

# SPECTRAL GAP AND EXPONENTIAL MIXING ON GEOMETRICALLY FINITE HYPERBOLIC MANIFOLDS

SAM EDWARDS AND HEE OH

ABSTRACT. Let  $\mathcal{M} = \Gamma \backslash \mathbb{H}^{d+1}$  be a geometrically finite hyperbolic manifold with critical exponent exceeding  $d/2$ . We obtain a precise asymptotic expansion of the matrix coefficients for the geodesic flow in  $L^2(\Gamma^1(\mathcal{M}))$ , with exponential error term essentially as good as the one given by the spectral gap for the Laplace operator on  $L^2(\mathcal{M})$  due to Lax and Phillips. Combined with the work of Bourgain, Gamburd, and Sarnak and its generalization by Golsefidy and Varju on expanders, this implies *uniform* exponential mixing for congruence covers of  $\mathcal{M}$  when  $\Gamma$  is a Zariski dense subgroup contained in an arithmetic subgroup of  $\mathrm{SO}^\circ(d, 1)$ .

## CONTENTS

1.	Introduction	1
2.	Preliminaries	6
3.	Complementary series representations and Eisenstein integrals	10
4.	Asymptotic expansions of matrix coefficients	18
5.	Leading term for Sobolev functions	28
6.	Proof of Theorem 1.1	30
	References	33

## 1. INTRODUCTION

Let  $\mathcal{M} = \Gamma \backslash \mathbb{H}^{d+1}$  be a geometrically finite hyperbolic manifold, where  $\Gamma$  is a discrete subgroup of  $G = \mathrm{SO}^\circ(d + 1, 1)$ . We suppose that the critical exponent  $\delta$  of  $\Gamma$  is strictly bigger than  $d/2$ .

Denoting by  $\Delta$  the negative of the Laplace operator on  $\mathcal{M}$ , the bottom eigenvalue of  $\Delta$  on  $L^2(\mathcal{M})$  is known to be simple and given as  $\lambda_0 = \delta(d - \delta)$  by Patterson [30] and Sullivan [40]. Lax and Phillips [24] showed that there are only finitely many eigenvalues of  $\Delta$  on  $L^2(\mathcal{M})$  in the interval  $[\lambda_0, \frac{d^2}{4})$ . We denote by  $\lambda_1$  the smallest eigenvalue of  $\Delta$  in  $(\lambda_0, \frac{d^2}{4})$ , and define  $s_1 =$

---

S. E. was supported by postdoctoral scholarship 2017.0391 from the Knut and Alice Wallenberg Foundation and H. O. was supported by NSF grants.

$\frac{d}{2} + \sqrt{(\frac{d}{2})^2 - \lambda_1}$ . Then  $s_1$  is the unique number in  $(\frac{d}{2}, \delta)$  such that

$$0 < \delta(d - \delta) = \lambda_0 < \lambda_1 = s_1(d - s_1) < d^2/4.$$

If there is no eigenvalue in the open interval  $(\lambda_0, \frac{d^2}{4})$ , we set  $\lambda_1 = d^2/4$ , i.e.,  $s_1 = d/2$ .

Denote by  $\mathcal{G}^t$  the geodesic flow on the unit tangent bundle  $T^1(\mathcal{M})$  and by  $dx$  the Liouville measure on  $T^1(\mathcal{M})$ ; this is a  $\mathcal{G}^t$ -invariant Borel measure which is infinite when  $\delta < d$ .

The main result of this paper is that the Lax-Philips spectral gap, that is,  $\lambda_1 - \lambda_0$ , or equivalently  $\delta - s_1$ , controls *rather precisely* the exponential convergence rate of the correlation function for the geodesic flow  $\mathcal{G}^t$  acting on  $L^2(T^1(\mathcal{M}), dx)$ :

**Theorem 1.1.** *Let  $\delta > d/2$ , and set  $\eta := \min\{\delta - s_1, 1\}$ . There exists  $m > d(d+1)/2$  such that for any  $\varepsilon > 0$  and functions  $\psi_1, \psi_2$  on  $T^1(\mathcal{M})$  with  $\mathcal{S}^m(\psi_1), \mathcal{S}^m(\psi_2) < \infty$ , we have, as  $t \rightarrow +\infty$ ,*

$$e^{(d-\delta)t} \int_{T^1(\mathcal{M})} \psi_1(\mathcal{G}^t(x)) \psi_2(x) dx = \frac{1}{m^{\text{BMS}}(T^1(\mathcal{M}))} m^{\text{BR}}(\psi_1) m^{\text{BR}^*}(\psi_2) + O_\varepsilon \left( e^{(-\eta+\varepsilon)t} \mathcal{S}^m(\psi_1) \mathcal{S}^m(\psi_2) \right),$$

where

- $m^{\text{BMS}}$  is the Bowen-Margulis-Sullivan measure on  $T^1(\mathcal{M})$ ;
- $m^{\text{BR}}$  and  $m^{\text{BR}^*}$  are, respectively, the unstable and stable Burger-Roblin measures on  $T^1(\mathcal{M})$ , which are defined compatibly with the choice of  $dx$  and  $dm^{\text{BMS}}$ ;
- $\mathcal{S}^m(\psi_i)$  denotes the  $L^2$ -Sobolev norm of  $\psi_i$  of degree  $m$  and the implied constant depends only on  $\varepsilon$ .

We note that  $|m^{\text{BR}}(\psi_i)| < \infty$  when  $\mathcal{S}^m(\psi_i) < \infty$  (see Lemma 2.2), and hence the main term above is well-defined. For  $\psi_i$  compactly supported, the asymptotic formula (without error term) holds for any  $\delta > 0$ , as was obtained by Roblin [33].

**Remark 1.1.** (1) When  $\mathcal{M}$  has finite volume, i.e., when  $\delta = d$ , exponential mixing of the geodesic flow is a classical result due to Ratner for  $d = 1$  and to Moore [28] for general  $d \geq 1$ . Combining Moore's proof with Shalom's trick [37, proof of Theorem 2.1] yields a rate of mixing  $\eta = \min\{d - s_1, 1\}$  as in Theorem 1.1 above. This proof makes explicit use of the fact [16] that there are no non-spherical complementary series representations  $\mathcal{U}(v, s)$  with  $s > d - 1$  (see below for notation). Our proof does not require this statement. However, combining the aforementioned fact with our proof improves the rate of mixing for the geodesic flow on finite-volume hyperbolic manifolds to  $\eta = \min\{d - s_1, 2\}$ . In addition, if there is no non-spherical complementary series representation  $\mathcal{U}(v, s)$  with  $s > s_1 + 1$  appearing

in  $L^2(\Gamma \backslash G)$ , then  $\eta$  can be taken to be  $d - s_1$  (see Remark 6.5 for details).

- (2) For  $d = 1$  and 2, Theorem 1.1 can be deduced from [7] and [43], respectively. For a general  $d \geq 1$ , the case of  $\delta > \max\{d - 1, d/2\}$  was obtained in [27] for *some*  $\eta > 0$  which is not explicit.
- (3) The main novelties of Theorem 1.1 are that it addresses *all* geometrically finite groups with  $\delta > d/2$  (even those with cusps), and that it gives an optimal value of  $\eta$  with respect to the dependency on the spectral gap  $\delta - s_1$ .
- (4) The order  $m$  of Sobolev norm required is in principle completely computable; it may be chosen independently of  $\Gamma$ , and satisfies  $m = O(d^2)$ .

The BMS measure  $m^{\text{BMS}}$  is known to be the unique measure of maximal entropy (which is  $\delta$ ) for the geodesic flow ([40], [31]). Babillot showed that the geodesic flow is mixing with respect to  $m^{\text{BMS}}$  [2]. Theorem 1.1 is known to imply the following exponential mixing for the BMS measure (see [27, Theorem 1.6] for compactly supported functions, and [18, Theorem 1.9] for general bounded functions):

**Theorem 1.2.** *There exist  $\beta > 0$  (explicitly computable, depending only on  $\eta$  in Theorem 1.1) and  $m > d(d + 1)/2$  such that for all bounded functions  $\psi_1, \psi_2$  on  $\mathbb{T}^1(\mathcal{M})$  supported on the one-neighbourhood of  $\text{supp}(m^{\text{BMS}})$ , we have, as  $t \rightarrow \infty$ ,*

$$\int_{\mathbb{T}^1(\mathcal{M})} \psi_1(\mathcal{G}^t(x)) \psi_2(x) dm^{\text{BMS}}(x) = \frac{1}{m^{\text{BMS}}(\mathbb{T}^1(\mathcal{M}))} m^{\text{BMS}}(\psi_1) m^{\text{BMS}}(\psi_2) + O\left(e^{-\beta t} \|\psi_1\|_{C^m} \|\psi_2\|_{C^m}\right)$$

where  $\|\psi_i\|_{C^m}$  denotes the  $C^m$ -norm of  $\psi_i$ .

**Remark 1.2.** When  $\Gamma$  is convex cocompact, Theorem 1.2, for some  $\beta > 0$  (which is not explicit), follows from the work of Stoyanov [38], which is based on symbolic dynamics and Dolgopyat operators (see also [35]). This result in turn implies Theorem 1.1 with implicit  $\eta > 0$  (see [29], [35]).

Let  $\Gamma$  be a Zariski dense subgroup of an arithmetic subgroup  $G(\mathbb{Z})$  of  $G$ . Denote by  $\Gamma_q$  the congruence subgroup of  $\Gamma$  of level  $q$ :  $\Gamma_q = \{\gamma \in \Gamma : \gamma = e \pmod{q}\}$ . When  $\pi^{-1}(\Gamma)$  satisfies the strong approximation property for the spin covering map  $\pi : \text{Spin}(d + 1, 1) \rightarrow G$ , the work of Bourgain-Gamburd-Sarnak [6] and its generalization by Golsefidy-Varju [41] on expanders imply that there exists a finite set  $S$  of primes such that the family  $\mathcal{F} := \{\Gamma_q : q \text{ is square-free with no factors in } S\}$  has a uniform spectral gap in the sense that

$$\inf_{\Gamma_q \in \mathcal{F}} \{\delta - s_1(q)\} > 0, \tag{1.3}$$

where  $s_1(q)(d - s_1(q))$  is the second smallest eigenvalue of the Laplacian  $\Delta$  on  $L^2(\Gamma_q \backslash \mathbb{H}^{d+1})$ <sup>1</sup>. This uses the transfer property from the combinatorial spectral gap to the archimedean spectral gap due to [5] (see also [19]). For certain families of subgroups the uniform spectral gap (1.3) may be explicitly bounded from below, see [8, 9, 13, 25, 36].

