}^{\text{PS}}(\phi)}{|m^{\text{BMS}}|} m^{\text{BR}}(\Psi) + O(c_\Gamma \|\Psi\|_{C^1} \|\phi\|_{C^1} e^{-\eta_1 t/2}),$$

where the implied constant depends only on \mathcal{Q} and $\alpha(\mathcal{Q}, \Lambda(\Gamma))$.

Proof. We deduce this proposition from Proposition 5.9 by comparing the two integrals on the left hand sides via transversal intersections.

Define $\phi_\epsilon^\pm \in C^1(yN)$ by

$$(5.9) \quad \phi_\epsilon^+(yn) = \sup_{n' \in N_\epsilon} \phi(ynn') \quad \text{and} \quad \phi_\epsilon^-(yn) = \inf_{n' \in N_\epsilon} \phi(ynn').$$

Fix $R_1 := \alpha(\mathcal{Q}, \Lambda(\Gamma)) + 1$. For each $p \in P_{\epsilon_0}$, let $N_p := \{n \in N : (pn)^+ = n_s^+ \text{ for some } |s| < R_1\}$; then $\mu_{xpN}^{\text{PS}}(xpN_p) > 0$, and the map $xp \mapsto \mu_{xpN}^{\text{PS}}(xpN_p)$ is a positive smooth function on xP_{ϵ_0} . Set $B'_{\epsilon_0} := \cup_{p \in P_{\epsilon_0}} pN_p$; we may assume that the map $g \rightarrow xg$ is injective on B'_{ϵ_0} by replacing ϵ_0 by a smaller number if necessary.

Define the finite set

$$P_x(t) := \{p \in P_{\epsilon_0} : xpn \in \text{supp}(\phi)a_t \text{ for some } n \in N_p\}.$$

Define functions ψ and Ψ' supported on xP_{ϵ_0} and xB'_{ϵ_0} , respectively,

$$\psi(xp) := \int_{xpN_{\epsilon_0}} \Psi(xpn)dn \quad \text{and} \quad \Psi'(xpn) := \frac{\psi(xp)}{\mu_{xpN}^{\text{PS}}(xpN_p)} \text{ for } pn \in B'_{\epsilon_0}.$$

We then have $m^{\text{BMS}}(\Psi') = \nu_{xP}(\psi) = m^{\text{BR}}(\Psi)$, and we can find C^1 -approximations $\Psi'_{\epsilon, -} \leq \Psi' \leq \Psi'_{\epsilon, +}$ such that $m^{\text{BMS}}(\Psi'_{\epsilon, \pm}) = m^{\text{BMS}}(\Psi') + O(\epsilon^\delta)$, and $\|\Psi'_{\epsilon, \pm}\|_{C^1} = O(\epsilon^{-1}\|\Psi\|_{C^1})$. The following computation holds for all small $0 < \epsilon \ll \epsilon_0$:

$$\begin{aligned} & e^{(1-\delta)t} \int_{yN} \Psi(yna_t)\phi(yn)dn \\ &= (1 + O(\epsilon))e^{-\delta t} \sum_{p \in P_x(t)} \psi(xp)\phi_{ce^{-t}\epsilon_0}^\pm(xpa_{-t}) \\ &= (1 + O(\epsilon)) \int_{yN} \Psi'(yna_t)\phi_{c'(\epsilon_0+R_1)e^{-t}}^\pm(yn)d\mu_{yN}^{\text{PS}}(yn) \\ &= (1 + O(\epsilon)) \int_{yN} \Psi'_{\epsilon, \pm}(yna_t)\phi_{c'(\epsilon_0+R_1)e^{-t}}^\pm(yn)d\mu_{yN}^{\text{PS}}(yn) \\ &= (1 + O(\epsilon) + O((\epsilon_0 + R_1)e^{-t})) \left(\frac{m^{\text{BR}}(\Psi)\mu_{yN}^{\text{PS}}(\phi)}{|m^{\text{BMS}}|} + O(c_\Gamma\epsilon^{-1}\|\Psi\|_{C^1}\|\phi\|_{C^1}e^{-\eta_1 t}) \right) \end{aligned}$$

by Proposition 5.9 (we refer the reader to [29] and [27] for details in this step).

Therefore taking $\epsilon = e^{-\eta_1 t/2}$,

$$e^{(1-\delta)t} \int_{yN} \Psi(yna_t)\phi(yn)dn = \frac{m^{\text{BR}}(\Psi)\mu_{yN}^{\text{PS}}(\phi)}{|m^{\text{BMS}}|} + O(c_\Gamma\|\Psi\|_{C^1}\|\phi\|_{C^1}e^{-\eta_1 t/2}),$$

where the implied constant depends only on ϵ_0 and R_1 , and hence only on \mathcal{Q} and $\alpha(\mathcal{Q}, \Lambda(\Gamma))$. □

In order to finish the proof of Theorem 5.8, we first observe that by the partition of unity argument, it suffices to prove the claim for Φ and Ψ supported on $xB_{\epsilon_0/2}$ for $x \in \mathcal{Q}$. We note that $dm^{\text{Haar}}(pn) = dpdn$ where dp is a left Haar measure on P , and hence

$$\int_{\Gamma \backslash G} \Psi(xa_t)\Phi(x)dm^{\text{Haar}}(x) = \int_{xp \in zP_{\epsilon_0}} \int_{xpN_{\epsilon_0}} \Psi(xpna_t)\Phi(xpn)dndp.$$

Hence applying Propositions 5.9 and 5.10 for each $y = xp \in xP_{\epsilon_0}$, we deduce that

$$\begin{aligned} & e^{(1-\delta)t} \int_{\Gamma \backslash G} \Psi(xa_t)\Phi(x)dm^{\text{Haar}}(x) \\ &= \int_{xp \in xP_{\epsilon_0}} \left(\frac{m^{\text{BR}}(\Psi)\mu_{xpN}^{\text{PS}}(\Phi|_{xpN_{\epsilon_0}})}{|m^{\text{BMS}}|} + O(c_\Gamma \|\Psi\|_{C^1} \|\Phi|_{xpN_{\epsilon_0}}\|_{C^1} e^{-\eta_1 t/2}) \right) dp \\ &= \frac{m^{\text{BR}}(\Psi)m^{\text{BR}*}(\Phi)}{|m^{\text{BMS}}|} + O(c_\Gamma \|\Psi\|_{C^1} \|\Phi\|_{C^1} e^{-\eta_1 t/2}), \end{aligned}$$

where the implied constant depends only on \mathcal{Q} and $\alpha(\mathcal{Q}, \Lambda(\Gamma))$. This finishes the proof.

6. ZERO-FREE REGION OF THE SELBERG ZETA FUNCTIONS

Let $\Gamma < \text{SL}_2(\mathbb{Z})$ be as in Theorem 1.1. In [26, 27], it was shown that Theorem 1.1 implies the following theorem.

Theorem 6.1. *There exist $C' > 0$ and $\epsilon_0 > 0$ such that for all square free $q \geq 1$ with $(q, q_0) = 1$,*

(1)

$$\begin{aligned} \mathcal{P}_q(T) &:= \#\{C : \text{primitive closed geodesic in } \Gamma(q) \backslash \text{PSL}_2(\mathbb{R}) \text{ with } \ell(C) < T\} \\ &= \text{li}(e^{\delta T}) + O(q^{C'} e^{(\delta - \epsilon_0)T}), \end{aligned}$$

where $\text{li}(x) = \int_2^x \frac{dx}{\log x}$ and $\ell(C)$ is the length of C ;

(2) for any $z, w \in \mathbb{H}^2$,

$$\begin{aligned} N_q(T; z, w) &:= \#\{\gamma \in \Gamma(q) : d(z, \gamma w) \leq T\} \\ &= C_q(z, w)e^{\delta T} + O(q^{C'} e^{(\delta - \epsilon_0)T}) \end{aligned}$$

for some constant $C_q(z, w) > 0$.