We then have:

**Corollary 1.3.** *In Theorems 1.1 and 1.2, the exponents  $\eta$  and  $\beta$  can be chosen uniformly over all congruence covers  $\Gamma^1(\Gamma_q \backslash \mathbb{H}^{d+1})$ ,  $\Gamma_q \in \mathcal{F}$ .*

When  $\Gamma$  is convex cocompact, Corollary 1.3 was obtained by the second-named author and Winter [29] for  $d = 1$  and by Sarkar [34] for a general  $d \geq 1$ , combining Dolgopyat's methods with the expander theory ([6], [41]).

Theorems 1.1-1.2 and Corollary 1.3 are known to have many immediate applications in number theory and geometry. To name a few, see ([11], [12], [3], [5], [27]) for effective counting and affine sieve, [26] for the prime geodesic theorem, and [18] for shrinking target problems.

Fix  $o \in \mathbb{H}^{d+1}$  and  $v_o \in \Gamma^1(\mathbb{H}^{d+1})$  based at  $o$ . Setting  $K = \text{Stab}_G(o) \simeq \text{SO}(d+1)$  and  $M = \text{Stab}_G(v_o) \simeq \text{SO}(d)$ , we can identify  $\mathcal{M} = \Gamma \backslash G/K$  and  $\Gamma^1(\mathcal{M}) = \Gamma \backslash G/M$ . We let  $\{a_t\}$  denote the one-parameter diagonalizable subgroup of  $G$  commuting with  $M$  whose right translation on  $\Gamma \backslash G/M$  corresponds to the geodesic flow  $\mathcal{G}^t$  on  $\Gamma^1(\mathcal{M})$ . As  $L^2(\Gamma^1(\mathcal{M}))$  can be identified with the space of  $M$ -invariant functions in  $L^2(\Gamma \backslash G)$ , Theorem 1.1 amounts to understanding the asymptotic behaviour of the matrix coefficients  $\langle a_t \psi_1, \psi_2 \rangle$  for  $M$ -invariant functions  $\psi_1, \psi_2 \in L^2(\Gamma \backslash G)$ .

The non-tempered part of the unitary dual  $\hat{G}$  consists of the trivial representation and the complementary series representations  $\mathcal{U}(v, s)$ , parameterized by a representation  $v$  in the unitary dual  $\hat{M}$  and a real number  $s \in \mathcal{I}_v$ , where  $\mathcal{I}_v \subset (d/2, d)$  is an interval depending on  $v$ . In this parameterization, the complementary series representation  $\mathcal{U}(v, s)$  is spherical if and only if  $v \in \hat{M}$  is trivial, and for  $v$  non-trivial,  $\mathcal{I}_v$  is contained in  $(d/2, d - 1)$  [16]; this was the main reason for the hypothesis  $\delta > d - 1$  in [27].

Our main work for the proof of Theorem 1.1 lies in the detailed analysis of the behavior of matrix coefficients of the complementary series representations  $\mathcal{U}(v, s)$ . We remark that Harish-Chandra's work ([46], [47]) does not give an asymptotic expansion for the matrix coefficients of  $\mathcal{U}(v, s)$  for all  $s$ , as it excludes finitely many (unknown) parameters  $s$ . Even for those  $\mathcal{U}(v, s)$  for which Harish-Chandra's expansion is given, it is hard to use his expansion directly, as it relies on various parameters and recursive formulas.

Denoting by  $\langle \cdot, \cdot \rangle_{\mathcal{U}(v, s)}$  the *unitary* inner product, for each  $m \in \mathbb{N}$  we define a Sobolev norm  $\|\cdot\|_{\mathcal{S}^m(v, s)}$  on  $\mathcal{U}(v, s)$  as in (2.3), and let  $\mathcal{S}^m(v, s)$  denote the space of vectors in  $\mathcal{U}(v, s)$  with finite  $\|\cdot\|_{\mathcal{S}^m(v, s)}$ -norm.

---

<sup>1</sup>In view of the recent preprint [15] and an upcoming work by He and de Saxcé, the square-free condition in  $\mathcal{F}$  may be removed.

Write  $\mathcal{U}(v, s) = \bigoplus_{\tau \in \hat{K}} \mathcal{U}(v, s)_\tau$  for the decomposition into different  $K$ -types and for  $d/2 < s < d$ , set

$$\eta_s := \min\{2s - d, 1\} > 0.$$

We show the following *concrete* asymptotic expansion of matrix coefficients with an optimal rate.

**Theorem 1.4.** *There exists  $m \in \mathbb{N}$  such that for any complementary series representation  $\mathcal{U}(v, s)$  containing a non-trivial  $M$ -invariant vector, for all  $\mathbf{u}, \mathbf{v} \in \mathcal{S}^m(v, s)$  and  $t \geq 0$ ,*

$$\begin{aligned} \langle \mathcal{U}(v, s)(a_t)\mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)} &= e^{(s-d)t} \left( \sum_{\tau_1, \tau_2 \in \hat{K}} \langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{P}_{\tau_1} \mathbf{u}, \mathbf{P}_{\tau_2} \mathbf{v} \rangle_{\mathcal{U}(v, s)} \right) \\ &\quad + O_s(e^{(s-d-\eta_s)t} \|\mathbf{u}\|_{\mathcal{S}^m(v, s)} \|\mathbf{v}\|_{\mathcal{S}^m(v, s)}) \end{aligned}$$

and the sum

$$\sum_{\tau_1, \tau_2 \in \hat{K}} \langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{P}_{\tau_1} \mathbf{u}, \mathbf{P}_{\tau_2} \mathbf{v} \rangle_{\mathcal{U}(v, s)} \quad (1.4)$$

converges absolutely. Here  $T_{\tau_1}^{\tau_2} : \mathcal{U}(v, s)_{\tau_1} \rightarrow \mathcal{U}(v, s)_{\tau_2}$  is given by (3.5),  $C_+(s)$  is the Harish-Chandra  $c$ -function given in (4.1), and  $\mathbf{P}_\tau$  is the orthogonal projection onto the  $K$ -type  $\tau$ . Moreover, the implied constant is uniformly bounded over  $s$  in compact subsets of the interval  $\mathcal{I}_v$ .

One of the key observations made in this paper is that if  $\mathcal{U}(v, s)$  is non-spherical and  $\mathbf{u} \in \mathcal{U}(v, s)$  is  $M$ -invariant, then the main term (1.4) vanishes (Corollary 3.5): for all  $\tau_1, \tau_2 \in \hat{K}$ ,

$$T_{\tau_1}^{\tau_2} C_+(s) \mathbf{P}_{\tau_1} \mathbf{u} = 0.$$

Furthermore, there is additional uniformity with respect to the  $s$ -variable when considering only  $M$ -invariant vectors.

**Corollary 1.5.** *There exists  $m \in \mathbb{N}$  such that if  $\mathcal{U}(v, s)$  is non-spherical, then for all  $M$ -invariant vectors  $\mathbf{u}, \mathbf{v} \in \mathcal{U}(v, s)$ ,*

$$|\langle \mathcal{U}(v, s)(a_t)\mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)}| \ll_s e^{(s-d-\eta_s)t} \|\mathbf{u}\|_{\mathcal{S}^m(v, s)} \|\mathbf{v}\|_{\mathcal{S}^m(v, s)}.$$

Moreover, the implied constant is uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ .

**Organization.** We start by recalling in Section 2 some key background facts and notation. In Section 3, we define the models of the complementary series that we will work in and the Eisenstein integrals that are key to understanding their matrix coefficients. The main technical work is done in Section 4; establishing the asymptotic expansion and bounds for matrix coefficients in complementary series representations. Section 5 is devoted to extending Roblin's mixing result [33, Theorem 3.4] from compactly supported functions to arbitrary Sobolev functions. Finally, the proof of Theorem 1.1 is given in Section 6, the main part of which consists of decomposing the

regular representation of  $G$  on  $L^2(\Gamma \backslash G)$  into irreducible representations and applying the bounds obtained in Section 4 to each individual component.

## 2. PRELIMINARIES

Let  $d \geq 1$  and  $G = \mathrm{SO}^\circ(d+1, 1)$  be the group of orientation-preserving isometries of  $\mathbb{H}^{d+1}$ . Let  $\Gamma < G$  be a torsion-free discrete subgroup of  $G$ , and let  $\mathcal{M} = \Gamma \backslash \mathbb{H}^{d+1}$  be the associated hyperbolic manifold.

**2.1. Structure and subgroups of  $G$ .** Let  $K = \mathrm{SO}(d+1) < G$  be a maximal compact subgroup, and let  $A = \{a_t : t \in \mathbb{R}\}$  be a one-parameter diagonalizable subgroup and  $M = \mathrm{SO}(d)$  the centralizer of  $A$  in  $K$ . We can identify  $\mathcal{M}$  with  $\Gamma \backslash G/K$  and the unit tangent bundle  $\mathrm{T}^1(\mathcal{M})$  with  $\Gamma \backslash G/M$  in the way that the geodesic flow on  $\mathrm{T}^1(\mathcal{M})$  is given by the right translation action of  $a_t$  on  $\Gamma \backslash G/M$ .

We denote by  $N$  and  $\bar{N}$  the contracting and expanding horospherical subgroups, respectively; i.e.

$$N = \{g \in G : a_{-t}ga_t \rightarrow e \text{ as } t \rightarrow +\infty\};$$

$$\bar{N} = \{g \in G : a_tga_{-t} \rightarrow e \text{ as } t \rightarrow +\infty\}.$$

Then  $\bar{N} = \omega N \omega^{-1}$ , where  $\omega \in N_K(A)$  is such that  $\omega a \omega^{-1} = a^{-1}$  for all  $a \in A$ .