Proof of Theorem 1.3. We use the well-known relation between the Poincaré series and the leading term for the resolvent of the Laplacian $R_q(s) = (\Delta_q - s(1 - s))^{-1}$. More precisely, there is a decomposition, valid on $\Re(s) > \delta$, of the resolvent as

$$(6.1) \quad R_q(s) = f(s)P_q(s) + K_q(s),$$

where $P_q(s)$ is the integral operator with kernel

$$P_q(s, z, w) := \sum_{\gamma \in \Gamma(q)} e^{-d(z, \gamma w)s} = s \int_0^\infty e^{-st} N_q(t; z, w) dt,$$

where $K_q(s)$ is holomorphic on $\Re(s) > \delta - 1$, and f is a ratio of Gamma functions holomorphic on $\Re(s) > 0$ (see [20, Proposition 2.2] and its proof). Applying the estimates on $N_q(t; z, w)$ from Theorem 6.1 (2), we see that the right hand side of (6.1) has an analytic extension to the half plane $\Re(s) > \delta - \epsilon_0$ (with ϵ_0 as in Theorem 6.1) except for a simple pole at $s = \delta$. □

ACKNOWLEDGMENTS

The authors are grateful to the referee for helpful remarks on the paper, especially for providing an alternative succinct argument in the deduction of Theorem 1.3 from Theorem 1.1.

REFERENCES

- [1] Artur Avila, Sébastien Gouëzel, and Jean-Christophe Yoccoz, *Exponential mixing for the Teichmüller flow*, Publ. Math. Inst. Hautes Études Sci. **104** (2006), 143–211, DOI 10.1007/s10240-006-0001-5. MR2264836 (2007j:37049)
- [2] Martine Babilot, *On the mixing property for hyperbolic systems*, Israel J. Math. **129** (2002), 61–76, DOI 10.1007/BF02773153. MR1910932 (2003g:37008)
- [3] Rufus Bowen and David Ruelle, *The ergodic theory of Axiom A flows*, Invent. Math. **29** (1975), no. 3, 181–202. MR0380889 (52 #1786)
- [4] David Borthwick, *Spectral theory of infinite-area hyperbolic surfaces*, Progress in Mathematics, vol. 256, Birkhäuser Boston, Inc., Boston, MA, 2007. MR2344504 (2008h:58056)
- [5] Rufus Bowen, *Markov partitions for Axiom A diffeomorphisms*, Amer. J. Math. **92** (1970), 725–747. MR0277003 (43 #2740)
- [6] B. H. Bowditch, *Geometrical finiteness with variable negative curvature*, Duke Math. J. **77** (1995), no. 1, 229–274, DOI 10.1215/S0012-7094-95-07709-6. MR1317633 (96b:53056)
- [7] Jean Bourgain and Alex Gamburd, *Uniform expansion bounds for Cayley graphs of $SL_2(\mathbb{F}_p)$* , Ann. of Math. (2) **167** (2008), no. 2, 625–642, DOI 10.4007/annals.2008.167.625. MR2415383 (2010b:20070)
- [8] Jean Bourgain, Alex Gamburd, and Peter Sarnak, *Generalization of Selberg’s $\frac{3}{16}$ theorem and affine sieve*, Acta Math. **207** (2011), no. 2, 255–290, DOI 10.1007/s11511-012-0070-x. MR2892611
- [9] Jean Bourgain, Alex Gamburd, and Peter Sarnak, *Affine linear sieve, expanders, and sum-product*, Invent. Math. **179** (2010), no. 3, 559–644, DOI 10.1007/s00222-009-0225-3. MR2587341 (2011d:11018)
- [10] Jean Bourgain and Alex Kontorovich, *On Zaremba’s conjecture*, Ann. of Math. (2) **180** (2014), no. 1, 137–196, DOI 10.4007/annals.2014.180.1.3. MR3194813
- [11] Jean Bourgain, Alex Kontorovich, and Michael Magee, *Thermodynamic expansion to arbitrary moduli*. Preprint.
- [12] Jean Bourgain, Alex Kontorovich, and Peter Sarnak, *Sector estimates for hyperbolic isometries*, Geom. Funct. Anal. **20** (2010), no. 5, 1175–1200, DOI 10.1007/s00039-010-0092-5. MR2746950
- [13] Jack Button, *All Fuchsian Schottky groups are classical Schottky groups*, The Epstein birthday schrift, Geom. Topol. Monogr., vol. 1, Geom. Topol. Publ., Coventry, 1998, pp. 117–125 (electronic), DOI 10.2140/gtm.1998.1.117. MR1668339 (2000e:20078)
- [14] Jean Bourgain and Péter P. Varjú, *Expansion in $SL_d(\mathbf{Z}/q\mathbf{Z})$, q arbitrary*, Invent. Math. **188** (2012), no. 1, 151–173, DOI 10.1007/s00222-011-0345-4. MR2897695
- [15] N. Chernov, *Invariant measures for hyperbolic dynamical systems*, Handbook of dynamical systems, Vol. 1A, North-Holland, Amsterdam, 2002, pp. 321–407, DOI 10.1016/S1874-575X(02)80006-6. MR1928521 (2003g:37047)
- [16] Yves Colin De Verdière, *Théorie spectrale des surfaces de Riemann d’aire infinie*, Astérisque **132** (1985), 259–275.
- [17] Dmitry Dolgopyat, *On decay of correlations in Anosov flows*, Ann. of Math. (2) **147** (1998), no. 2, 357–390, DOI 10.2307/121012. MR1626749 (99g:58073)
- [18] Boris Hasselblatt and Anatole Katok (eds.), *Handbook of Dynamical Systems, 1A*, North-Holland, Amsterdam, 2002.
- [19] A. Salehi Golsefidy and Péter P. Varjú, *Expansion in perfect groups*, Geom. Funct. Anal. **22** (2012), no. 6, 1832–1891, DOI 10.1007/s00039-012-0190-7. MR3000503
- [20] Colin Guillarmou and Frédéric Naud, *Wave decay on convex co-compact hyperbolic manifolds*, Comm. Math. Phys. **287** (2009), no. 2, 489–511, DOI 10.1007/s00220-008-0706-z. MR2481747 (2009m:58060)