We note that  $N$  and  $\bar{N}$  are both abelian subgroups, isomorphic to  $\mathbb{R}^d$  via the logarithm map, which are normalized by  $AM$ . Under this isomorphism  $N \simeq \mathbb{R}^d$ , conjugation by an element  $m \in M = \mathrm{SO}(d)$  is an isometry and conjugation by  $a_t$  corresponds to scaling by  $e^t$  on  $N$ , and  $e^{-t}$  on  $\bar{N}$ . As usual, the Lie algebras of  $G, K, A, N, \bar{N}$  are denoted by  $\mathfrak{g}, \mathfrak{k}, \mathfrak{a}, \mathfrak{n}$  and  $\bar{\mathfrak{n}}$ , respectively. This gives

$$\mathrm{Ad}(a_t)|_{\mathfrak{n}} = e^t \times \mathrm{Id} \quad \text{and} \quad \mathrm{Ad}(a_t)|_{\bar{\mathfrak{n}}} = e^{-t} \times \mathrm{Id}.$$

We have an Iwasawa decomposition  $G = KAN$ , where the product map  $K \times A \times N \rightarrow G$  is a diffeomorphism and the projection to each individual factor is a smooth map [21, Chapter VI.4]. For  $g \in G$ , we write

$$g = \kappa(g) \exp(H(g))n_g,$$

where  $\kappa(g) \in K$ ,  $H(g) \in \mathfrak{a}$ , and  $n_g \in N$ . Letting  $\alpha \in \mathfrak{a}_{\mathbb{C}}^*$  be the unique element such that

$$\alpha(H(a_t)) = t,$$

the map  $s \mapsto s \cdot \alpha$  defines an identification  $\mathbb{C} \simeq \mathfrak{a}_{\mathbb{C}}^*$ . In view of this, we write

$$e^{sH(g)} = e^{s \cdot \alpha(H(g))} \quad \text{for } g \in G \text{ and } s \in \mathbb{C}.$$

**2.2. Various measures on  $T^1(\mathcal{M})$ .** Let  $\Lambda \subset \partial\mathbb{H}^{d+1}$  denote the limit set of  $\Gamma$ . We assume that  $\Gamma$  is non-elementary, or equivalently,  $\Lambda$  has at least 3 points. Throughout the paper, we assume that

*$\Gamma$  is geometrically finite,*

that is, the unit neighborhood of the convex core of  $\mathcal{M}$ , given by  $\Gamma \backslash \text{hull}(\Lambda)$ , has finite Riemannian volume. We let  $\delta$  denote the critical exponent of  $\Gamma$ , which is known to be equal to the Hausdorff dimension of  $\Lambda$ . We remark that  $\delta = d$  if and only if  $\partial\mathbb{H}^n = \Lambda$  if and only if  $\Gamma$  is a lattice in  $G$ .

We fix  $o \in \mathbb{H}^{d+1}$  whose stabilizer is given by  $K$ , and  $\nu_o \in T_o(\mathbb{H}^{d+1})$  the unit vector whose stabilizer is  $M$ . Let  $\nu_o$  be the Patterson-Sullivan measure on  $\Lambda$  which is unique up to a constant multiple; this is characterized by the condition that

$$\gamma^* \nu_o = |\gamma'|^\delta \cdot \nu_o$$

for all  $\gamma \in \Gamma$ , where  $|\gamma'|$  denotes the derivative of  $\gamma$  in the spherical metric on  $\partial\mathbb{H}^{d+1}$  with respect to  $o$ . We also let  $m_o$  denote the  $K$ -invariant probability measure on  $\partial\mathbb{H}^{d+1}$ .

Let  $\pi : T^1(\mathbb{H}^{d+1}) \rightarrow \mathbb{H}^{d+1}$  be the base point projection. For  $u \in T^1(\mathbb{H}^{d+1})$ , we denote by  $u^\pm \in \partial\mathbb{H}^{d+1}$  the forward and the backward endpoints of the geodesic determined by  $u$ . Consider the Hopf parameterization of  $T^1(\mathbb{H}^{d+1})$  given by:

$$u \mapsto (u^+, u^-, s = \beta_{u^-}(o, \pi(u))),$$

where  $\beta$  denotes the Busemann function. Using the Hopf coordinates, the Bowen-Margulis-Sullivan measure  $m^{\text{BMS}}$ , the Liouville measure  $du$ , and the Burger-Roblin measure  $m^{\text{BR}}$  on  $T^1(\mathbb{H}^{d+1})$  are respectively given as follows:

- (1)  $dm^{\text{BMS}}(u) = e^{\delta\beta_{u^+}(o, \pi(u))} e^{\delta\beta_{u^-}(o, \pi(u))} d\nu_o(u^+) d\nu_o(u^-) ds$ ;
- (2)  $du = dm^{\text{Liouville}}(u) = e^{d\beta_{u^+}(o, \pi(u))} e^{d\beta_{u^-}(o, \pi(u))} dm_o(u^+) dm_o(u^-) ds$ ;
- (3)  $dm^{\text{BR}}(u) = e^{d\beta_{u^+}(o, \pi(u))} e^{\delta\beta_{u^-}(o, \pi(u))} dm_o(u^+) d\nu_o(u^-) ds$ .

The BR measure  $m^{\text{BR}}$  is a Lebesgue measure on each  $\overline{N}$ -leaf, so we call it the unstable BR measure. The stable BR-measure  $m^{\text{BR}*}$  is defined similarly to  $m^{\text{BR}}$  by exchanging the roles of  $u^+$  and  $u^-$ .

These measures are all left  $\Gamma$ -invariant, and hence induce Borel measures on  $T^1(\mathcal{M}) = \Gamma \backslash T^1(\mathbb{H}^{d+1})$  for which we use the same notation. We remark that  $m^{\text{BMS}}$  is a finite measure, invariant under the geodesic flow. We will normalize the Patterson-Sullivan measure  $\nu_o$  so that  $m^{\text{BMS}}(T^1(\mathcal{M})) = 1$ . The other three measures are infinite, unless  $\delta = d$ .

We will sometimes consider these measures as measures on  $\Gamma \backslash G$  by putting: for  $\psi \in C_c(\Gamma \backslash G)$ , and for  $\star = \text{BMS, Liouville, BR}$ ,

$$m^\star(\psi) = m^\star \left( \int_M \psi dm \right),$$

where  $dm$  denotes the Haar probability measure on  $M$ . Note that the Liouville measure, considered as a measure on  $\Gamma \backslash G$ , is  $G$ -invariant; we will denote this by  $dg$ .

**2.3. The base eigenfunction  $\phi_0$ .** Throughout the article, we assume

$$\delta > \frac{d}{2},$$

which is a necessary and sufficient condition for the existence of an eigenvalue of  $\Delta$  on  $L^2(\mathcal{M})$ . The smallest eigenvalue of  $\Delta$  on  $L^2(\mathcal{M})$  is given by  $\lambda_0 = \delta(d - \delta)$ , and is known to be simple. We denote by  $\phi_0$  the the unit eigenfunction in  $L^2(\mathcal{M})$  with eigenvalue  $\lambda_0$ , and call it the *base eigenfunction*. Up to a constant multiple,  $\phi_0$  is given by: for  $x \in \mathcal{M}$ ,

$$\phi_0(x) = \int_{\xi \in \Lambda} e^{-\delta\beta_\xi(o,x)} d\nu_o(\xi).$$

The fact that the base eigenfunction  $\phi_0$  is square-integrable when  $\delta > d/2$  is a key result, which allows the unitary representation theory of  $G$ , specifically the right translation action of  $G$  on  $L^2(\Gamma \backslash G)$ , to be used to prove dynamical results for the geodesic flow.

**Theorem 2.1** (Lax-Phillips [24]). *The intersection of the interval  $[0, \frac{d^2}{4})$  with the spectrum of  $\Delta$ , viewed as an unbounded operator on  $L^2(\mathcal{M})$ , consists of a finite set of eigenvalues  $\{\lambda_i = s_i(d - s_i)\}_{0 \leq i \leq \ell}$ , satisfying*

$$0 < \lambda_0 = \delta(d - \delta) < \lambda_1 \leq \dots \leq \lambda_\ell < \frac{d^2}{4}.$$

**2.4. The quasi-regular representation  $L^2(\Gamma \backslash G)$ .** We denote  $L^2(\Gamma \backslash G)$  the space of square-integrable functions on  $\Gamma \backslash G$  with respect to  $dg$ . The  $G$ -invariance of  $dg$  gives rise to a unitary representation  $(\rho, L^2(\Gamma \backslash G))$  of  $G$ , where  $\rho$  is given by the right translation:

$$[\rho(g)f](x) = f(xg)$$

for all  $f \in L^2(\Gamma \backslash G)$ ,  $g \in G$ ,  $x \in \Gamma \backslash G$ .

As a number of inner products on different vector spaces will show up throughout the article, we reserve now  $\langle \cdot, \cdot \rangle$  to mean the  $\rho(G)$ -invariant inner product on  $L^2(\Gamma \backslash G)$ . All other inner products will have some additional notation to distinguish them. The subspace of  $\rho(K)$ -invariant vectors in  $L^2(\Gamma \backslash G)$  is denoted  $L^2(\Gamma \backslash G)^K$ . Similarly,  $L^2(\Gamma \backslash G)^M$  denotes the subspace of  $\rho(M)$ -invariant vectors. We use  $L^2(\Gamma \backslash G)^K$  to construct subrepresentations of  $(\rho, L^2(\Gamma \backslash G))$  as follows: define

$$L^2(\Gamma \backslash G)_{\text{sph}} := \text{the closure of } \{\rho(g)f : f \in L^2(\Gamma \backslash G)^K, g \in G\}.$$

Then

$$(\rho, L^2(\Gamma \backslash G)) = (\rho, L^2(\Gamma \backslash G)_{\text{sph}}) \oplus (\rho, L^2(\Gamma \backslash G)_{\text{sph}}^\perp);$$

note that both  $L^2(\Gamma \backslash G)_{\text{sph}}$  and  $L^2(\Gamma \backslash G)_{\text{sph}}^\perp$  are  $\rho(G)$ -invariant closed subspaces. Viewing the base eigenfunction  $\phi_0$  as an element of  $L^2(\Gamma \backslash G)^K$ , we define in a similar manner

$$\mathcal{B}_\delta := \text{the closure of the span of } \{\rho(g)\phi_0 : g \in G\} \subset L^2(\Gamma \backslash G)_{\text{sph}},$$



and

$$(\rho, L^2(\Gamma \backslash G)_{\text{sph}}) = (\rho, \mathcal{B}_\delta) \oplus (\rho, \mathcal{W}),$$

$\mathcal{W}$  being the orthogonal complement of  $\mathcal{B}_\delta$  in  $L^2(\Gamma \backslash G)_{\text{sph}}$ . It will be of importance later that  $(\rho, \mathcal{B}_\delta)$  is an (irreducible) complementary series representation and that Theorem 2.1 gives a complete understanding of the complementary series representations contained in  $(\rho, \mathcal{W})$ . A useful fact (cf. [23]) that we make of in Section 5 is that for all  $f \in L^2(\Gamma \backslash G)^K$ ,

$$m^{\text{BR}}(f) = m^{\text{BR}^*}(f) = \langle f, \phi_0 \rangle. \quad (2.1)$$

Finally, the direct integral decomposition of  $(\rho, L^2(\Gamma \backslash G))$  reads

$$(\rho, L^2(\Gamma \backslash G)) \cong \int_{\mathcal{Z}}^{\oplus} (\pi_\zeta, \mathcal{H}_\zeta) d\mu_{\mathcal{Z}}(\zeta), \quad (2.2)$$

where  $(\pi_\zeta, \mathcal{H}_\zeta)$  is an irreducible unitary representation of  $G$  for  $\mu_{\mathcal{Z}}$ -a.e.  $\zeta$  (cf. [45, Corollary 14.9.5]).