- [21] Laurent Guillopé, Kevin K. Lin, and Maciej Zworski, *The Selberg zeta function for convex co-compact Schottky groups*, Comm. Math. Phys. **245** (2004), no. 1, 149–176, DOI 10.1007/s00220-003-1007-1. MR2036371 (2005f:11193)
- [22] Min Lee and Hee Oh, *Effective circle count for Apollonian packings and closed horospheres*, Geom. Funct. Anal. **23** (2013), no. 2, 580–621, DOI 10.1007/s00039-013-0217-8. MR3053757
- [23] Peter D. Lax and Ralph S. Phillips, *The asymptotic distribution of lattice points in Euclidean and non-Euclidean spaces*, J. Funct. Anal. **46** (1982), no. 3, 280–350, DOI 10.1016/0022-1236(82)90050-7. MR661875 (83j:10057)
- [24] Michael Magee, Hee Oh, and Dale Winter, *Expanding maps and continued fractions*, available at [arXiv:1412.4284](https://arxiv.org/abs/1412.4284).
- [25] Rafe R. Mazzeo and Richard B. Melrose, *Meromorphic extension of the resolvent on complete spaces with asymptotically constant negative curvature*, J. Funct. Anal. **75** (1987), no. 2, 260–310, DOI 10.1016/0022-1236(87)90097-8. MR916753 (89c:58133)
- [26] Gregory Margulis, Amir Mohammadi, and Hee Oh, *Closed geodesics and holonomies for Kleinian manifolds*, Geom. Funct. Anal. **24** (2014), no. 5, 1608–1636, DOI 10.1007/s00039-014-0299-y. MR3261636
- [27] Amir Mohammadi and Hee Oh, *Matrix coefficients, counting and primes for orbits of geometrically finite groups*, J. Eur. Math. Soc. (JEMS) **17** (2015), no. 4, 837–897, DOI 10.4171/JEMS/520. MR3336838
- [28] Frédéric Naud, *Expanding maps on Cantor sets and analytic continuation of zeta functions* (English, with English and French summaries), Ann. Sci. Éc. Norm. Supér. (4) **38** (2005), no. 1, 116–153, DOI 10.1016/j.ansens.2004.11.002. MR2136484 (2006e:37033)
- [29] Hee Oh and Nimish A. Shah, *Equidistribution and counting for orbits of geometrically finite hyperbolic groups*, J. Amer. Math. Soc. **26** (2013), no. 2, 511–562, DOI 10.1090/S0894-0347-2012-00749-8. MR3011420
- [30] William Parry and Mark Pollicott, *Zeta functions and the periodic orbit structure of hyperbolic dynamics* (English, with French summary), Astérisque **187-188** (1990), 268. MR1085356 (92f:58141)
- [31] S. J. Patterson, *The limit set of a Fuchsian group*, Acta Math. **136** (1976), no. 3-4, 241–273. MR0450547 (56 #8841)
- [32] S. J. Patterson, *On a lattice-point problem in hyperbolic space and related questions in spectral theory*, Ark. Mat. **26** (1988), no. 1, 167–172, DOI 10.1007/BF02386116. MR948288 (89g:11093)