**2.5. Sobolev norms on unitary representations of  $K$ .** Given a unitary representation  $(\pi, V)$  of  $K$  with invariant inner product  $(\cdot, \cdot)_V$ , a basis  $\{X_j\}$  of  $\mathfrak{k}$  and  $m \in \mathbb{N}$ , we define a Sobolev norm  $\|\cdot\|_{\mathcal{S}^m(V)}$  on  $V$  by

$$\|\mathbf{v}\|_{\mathcal{S}^m(V)}^2 := \|\mathbf{v}\|_V^2 + \sum_U \|d\pi(U)\mathbf{v}\|_V^2 \quad \text{for any } \mathbf{v} \in V, \quad (2.3)$$

where  $\|\cdot\|_V$  denotes the corresponding norm and the sum runs over all monomials in  $\{X_j\}$  of order up to  $m$ .

We set  $\mathcal{S}^m(V) := \{\mathbf{v} \in V : \|\mathbf{v}\|_{\mathcal{S}^m(V)} < \infty\}$ . Observe that different choices of the basis  $\{X_j\}$  give rise to equivalent norms, and that in the case when  $V$  is finite-dimensional, we have  $\mathcal{S}^m(V) = V$  for any  $m \geq 0$ .

Viewing  $(\rho, L^2(\Gamma \backslash G))$  as a unitary representation of  $K$ , given a function  $f \in L^2(\Gamma \backslash G)$ , we let either  $\|f\|_{\mathcal{S}^m(\Gamma \backslash G)}$  or simply  $\mathcal{S}^m(f)$  denote the norm  $\|f\|_{\mathcal{S}^m(L^2(\Gamma \backslash G))}$  as defined by (2.3) above. We denote by  $\mathcal{S}^m(\Gamma \backslash G)$  the space of all functions  $f \in L^2(\Gamma \backslash G)$  with  $\mathcal{S}^m(f) < \infty$ .

For  $f \in L^2(\Gamma \backslash G)$ ,  $m^{\text{BR}}(f)$  may be infinite in general when  $\delta < d$ . However we have the following lemma:

**Lemma 2.2.** *If  $m > \frac{(d+1)d}{4}$ , then for any  $f \in \mathcal{S}^m(\Gamma \backslash G)$ ,*

$$m^{\text{BR}}(f), m^{\text{BR}^*}(f) \ll \mathcal{S}^m(f).$$

*Proof.* Given  $f \in C(\Gamma \backslash G) \cap L^2(\Gamma \backslash G)$ , define

$$f_K(x) := \max_{k \in K} |f(xk)| \quad \text{for } x \in \Gamma \backslash G.$$

By construction,  $f_K$  is  $\rho(K)$ -invariant and  $|f(x)| \leq f_K(x)$  for all  $x \in \Gamma \backslash G$ , hence  $|m^{\text{BR}}(f)| \leq m^{\text{BR}}(f_K)$ . In view of (2.1), the Sobolev embedding theorem on the compact manifold  $K$  gives

$$|m^{\text{BR}}(f)| \leq m^{\text{BR}}(f_K) = \langle f_K, \phi_0 \rangle \leq \|f_K\| \ll \|f\|_{\mathcal{S}^m(\Gamma \backslash G)}$$

for all  $m > \frac{\dim(K)}{2}$  (cf. [1, Theorem 2.30]). Since  $\dim K = \frac{d(d+1)}{2}$ , the chain of inequalities above holds for any  $m > \frac{(d+1)d}{4}$ . The claim for  $m^{\text{BR}^*}$  follows similarly.  $\square$

### 3. COMPLEMENTARY SERIES REPRESENTATIONS AND EISENSTEIN INTEGRALS

In this section we recall the definition and  $K$ -type structure of the complementary series representations of  $G = \text{SO}^\circ(d+1, 1)$ , as well as the Eisenstein integral representations of matrix coefficients for these representations. We start by reviewing the representation theory of the special orthogonal groups  $\text{SO}(n)$ .

**3.1. Irreducible representations of  $K$  and  $M$ .** The primary reference for the first facts listed below is [4, pp. 272-277]. Let  $m \geq 1$ . The irreducible representations of  $\text{SO}(2m)$  are parameterized by  $m$ -tuples  $\tau = (\tau_1, \tau_2, \dots, \tau_m) \in \mathbb{Z}^m$  such that

$$\tau_1 \geq \tau_2 \geq \dots \geq \tau_{m-1} \geq \tau_m \geq |\tau_m|.$$

Similarly, the irreducible representations of  $\text{SO}(2m+1)$  are parameterized by  $m$ -tuples  $\tau = (\tau_1, \tau_2, \dots, \tau_m) \in \mathbb{Z}^m$  such that

$$\tau_1 \geq \tau_2 \geq \dots \geq \tau_{m-1} \geq \tau_m \geq 0.$$

We recall that all irreducible representations of  $\text{SO}(2m+1)$  and  $\text{SO}(4m)$  are self-dual, and that the dual  $\tau^*$  of an irreducible representation  $\tau = (\tau_1, \tau_2, \dots, \tau_{2m+1})$  of  $\text{SO}(4m+2)$  is given by  $\tau^* = (\tau_1, \tau_2, \dots, -\tau_{2m+1})$ ; the dual  $\tau^*$  is defined by  $\tau^*(k) = \tau(k^{-1})$ , acting on the dual (or complex conjugate) of the underlying vector space.

A key result that we will make repeated use of is the branching law for restrictions of representations of  $K$  to  $M$  which we recall (cf. [21, Chapter IX.3] or [48, pp. 377-380]). We henceforth call irreducible representations of  $K$  or  $M$  a  $K$  or  $M$ -type, respectively. Let  $\tau$  be a  $K$ -type. Then the decomposition of  $\tau|_M$  reads

$$\tau|_M = \bigoplus_{\sigma \in \hat{M}} m_{\sigma, \tau} \cdot \sigma,$$

where  $m_{\sigma, \tau} = 1$  if  $\sigma = (\sigma_1, \dots, \sigma_{\lfloor \frac{d}{2} \rfloor})$  satisfies the interlacing property

$$\tau_1 \geq \sigma_1 \geq \tau_2 \geq \sigma_2 \geq \dots \geq \sigma_{m-1} \geq |\tau_m| \quad \text{for } d = 2m - 1,$$

and

$$\tau_1 \geq \sigma_1 \geq \tau_2 \geq \sigma_2 \geq \tau_m \geq |\sigma_m| \quad \text{for } d = 2m,$$

and  $m_{\sigma, \tau} = 0$  otherwise. If  $m_{\sigma, \tau} = 1$ , we say that  $\tau$  contains  $\sigma$ , and write  $\sigma \subset \tau$ . The fact that each  $M$ -type occurs at most once in  $\tau$  will allow us to make repeated use of Schur's lemma in the following manner: if a linear operator on  $V_\tau$  ( $V_\tau$  being a vector space on which  $\tau$  is realized) commutes

with  $\tau(M)$ , then it acts as a scalar on each  $\sigma \subset \tau$ . From now on we write  $\dim(\tau)$  and  $\dim(\sigma)$  for  $\dim(V_\tau)$  and  $\dim(V_\sigma)$ , respectively.

We connect the dimensions of  $K$ -types with the Sobolev norms introduced in Section 2.5. Firstly, let  $(\pi, \mathcal{H})$  be a unitary representation of  $G$ . For each  $\tau \in \hat{K}$ , we define

$$\chi_\tau(k) = \dim(V_\tau) \cdot \text{tr}(\tau(k))$$

(the trace being defined with respect to any invariant inner product on any realization of  $\tau$ ) and an operator  $P_\tau$  by

$$P_\tau = \int_K \overline{\chi_\tau(k)} \pi(k) dk.$$

Note that  $P_\tau$  is the orthogonal projection onto the space  $\mathcal{H}_\tau$ , where

$$\mathcal{H}_\tau := \{v \in \mathcal{H} : P_\tau v = v\}.$$

This gives rise to a decomposition of  $\mathcal{H}$  as the orthogonal direct sum

$$\mathcal{H} = \bigoplus_{\tau \in \hat{K}} \mathcal{H}_\tau.$$

If  $(\pi, \mathcal{H})$  is irreducible then each  $\mathcal{H}_\tau$  is finite-dimensional. Each  $\mathcal{H}_\tau$  has the property that  $\pi(k)|_{\mathcal{H}_\tau} \cong \tau(k)$  for all  $k \in K$ . If  $\mathcal{H}_\tau \neq 0$ , we say that  $\pi$  contains  $\tau$ , and write  $\tau \subset \pi$ . There is a similar decomposition of  $\mathcal{H}$  with respect to  $M$ -types, and projection operators  $P_\sigma = \int_M \overline{\chi_\sigma(m)} \pi(m) dm$  for  $\sigma \in \hat{M}$ . As usual,  $\lceil \beta \rceil$  denotes the smallest integer greater than or equal to  $\beta$ .

**Lemma 3.1.** *There exists  $m_0 \in \mathbb{N}$  depending only on  $K$  such that for any unitary representation  $(\pi, V)$  of  $K$ ,  $\alpha > 0$ ,  $m \in \mathbb{N}$ , and for all  $v \in V$ ,*

$$\sum_{\tau \subset V} \dim(\tau)^\alpha \|P_\tau v\|_{S^m(V)} \ll \|v\|_{S^{m+m_0 \lceil \frac{\alpha+1}{2} \rceil}(V)},$$

where the implied constant depends only on  $K$ .

*Proof.* Though the argument is standard (cf. [46, Chapter 4.4.2] or [20, Lemmas 10.3 and 10.4]), we briefly recount it: using the Harish-Chandra isomorphism and the highest weight theory, we find an element  $\omega_K$  in the center of the universal enveloping algebra of  $\mathfrak{k}$  such that for each  $\tau \in \hat{K}$ ,  $d\tau(\omega_K)$  acts as a scalar  $c_\tau$  on  $\tau$ , where  $|c_\tau| \geq \dim(\tau)^2$ . Letting  $m_0$  be the

order of  $\omega_K$ , for any  $\ell \in \mathbb{N}$ , we have (using the Cauchy-Schwarz inequality),

$$\begin{aligned} \sum_{\tau \subset V} \dim(\tau)^\alpha \|P_\tau \mathbf{v}\|_{\mathcal{S}^m(V)} &= \sum_{\tau \subset V} \frac{\dim(\tau)^\alpha}{|c_\tau|^\ell} \|P_\tau d\pi(\Omega_K^\ell) \mathbf{v}\|_{\mathcal{S}^m(V)} \\ &\leq \left( \sum_{\tau \subset V} \dim(\tau)^{2\alpha-4\ell} \sum_{\tau \subset V} \|P_\tau d\pi(\Omega_K^\ell) \mathbf{v}\|_{\mathcal{S}^m(V)}^2 \right)^{1/2} \\ &\ll \left( \sum_{\tau \subset V} \dim(\tau)^{2\alpha-4\ell} \right)^{1/2} \|\mathbf{v}\|_{\mathcal{S}^{m+\ell m_0}(V)}. \end{aligned}$$

The sum  $\sum_{\tau} \dim(\tau)^{-2}$  is finite by [20, Lemma 13], so choosing  $\ell = \lceil \frac{\alpha+1}{2} \rceil$  gives  $\sum_{\tau \subset V} \dim(\tau)^{2\alpha-4\ell} < \infty$ .  $\square$

**3.2. Complementary series representations.** We will now recall Hirai's classification [16] of the non-tempered unitary dual of  $G$  and construct the models of the complementary series that we will work with, see [46, Chapter 5.5] (cf. also [27, Sections 3.1-3.3]).