- [33] Mark Pollicott, *On the rate of mixing of Axiom A flows*, Invent. Math. **81** (1985), no. 3, 413–426, DOI 10.1007/BF01388579. MR807065 (87i:58148)
- [34] Vesselin Pektov and Luchezar Stoyanov, *Spectral estimates for Ruelle transfer operators with two parameters and applications*, available at [arXiv:1409.0721](https://arxiv.org/abs/1409.0721).
- [35] Thomas Roblin, *Ergodicité et équidistribution en courbure négative*, Mém. Soc. Math. Fr. (N.S.) **95** (2003), 96p.
- [36] M. Ratner, *Markov partitions for Anosov flows on n -dimensional manifolds*, Israel J. Math. **15** (1973), 92–114. MR0339282 (49 #4042)
- [37] Daniel J. Rudolph, *Ergodic behaviour of Sullivan’s geometric measure on a geometrically finite hyperbolic manifold*, Ergodic Theory Dynam. Systems **2** (1982), no. 3-4, 491–512 (1983), DOI 10.1017/S0143385700001735. MR721736 (85i:58101)
- [38] Barbara Schapira, *Equidistribution of the horocycles of a geometrically finite surface*, Int. Math. Res. Not. IMRN **40** (2005), 2447–2471, DOI 10.1155/IMRN.2005.2447. MR2180113 (2006i:37073)
- [39] Luchezar Stoyanov, *Spectra of Ruelle transfer operators for axiom A flows*, Nonlinearity **24** (2011), no. 4, 1089–1120, DOI 10.1088/0951-7715/24/4/005. MR2776112 (2012f:37060)
- [40] Luchezar Stoyanov, *On the Ruelle-Perron-Frobenius theorem*, Asymptot. Anal. **43** (2005), no. 1-2, 131–150. MR2148129 (2007d:37023)
- [41] Dennis Sullivan, *The density at infinity of a discrete group of hyperbolic motions*, Inst. Hautes Études Sci. Publ. Math. **50** (1979), 171–202. MR556586 (81b:58031)
- [42] Dennis Sullivan, *Entropy, Hausdorff measures old and new, and limit sets of geometrically finite Kleinian groups*, Acta Math. **153** (1984), no. 3-4, 259–277, DOI 10.1007/BF02392379. MR766265 (86c:58093)

- [43] Ilya Vinogradov, *Effective bisector estimate with application to Apollonian circle packings*, Int. Math. Res. Not. IMRN **12** (2014), 3217–3262. MR3217660
- [44] David Vernon Widder, *The Laplace Transform*, Princeton Mathematical Series, v. 6, Princeton University Press, Princeton, NJ, 1941. MR0005923 (3,232d)

MATHEMATICS DEPARTMENT, YALE UNIVERSITY, NEW HAVEN, CONNECTICUT 06511 AND KOREA INSTITUTE FOR ADVANCED STUDY, SEOUL, KOREA
E-mail address: `hee.oh@yale.edu`

DEPARTMENT OF MATHEMATICS, BROWN UNIVERSITY, PROVIDENCE, RHODE ISLAND 02906
E-mail address: `dale_winter@brown.edu`