Given  $s \in \mathbb{C}$ , we define the standard representation  $U^s$  of  $G$  on  $L^2(K)$  by

$$[U^s(g)\mathbf{v}](k) = e^{-sH(g^{-1}k)} \mathbf{v}(\kappa(g^{-1}k)) \quad (3.1)$$

for all  $\mathbf{v} \in L^2(K)$ ,  $g \in G$ , and  $k \in K$ . We let  $\lambda$  and  $\rho$  denote the left- and right-regular representations of  $K$  on  $L^2(K)$  respectively. Observe that  $U^s|_K = \lambda$ ; for this reason, it is practical to always view  $L^2(K)$  as the unitary representation  $(\lambda, L^2(K))$  of  $K$ . The decomposition of  $L^2(K)$  into  $K$ -types reads

$$L^2(K) = \bigoplus_{\tau \in \hat{K}} L^2(K)_\tau,$$

where each  $L^2(K)_\tau$  is isomorphic to  $\dim(\tau)$  copies of  $\tau$ . Given  $v \in \hat{M}$ , we define

$$L^2(K : v) := \left\{ \mathbf{v} \in L^2(K) : \int_M \overline{\chi_v(m)} \rho(m) \mathbf{v} dm = \mathbf{v} \right\}.$$

The fact that  $U^s(g)$  and  $\rho(m)$  commute for all  $g \in G$  and  $m \in M$  shows that  $(U^s, L^2(K : v))$  is a representation of  $G$ ; in fact  $(U^s, L^2(K : v))$  is isomorphic to  $\dim(v)$  copies of the representation  $\text{ind}_{MAN}^G(v \otimes s \otimes 1)$  (cf. [20, Chapter VII]). Note that  $L^2(K : v)_\tau \subset L^2(K)_\tau$  for all  $\tau \in \hat{K}$ . From Hirai's classification of the unitary dual of  $G$  [16], for each  $v \in \hat{M}$ , there exists an interval  $\mathcal{I}_v \subset (\frac{d}{2}, d)$  such that if  $s \in \mathcal{I}_v$ , then there exists an irreducible, unitarizable subrepresentation  $\mathcal{U}(v, s)$  of  $(U^s, L^2(K : v))$ . Furthermore, every non-tempered representation of  $G$  may be realized as  $\mathcal{U}(v, s)$  for some  $v \in \hat{M}$  and  $s \in \mathcal{I}_v$  in this way, and each  $\tau \in \hat{K}$  that contains  $v$  occurs exactly once in  $\mathcal{U}(v, s)$  (cf. [46, Theorem 5.5.1.5] and [20, Theorem 8.37]).

The inner product that makes  $\mathcal{U}(v, s)$  a unitary representation is denoted by  $\langle \cdot, \cdot \rangle_{\mathcal{U}(v, s)}$ . Note that for  $g \in G$ , and  $\mathbf{v}, \mathbf{u} \in \mathcal{U}(v, s)$ ,

$$\langle U^s(g)\mathbf{v}, \mathbf{u} \rangle_{\mathcal{U}(v, s)} = \langle \mathcal{U}(v, s)(g)\mathbf{v}, \mathbf{u} \rangle_{\mathcal{U}(v, s)}.$$

We conclude this section by recalling two general facts regarding complementary series representations that let us classify the representation  $\mathcal{B}_\delta$  (cf. Section 2.4). Firstly, the spherical complementary series representations are the  $\mathcal{U}(1, s)$ ; i.e. precisely those where the representation  $v$  is trivial. Secondly, the Casimir operator  $\mathcal{C}$  of  $G$  acts as the scalar  $-s(d-s)$  on the smooth vectors of  $\mathcal{U}(v, s)$ . Combining these facts with the observation that the restriction of  $\mathcal{C}$  to right  $K$ -invariant functions on  $G$  is the Laplace-Beltrami operator lets one conclude that the representation  $\mathcal{B}_\delta$  is isomorphic to  $\mathcal{U}(1, \delta)$ .

**3.3. Eisenstein integrals.** Here we develop the Eisenstein integrals needed to obtain the desired asymptotic expansion of matrix coefficients (cf. [27, Section 3.3], [47, Chapter 6]). Before starting, we make the following elementary but important observation:

**Lemma 3.2.** *If  $\sigma \in \hat{M} \setminus \{v^*\}$ , then for any  $\mathbf{v} \in L^2(K : v)_\sigma$ , we have*

$$\mathbf{v}(m) = 0 \quad \text{for all } m \in M.$$

*Proof.* Since  $\mathbf{v} \in L^2(K : v)$ , we have  $\mathbf{v} = \int_M \rho(m) \mathbf{v} \overline{\chi_v(m)} dm$ . Hence for any  $m \in M$ ,

$$\begin{aligned} \mathbf{v}(m) &= \int_M \mathbf{v}(mm_1) \overline{\chi_v(m_1)} dm_1 \\ &= \int_M \mathbf{v}(m_1m) \overline{\chi_v(m^{-1}m_1m)} dm_1 \\ &= \int_M \mathbf{v}(m_1^{-1}m) \overline{\chi_v(m_1^{-1})} dm_1 \\ &= \int_M \mathbf{v}(m_1^{-1}m) \overline{\chi_{v^*}(m_1)} dm_1 \\ &= [\mathbf{P}_{v^*}\mathbf{v}](m). \end{aligned}$$

Since  $\mathbf{v} \in L^2(K : v)_\sigma$  for  $\sigma \in \hat{M} \setminus \{v^*\}$ , the orthogonality of the  $M$ -types of  $L^2(K)$  gives  $\mathbf{P}_{v^*}\mathbf{v} = 0$ , yielding the claim.  $\square$

Before stating the main result of this section, we introduce some more notation. Firstly, we let  $\langle \cdot, \cdot \rangle_K$  denote the usual inner product on  $L^2(K)$ . The corresponding norm on  $L^2(K)$  is denoted  $\| \cdot \|_K$ , and similarly for the operator norm defined with respect to it.

In the rest of this section, we fix a complementary series representation  $\mathcal{U}(v, s)$ ,  $s \in \mathcal{I}_v$ . For each  $K$ -type  $\tau$  of  $\mathcal{U}(v, s)$  we define a vector  $\chi_\tau \in \mathcal{U}(v, s)_\tau$

by

$$\chi_\tau = \sum_{i=1}^{\dim(\tau)} \overline{\mathbf{v}_i(e)} \mathbf{v}_i, \quad (3.2)$$

where  $\{\mathbf{v}_i\}_{i=1}^{\dim(\tau)}$  is an orthonormal basis of  $\mathcal{U}(v, s)_\tau$  (recall that  $\tau$  has multiplicity one in  $\mathcal{U}(v, s)$ ). Note that  $\chi_\tau$  is independent of the choice of basis, and for all  $\mathbf{v} \in \mathcal{U}(v, s)_\tau$  and  $k \in K$ , we have

$$\mathbf{v}(k) = \langle \mathbf{v}, U^s(k)\chi_\tau \rangle_K. \quad (3.3)$$

**Lemma 3.3.** *Let  $\tau$  be a  $K$ -type of  $\mathcal{U}(v, s)$ .*

(1) *For all  $\mathbf{v} \in \mathcal{U}(v, s)_\tau \cap \mathcal{U}(v, s)_{v^*}$ , we have*

$$\int_M |\mathbf{v}(m)|^2 dm = \frac{\dim(\tau) \|\mathbf{v}\|_K^2}{\dim(v)}.$$

(2)  $\int_M |\chi_\tau(m)|^2 dm = \frac{\dim(\tau)^2}{\dim(v)}$ .

*Proof.* Let  $\{\mathbf{v}_j\}$  be an orthonormal basis of  $\mathcal{U}(v, s)_\tau \cap \mathcal{U}(v, s)_{v^*}$  with respect to  $\langle \cdot, \cdot \rangle_K$ . Using (3.2) and the Schur orthogonality relations, we get

$$\begin{aligned} \int_M |\mathbf{v}(m)|^2 dm &= \int_M |\langle \mathbf{v}, U^s(m)\chi_\tau \rangle_K|^2 dm \\ &= \sum_{i,j} \overline{\mathbf{v}_i(e)} \mathbf{v}_j(e) \int_M \langle U^s(m)\mathbf{v}, \mathbf{v}_i \rangle_K \overline{\langle U^s(m)\mathbf{v}, \mathbf{v}_j \rangle_K} dm \\ &= \frac{\|\mathbf{v}\|_K^2 \sum_i |\mathbf{v}_i(e)|^2}{\dim(v)}. \end{aligned}$$

Complete  $\{\mathbf{v}_j\}$  into an orthonormal basis  $\{\mathbf{v}_j\} \cup \{\mathbf{w}_i\}$  for all of  $\mathcal{U}(v, s)_\tau$ . Note that  $\mathbf{w}_i(e) = 0$  for all  $i$  by Lemma 3.2. Therefore we have

$$\sum_j |\mathbf{v}_j(e)|^2 = \sum_j |\mathbf{v}_j(e)|^2 + \sum_i |\mathbf{w}_i(e)|^2 = \dim(\tau). \quad (3.4)$$

Hence the claim (1) follows. Claim (2) follows from (1).  $\square$

For  $K$ -types  $\tau_1, \tau_2$  of  $\mathcal{U}(v, s)$ , we define an operator

$$T_{\tau_1}^{\tau_2} : \mathcal{U}(v, s)_{\tau_1} \rightarrow \mathcal{U}(v, s)_{\tau_2}$$

by

$$T_{\tau_1}^{\tau_2} \mathbf{v} = \int_M \mathbf{v}(m) U^s(m) \chi_{\tau_2} dm \quad \text{for all } \mathbf{v} \in \mathcal{U}(v, s)_{\tau_1}. \quad (3.5)$$

We have the following interpretation of the Eisenstein integral (a similar formula appears in [27, Theorem 3.4]):

**Theorem 3.4.** *For any  $K$ -types  $\tau_1, \tau_2$  of  $\mathcal{U}(v, s)$  and  $g \in G$ , we have*

$$\mathsf{P}_{\tau_2} U^s(g) \mathsf{P}_{\tau_1} = \int_K e^{(s-d)H(gk)} U^s(\kappa(gk)) T_{\tau_1}^{\tau_2} U^s(k^{-1}) dk,$$

where  $\mathsf{P}_{\tau_2} U^s(g) \mathsf{P}_{\tau_1}$  is viewed as an operator from  $\mathcal{U}(v, s)_{\tau_1}$  to  $\mathcal{U}(v, s)_{\tau_2}$ .

*Proof.* Following [47, Theorem 6.2.2.4, pp. 42-43], the key fact is that for any  $g \in G$ ,

$$dk = e^{dH(gk)} d(\kappa(gk)).$$

Let  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_1}$  and  $\mathbf{w} \in \mathcal{U}(v, s)_{\tau_2}$ . Then

$$\begin{aligned} \langle U^s(g)\mathbf{v}, \mathbf{w} \rangle_K &= \int_K e^{-sH(g^{-1}k)} \mathbf{v}(\kappa(g^{-1}k)) \overline{\mathbf{w}(k)} dk \\ &= \int_K e^{-sH(g^{-1}k)} \mathbf{v}(\kappa(g^{-1}k)) \overline{\mathbf{w}(k)} e^{dH(g^{-1}k)} d(\kappa(g^{-1}k)). \end{aligned}$$

We now carry out the change of variables  $\tilde{k} = \kappa(g^{-1}k)$ ; this gives  $k = \kappa(g\tilde{k})$ , allowing us to rewrite the above integral as follows

$$\begin{aligned} &= \int_K e^{(d-s)H(g^{-1}\kappa(gk))} \mathbf{v}(k) \overline{\mathbf{w}(\kappa(gk))} dk \\ &= \int_K e^{(s-d)H(gk)} \mathbf{v}(k) \overline{\mathbf{w}(\kappa(gk))} dk, \end{aligned}$$

where we used the identity

$$g^{-1}\kappa(gk) = k \exp(-H(gk)) (\exp(H(gk)) n_{kg}^{-1} \exp(-H(gk)))$$

to obtain  $H(g^{-1}\kappa(gk)) = -H(gk)$ .

By (3.3), we get

$$\begin{aligned} \mathbf{v}(k) \overline{\mathbf{w}(\kappa(gk))} &= \langle \mathbf{v}, U^s(k) \chi_{\tau_1} \rangle_K \cdot \overline{\langle \mathbf{w}, U^s(\kappa(gk)) \chi_{\tau_2} \rangle_K} \\ &= \langle U^s(k^{-1}) \mathbf{v}, \chi_{\tau_1} \rangle_K \cdot \langle U^s(\kappa(gk)) \chi_{\tau_2}, \mathbf{w} \rangle_K \\ &= \left\langle \langle U^s(k^{-1}) \mathbf{v}, \chi_{\tau_1} \rangle \cdot U^s(\kappa(gk)) \chi_{\tau_2}, \mathbf{w} \right\rangle_K. \end{aligned}$$

Hence

$$\begin{aligned} \langle U^s(g)\mathbf{v}, \mathbf{w} \rangle_K &= \int_K e^{(s-d)H(gk)} \left\langle \langle U^s(k^{-1}) \mathbf{v}, \chi_{\tau_1} \rangle \cdot U^s(\kappa(gk)) \chi_{\tau_2}, \mathbf{w} \right\rangle_K dk \\ &= \left\langle \int_K e^{(s-d)H(gk)} \langle U^s(k^{-1}) \mathbf{v}, \chi_{\tau_1} \rangle \cdot U^s(\kappa(gk)) \chi_{\tau_2} dk, \mathbf{w} \right\rangle_K. \end{aligned}$$

We define  $T_0 : \mathcal{U}(v, s)_{\tau_1} \rightarrow \mathcal{U}(v, s)_{\tau_2}$  by

$$T_0 \mathbf{v} = \langle \mathbf{v}, \chi_{\tau_1} \rangle_K \cdot \chi_{\tau_2} = \mathbf{v}(e) \chi_{\tau_2}.$$

Then

$$\mathbf{P}_{\tau_2} U^s(g) \mathbf{P}_{\tau_1} = \int_K e^{(s-d)H(gk)} U^s(\kappa(gk)) T_0 U^s(k^{-1}) dk. \quad (3.6)$$

Using (3.1) and (3.5), we observe that

$$T_{\tau_1}^{\tau_2} = \int_M U^s(m) T_0 U^s(m^{-1}) dm.$$

Writing  $g = k_1 a_t k_2$ , we then have

$$\begin{aligned} \mathbf{P}_{\tau_2} U^s(g) \mathbf{P}_{\tau_1} &= U^s(k_1) \mathbf{P}_{\tau_2} U^s(a_t) \mathbf{P}_{\tau_1} U^s(k_2) \\ &= \int_M U^s(k_1) U^s(m) \mathbf{P}_{\tau_2} U^s(a_t) \mathbf{P}_{\tau_1} U^s(m^{-1}) U^s(k_2) dm, \end{aligned}$$

where we used the facts that  $M = Z_K(A)$  and  $M \subset K$  to commute  $U^s(m)$  past  $U^s(a_t)$  and  $\mathbf{P}_{\tau_i}$ . Now using (3.6) gives

$$\begin{aligned} &= U^s(k_1) \left( \int_M \int_K e^{(s-d)H(a_t k)} U^s(m \kappa(a_t k)) T_0 U^s(k^{-1} m^{-1}) dk dm \right) U^s(k_2) \\ &= \int_K e^{(s-d)H(a_t k)} U^s(k_1 \kappa(a_t k)) \left( \int_M U^s(m) T_0 U^s(m^{-1}) dm \right) U^s(k^{-1} k_2) dk. \end{aligned}$$

Finally, using the identities  $k_1 \kappa(a_t k) = \kappa(k_1 a_t k)$  and  $H(a_t k) = H(k_1 a_t k)$  together with the change of variables  $k = k_2 \tilde{k}$  gives

$$\mathbf{P}_{\tau_2} U^s(g) \mathbf{P}_{\tau_1} = \int_K e^{(s-d)H(g \tilde{k})} U^s(\kappa(g \tilde{k})) T_{\tau_1}^{\tau_2} U^s(\tilde{k}^{-1}) d\tilde{k},$$

proving the first identity claimed in the theorem. The second identity follows from the fact that  $\langle U^s(m^{-1}) \mathbf{v}, \chi_{\tau_1} \rangle_K = \mathbf{v}(m)$  for all  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_1}$  and  $m \in M$ .  $\square$

**3.4. On the operator  $T_{\tau_1}^{\tau_2}$ .** Using Lemma 3.2, we deduce:

**Proposition 3.5.** *Let  $\tau_1, \tau_2$  be  $K$ -types of  $\mathcal{U}(v, s)$ . If  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_1}$  is orthogonal to  $\mathcal{U}(v, s)_{v^*}$ , then*

$$T_{\tau_1}^{\tau_2} \mathbf{v} = 0.$$

Consequently,  $T_{\tau_1}^{\tau_2} \mathcal{U}(v, s)_{\tau_1} \subset \mathcal{U}(v, s)_{\tau_2} \cap \mathcal{U}(v, s)_{v^*}$ .

*Proof.* We have  $T_{\tau_1}^{\tau_2} \mathbf{v} = \int_M \mathbf{v}(m) U^s(m) \chi_{\tau_2} dm$ . By Lemma 3.2, if  $\mathbf{v}$  is orthogonal to  $\mathcal{U}(v, s)_{v^*}$ , then  $\mathbf{v}(m) = 0$  for all  $m \in M$  and thus  $T_{\tau_1}^{\tau_2} \mathbf{v} = 0$ . A direct computation using (3.5) shows that  $T_{\tau_1}^{\tau_2}$  commutes with  $U^s(m)$  for all  $m \in M$ , so Schur's lemma then gives that if  $T_{\tau_1}^{\tau_2} \mathbf{v}$  is non-zero, it must be contained in  $\mathcal{U}(v, s)_{\tau_2} \cap \mathcal{U}(v, s)_{v^*}$ .  $\square$

If  $\mathbf{v}$  is  $M$ -invariant and  $v \in \hat{M}$  is non-trivial, then  $v^*$  is non-trivial, and hence  $\mathbf{v}$  is orthogonal to all  $\mathcal{U}(v, s)_{v^*}$ . Therefore we deduce the following corollary:

**Corollary 3.6.** *If  $v \in \hat{M}$  is non-trivial, then for any  $M$ -invariant  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_1}$ ,*

$$T_{\tau_1}^{\tau_2} \mathbf{v} = 0.$$

**Lemma 3.7.** *For any  $K$ -types  $\tau_1, \tau_2$  of  $\mathcal{U}(v, s)$ , we have*

$$(T_{\tau_1}^{\tau_2})^* = T_{\tau_2}^{\tau_1},$$

where the adjoint is defined with respect to  $\langle \cdot, \cdot \rangle_K$ .



*Proof.* For any  $\mathbf{u} \in \mathcal{U}(v, s)_{\tau_1}$  and  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_2}$ , we have

$$\begin{aligned} \langle T_{\tau_1}^{\tau_2} \mathbf{u}, \mathbf{v} \rangle_K &= \int_M \mathbf{u}(m) \langle U^s(m) \chi_{\tau_2}, \mathbf{v} \rangle_K dm \\ &= \int_M \langle \mathbf{u}, U^s(m) \chi_{\tau_1} \rangle_K \langle U^s(m) \chi_{\tau_2}, \mathbf{v} \rangle_K dm \\ &= \left\langle \mathbf{u}, \int_M \langle U^s(m^{-1}) \mathbf{v}, \chi_{\tau_2} \rangle_K \cdot U^s(m) \chi_{\tau_1} dm \right\rangle \\ &= \langle \mathbf{u}, T_{\tau_2}^{\tau_1} \mathbf{v} \rangle_K. \end{aligned}$$

□

**Corollary 3.8.** *For any  $K$ -types  $\tau_1, \tau_2$  of  $\mathcal{U}(v, s)$ ,*

$$\|T_{\tau_1}^{\tau_2}\|_K \leq \frac{\sqrt{\dim(\tau_1) \dim(\tau_2)}}{\dim(v)}.$$

*Proof.* For any unit vector  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_1} \cap \mathcal{U}(v, s)_{v^*}$ , we have

$$\begin{aligned} \langle T_{\tau_1}^{\tau_2} \mathbf{v}, T_{\tau_1}^{\tau_2} \mathbf{v} \rangle_K &= \int_M \overline{\mathbf{v}(m)} \langle T_{\tau_1}^{\tau_2} \mathbf{v}, U^s(m) \chi_{\tau_2} \rangle_K dm \\ &= \int_M \int_M \overline{\mathbf{v}(m_1)} \mathbf{v}(m_2) \langle U^s(m_2) \chi_{\tau_2}, U^s(m_1) \chi_{\tau_2} \rangle_K dm_2 dm_1 \\ &= \int_M \int_M \overline{\langle \mathbf{v}, U^s(m_1) \chi_{\tau_1} \rangle_K} \langle \mathbf{v}, U^s(m_2) \chi_{\tau_1} \rangle_K \chi_{\tau_2}(m_2^{-1} m_1) dm_2 dm_1 \\ &= \int_M \left( \int_M \overline{\langle \mathbf{v}, U^s(m_2 m_3) \chi_{\tau_1} \rangle_K} \langle \mathbf{v}, U^s(m_2) \chi_{\tau_1} \rangle_K dm_2 \right) \chi_{\tau_2}(m_3) dm_3. \end{aligned}$$

Since  $\mathbf{v}$  and  $\chi_{\tau_1}$  are both in  $\mathcal{U}(v, s)_{\tau_1} \cap \mathcal{U}(v, s)_{v^*}$ , the Schur orthogonality relations for  $M$  give

$$\begin{aligned} &\int_M \overline{\langle \mathbf{v}, U^s(m_2 m_3) \chi_{\tau_1} \rangle_K} \langle \mathbf{v}, U^s(m_2) \chi_{\tau_1} \rangle_K dm_2 \\ &= \frac{\|\mathbf{v}\|^2}{\dim(v)} \overline{\langle \chi_{\tau_1}, U^s(m_3) \chi_{\tau_1} \rangle} = \frac{\overline{\chi_{\tau_1}(m_3)}}{\dim(v)}. \end{aligned}$$

We thus have

$$\begin{aligned} \langle T_{\tau_1}^{\tau_2} \mathbf{v}, T_{\tau_1}^{\tau_2} \mathbf{v} \rangle_K &= \frac{1}{\dim(v)} \int_M \overline{\chi_{\tau_1}(m)} \chi_{\tau_2}(m) dm \\ &\leq \frac{1}{\dim(v)} \sqrt{\int_M |\chi_{\tau_1}(m)|^2 dm \int_M |\chi_{\tau_2}(m)|^2 dm} \\ &= \frac{\dim(\tau_1) \dim(\tau_2)}{\dim(v)^2} \end{aligned}$$

by Lemma 3.3 (2). Combining this estimate with Proposition 3.5, the claim follows. □

**Corollary 3.9.** *For any  $K$ -type  $\tau$  of  $\mathcal{U}(v, s)$ , we have*

$$T_\tau^\tau = \frac{\dim(\tau)}{\dim(v)} \mathbf{P}_{v^*}.$$

*Proof.* Following the proof of Corollary 3.8, we obtain

$$\|T_\tau^\tau\|_K^2 = \frac{1}{\dim(v)} \int_M |\chi_\tau(m)|^2 dm = \frac{\dim(\tau)^2}{\dim(v)^2}.$$

So by Proposition 3.5, we have

$$T_\tau^\tau = c \cdot \mathbf{P}_{v^*},$$

where  $c$  is one of  $\pm \frac{\dim(\tau)}{\dim(v)}$ . Since

$$c = \langle T_\tau^\tau \mathbf{v}, \mathbf{v} \rangle_K = \int_M |\mathbf{v}(m)|^2 dm \geq 0$$

for any unit vector  $\mathbf{v} \in \mathcal{U}(v, s)_\tau \cap \mathcal{U}(v, s)_{v^*}$ , we get  $c = \frac{\dim(\tau)}{\dim(v)}$ .  $\square$

#### 4. ASYMPTOTIC EXPANSIONS OF MATRIX COEFFICIENTS

We fix a complementary series representation  $\mathcal{U}(v, s)$  of  $G$  for some  $v \in \hat{M}$  and  $s \in \mathcal{I}_v$ . The main goal of this section is to obtain effective expansions of matrix coefficients for  $\mathcal{U}(v, s)$ . More precisely, we will write a matrix coefficient  $\langle U^s(a_t) \mathbf{v}, \mathbf{u} \rangle_{\mathcal{U}(v, s)}$  as a main term that decays like  $e^{(s-d)t}$  as  $t \rightarrow \infty$  and an error term that decays exponentially faster depending only on  $2s - d > 0$ . To do this, we first work out the asymptotics of matrix coefficients with respect to  $\langle \cdot, \cdot \rangle_K$  and then use explicit formulas for intertwining operators to convert our results into statements for  $\langle \cdot, \cdot \rangle_{\mathcal{U}(v, s)}$ .

**4.1. Matrix coefficients with respect to  $\langle \cdot, \cdot \rangle_K$ .** We start by proving a bound on how far elements of  $K$  move vectors in representations of  $K$ . We let  $\text{dist}_K$  denote the bi-invariant Riemannian metric on  $K$  induced from the negative of the Killing form on  $\mathfrak{k}$ , and let  $|\cdot|_K$  denote the corresponding norm on  $\mathfrak{k}$ .

**Lemma 4.1.** *Let  $(\pi, V)$  be a finite-dimensional unitary representation of  $K$  with invariant inner product  $(\cdot, \cdot)_V$ . Then for all  $\mathbf{v} \in V$  and  $k \in K$ , we have*

$$\|\pi(k)\mathbf{v} - \mathbf{v}\|_V \ll \text{dist}_K(k, e) \|\mathbf{v}\|_{\mathcal{S}^1(V)},$$

where the implied constant depends solely on the choice of basis defining  $\|\cdot\|_{\mathcal{S}^m(V)}$ .

*Proof.* Since  $\text{dist}_K$  is bi-invariant, there exists  $J \in \mathfrak{k}$  with  $|J|_K = 1$  such that the  $k = \exp(\text{dist}_K(k, e)J)$ . This gives

$$\begin{aligned} \pi(k)\mathbf{v} - \mathbf{v} &= \int_0^{\text{dist}_K(k, e)} \frac{d}{dt} \pi(\exp(tJ)) \mathbf{v} dt \\ &= \int_0^{\text{dist}_K(k, e)} \pi(\exp(tJ)) d\pi(\text{Ad}(\exp(-tJ))J) \mathbf{v} dt. \end{aligned}$$

Since  $\pi$  is unitary,  $|J|_K = 1$ , and  $K$  is compact, we get

$$\|\pi(\exp(tJ))d\pi(\text{Ad}(\exp(-tJ))J)\mathbf{v}\|_V \leq \max_{\substack{X \in \mathfrak{t} \\ |X|_K=1}} \|d\pi(X)\mathbf{v}\|_V \ll \|\mathbf{v}\|_{\mathcal{S}^1(V)}.$$

□

**Definition 4.1.** The Harish-Chandra  $c$ -function  $C_+(s) : \mathcal{U}(v, s) \rightarrow \mathcal{U}(v, s)$  is defined as follows:

$$C_+(s) = \int_{\bar{N}} U^s(\kappa(\bar{n})^{-1})e^{-sH(\bar{n})} d\bar{n}. \quad (4.1)$$

Since  $U^s$  is unitary on  $\mathcal{U}(v, s)$  and  $\int_{\bar{N}} e^{-sH(\bar{n})} d\bar{n} < \infty$  for all  $s > \frac{d}{2}$  (cf. [20, Proposition 7.6]),  $C_+(s)$  is a well-defined bounded operator on  $\mathcal{U}(v, s)$ . Since  $C_+(s)$  is defined using only the restriction of  $U^s$  to  $K$ , it preserves the  $K$ -types of  $\mathcal{U}(v, s)$ .

**Lemma 4.2.**  $C_+(s)$  preserves the  $M$ -types of  $\mathcal{U}(v, s)$ .

*Proof.* Note that  $M$  normalizes  $\bar{N}$ ,  $d\bar{n} = d(m\bar{n}m^{-1})$  and  $H(mgm^{-1}) = H(g)$  for all  $m \in M$ ,  $\bar{n} \in \bar{N}$ , and  $g \in G$ . Therefore for all  $m \in M$ ,

$$\begin{aligned} C_+(s)U^s(m) &= \int_{\bar{N}} U^s(\kappa(\bar{n})^{-1}m)e^{-sH(\bar{n})} d\bar{n} \\ &= \int_{\bar{N}} U^s(m\kappa(m^{-1}\bar{n}m)^{-1})e^{-sH(\bar{n})} d\bar{n} = U^s(m)C_+(s); \end{aligned}$$

hence the claim follows. □

For  $d/2 < s < d$ , set

$$\eta_s = \min\{2s - d, 1\} > 0. \quad (4.2)$$

We remark that the following theorem was shown in [27, Theorem 3.23] for some  $\eta_s$ , based on Harish-Chandra's expansion formula and the maximum modulus principle. We give a more direct argument, with the explicit  $\eta_s$  given in (4.2). Let  $\mathcal{S}^1(K)$  denote the Sobolev norm  $\mathcal{S}^1(L^2(K))$  defined in Section 2.5.

**Theorem 4.3.** Let  $\tau_1, \tau_2$  be  $K$ -types of  $\mathcal{U}(v, s)$ . For all  $\mathbf{u} \in \mathcal{U}(v, s)_{\tau_1}$  and  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_2}$ , we have for any  $t \geq 0$ ,

$$\begin{aligned} \langle U^s(a_t)\mathbf{u}, \mathbf{v} \rangle_K &= e^{(s-d)t} \langle T_{\tau_1}^{\tau_2} C_+(s)\mathbf{u}, \mathbf{v} \rangle_K \\ &\quad + O_s \left( e^{(s-d-\eta_s)t} \|T_{\tau_1}^{\tau_2}\|_K \|\mathbf{u}\|_K \|\mathbf{v}\|_{\mathcal{S}^1(K)} \right), \end{aligned}$$

where the implied constant is uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ .

*Proof.* Applying Theorem 3.4, we have

$$\langle U^s(a_t)\mathbf{u}, \mathbf{v} \rangle_K = \int_K e^{(s-d)H(a_t k)} \langle U^s(\kappa(a_t k)) T_{\tau_1}^{\tau_2} U^s(k^{-1})\mathbf{u}, \mathbf{v} \rangle_K dk.$$

Since the function  $k \mapsto e^{(s-d)H(a_t k)} \langle U^s(\kappa(a_t k)) T_{\tau_1}^{\tau_2} U^s(k^{-1}) \mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v,s)}$  is right  $M$ -invariant, we may use the integration formula [20, Consequence 3, p. 147] to obtain

$$\begin{aligned} & \langle U^s(a_t) \mathbf{u}, \mathbf{v} \rangle_K \\ &= \int_{\bar{N}} e^{(s-d)H(a_t \kappa(\bar{n}))} \langle U^s(\kappa(a_t \kappa(\bar{n}))) T_{\tau_1}^{\tau_2} U^s(\kappa(\bar{n})^{-1}) \mathbf{u}, \mathbf{v} \rangle_K e^{-dH(\bar{n})} d\bar{n}. \end{aligned}$$

The identities

$$\begin{aligned} H(a_t \kappa(\bar{n})) &= H(a_t \bar{n} a_{-t}) + H(a_t) - H(\bar{n}) \quad \text{and} \\ \kappa(a_t \kappa(\bar{n})) &= \kappa(a_t \bar{n} a_{-t}) \end{aligned}$$

then give that the previous integral is equal to

$$e^{(s-d)t} \int_{\bar{N}} \langle T_{\tau_1}^{\tau_2} U^s(\kappa(\bar{n})^{-1}) \mathbf{u}, U^s(\kappa(a_t \bar{n} a_{-t})^{-1}) \mathbf{v} \rangle_K e^{(s-d)H(a_t \bar{n} a_{-t}) - sH(\bar{n})} d\bar{n}.$$

We now use the identification of  $\bar{N}$  with  $\mathbb{R}^d$  to again rewrite:

$$e^{(s-d)t} \int_{\mathbb{R}^d} \langle T_{\tau_1}^{\tau_2} U^s(\kappa(\bar{n}_{\mathbf{x}})^{-1}) \mathbf{u}, U^s(\kappa(\bar{n}_{e^{-t}\mathbf{x}})^{-1}) \mathbf{v} \rangle_K e^{(s-d)H(\bar{n}_{e^{-t}\mathbf{x}}) - sH(\bar{n}_{\mathbf{x}})} d\mathbf{x}; \quad (4.3)$$

note that the integral is absolutely convergent due to  $s > \frac{d}{2}$ . Using the fact that  $e^{H(\bar{n}_{\mathbf{x}})} = 1 + \|\mathbf{x}\|^2$  (cf. [21, p. 486 and p. 564]) gives

$$\begin{aligned} \langle U^s(a_t) \mathbf{u}, \mathbf{v} \rangle_K &= e^{(s-d)t} \langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{u}, \mathbf{v} \rangle_K + \\ & e^{(s-d)t} \int_{\mathbb{R}^d} \langle T_{\tau_1}^{\tau_2} U^s(\kappa(\bar{n}_{\mathbf{x}})^{-1}) \mathbf{u}, (1 + \|e^{-t}\mathbf{x}\|^2)^{s-d} U^s(\kappa(\bar{n}_{e^{-t}\mathbf{x}})^{-1}) \mathbf{v} - \mathbf{v} \rangle_K \frac{d\mathbf{x}}{(1 + \|\mathbf{x}\|^2)^s}. \end{aligned}$$

Changing to spherical coordinates, we let  $\mathbf{x} = (r, \boldsymbol{\theta})$ ,  $r \geq 0$ ,  $\boldsymbol{\theta} \in \mathbb{S}^{d-1}$  and set

$$\mathbf{w}(r, \boldsymbol{\theta}) = (1 + r^2)^{s-d} U^s(\kappa(\bar{n}_{(r, \boldsymbol{\theta})})^{-1}) \mathbf{v} - \mathbf{v}.$$

Using this, we deduce

$$\begin{aligned} \langle U^s(a_t) \mathbf{u}, \mathbf{v} \rangle_K &= e^{(s-d)t} \langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{u}, \mathbf{v} \rangle_K + \\ & O\left(e^{(s-d)t} \int_0^\infty \int_{\mathbb{S}^{d-1}} \langle T_{\tau_1}^{\tau_2} U^s(\kappa(\bar{n}_{(r, \boldsymbol{\theta})})^{-1}) \mathbf{u}, \mathbf{w}(e^{-t}r, \boldsymbol{\theta}) \rangle_K dm(\boldsymbol{\theta}) \frac{r^{d-1}}{(1+r^2)^s} dr\right), \end{aligned}$$

where  $dm(\boldsymbol{\theta})$  denotes the spherical measure.

Since the map  $\bar{n} \mapsto \kappa(\bar{n})$  is smooth, we get

$$d_K(\kappa(\bar{n}_{(r, \boldsymbol{\theta})}), e) \ll r \quad \text{for all } r > 0, \boldsymbol{\theta} \in \mathbb{S}^{d-1}.$$

Using  $(1 + r^2)^{s-d} = 1 - O_s(r)$  (with the implied constant depending continuously on  $s$ ) and Lemma 4.1, we have

$$\|\mathbf{w}(r, \boldsymbol{\theta})\|_K \ll_s \min\{1, r\} \|\mathbf{v}\|_{S^1(K)}.$$

This gives

$$\begin{aligned} \langle U^s(a_t)\mathbf{u}, \mathbf{v} \rangle_K &= e^{(s-d)t} \langle T_{\tau_1}^{\tau_2} C_+(s)\mathbf{u}, \mathbf{v} \rangle_K \\ &+ O_s \left( \|T_{\tau_1}^{\tau_2}\|_K \|\mathbf{u}\|_K \|\mathbf{v}\|_{S^1(K)} e^{(s-d)t} \int_0^\infty \min\{1, e^{-t}r\} (1+r^2)^{-s} r^{d-1} dr \right). \end{aligned}$$

The proof is completed by writing the integral  $\int_0^\infty$  as  $\int_0^1 + \int_1^{e^t} + \int_{e^t}^\infty$  to obtain

$$\begin{aligned} &\int_0^\infty \min\{1, e^{-t}r\} (1+r^2)^{-s} r^{d-1} dr \\ &\leq e^{-t} + e^{-t} \int_1^{e^t} r \cdot r^{-2s} \cdot r^{d-1} dr + \int_{e^t}^\infty r^{-2s} \cdot r^{d-1} dr \\ &\ll_s e^{-t} + e^{(d-2s)t}. \end{aligned}$$

□

**4.2. The invariant inner product on  $\mathcal{U}(v, s)$ .** The intertwining operator  $\mathcal{A}(v, s)$  on  $\mathcal{U}(v, s)$  is defined so that

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)} = \langle \mathbf{u}, \mathcal{A}(v, s)\mathbf{v} \rangle_K$$

for all  $K$ -finite vectors  $\mathbf{u}, \mathbf{v} \in \mathcal{U}(v, s)$ . The key intertwining property of  $\mathcal{A}(v, s)$  reads (cf. [22, Lemmas 22 and 23])

$$\mathcal{A}(v, s)U^s(g) = U^{d-s}(g)\mathcal{A}(v, s) \quad \text{for all } g \in G. \quad (4.4)$$

In particular,  $\mathcal{A}(v, s)$  commutes with  $U^s(k)$  for all  $k \in K$ . Since each  $K$ -type occurs at most once in  $\mathcal{U}(v, s)$ , by Schur's lemma,  $\mathcal{A}(v, s)$  acts as a scalar  $a(v, s, \tau)$  on each  $K$ -type  $\tau$  of  $\mathcal{U}(v, s)$ :

$$\mathcal{A}(v, s) = \sum_{\tau \supset v} a(v, s, \tau) P_\tau. \quad (4.5)$$

The positive definiteness of the inner product  $\langle \cdot, \cdot \rangle_{\mathcal{U}(v, s)}$  implies that for all  $\tau \in \hat{K}$  contained in  $\mathcal{U}(v, s)$ , we have

$$a(v, s, \tau) > 0.$$

Recalling the parameterization of  $K$  and  $M$  types given in Section 3.1, we now assume that  $\mathcal{U}(v, s)$  has a non-trivial  $M$ -invariant vector. There is then a  $K$ -type  $\sigma = (\sigma_1, \sigma_2, \dots)$  of  $\mathcal{U}(v, s)$  that contains the trivial representation of  $M$ . Thus  $(\sigma_1, \sigma_2, \dots)$  satisfies the interlacing relation with the trivial representation  $(0, 0, 0, \dots, 0)$  of  $M$ :

$$\sigma_1 \geq 0 \geq \sigma_2 \geq 0 \geq \sigma_3 \geq \dots$$

From this, we conclude that  $\sigma = (\sigma_1, 0, 0, \dots, 0)$ . Now writing  $v = (v_1, v_2, \dots)$ , the classification of  $K$ -types of  $\mathcal{U}(v, s)$  ensures that  $v$  is a subrepresentation of  $\sigma$ , and so  $(v_1, v_2, \dots)$  must satisfy the interlacing relation with  $(\sigma_1, 0, 0, 0, \dots)$ . We therefore see that  $v = (v_1, 0, 0, 0, \dots)$ . For notational convenience we will simply write  $v = (v, 0, 0, 0, \dots) \in \mathbb{Z}^{\lfloor \frac{d}{2} \rfloor}$ , where

$v \geq 0$ . Note that if  $d = 1$  or  $2$ , then by the classification of the unitary dual of  $\mathrm{SO}(2, 1)$  and  $\mathrm{SO}(3, 1)$ ,  $v = 0$  [16]. Combining the fact that each  $K$ -type of  $\mathcal{U}(v, s)$  must have  $v$  as a subrepresentation with the interlacing relation gives that all  $K$ -types  $\tau$  of  $\mathcal{U}(v, s)$  may be parameterized as vectors  $\tau = (t_1, t_2, 0, \dots, 0) \in \mathbb{Z}^{\lfloor \frac{d+1}{2} \rfloor}$  with  $t_1 \geq v \geq t_2 \geq 0$  if  $d > 3$ ,  $t_1 \geq v \geq |t_2|$  if  $d = 3$ ,  $t_1 \geq 0$  if  $d = 2$ , and just  $t_1 \in \mathbb{Z}$  if  $d = 1$ . We will thus write  $\tau = (t_1, t_2)$  and  $a(v, s, \tau) = a(v, s, t_1, t_2)$ , with  $v = t_2 = 0$  for  $d = 1, 2$ . Again using the classification of the unitary dual of  $G$ , observe that if  $s \geq d - 1$ , then  $v = t_2 = 0$ . Finally, we let  $\Omega_K \in Z(\mathfrak{k}_{\mathbb{C}})$  denote the Casimir operator of  $K$ .

The following lemma is the main technical result needed to establish the bounds on the quotients  $\frac{a(v, s, \tau_1)}{a(v, s, \tau_2)}$  given in Propositions 4.5 and 4.6.

**Lemma 4.4.** *Assume that  $\mathcal{U}(v, s)$  has a non-trivial  $M$ -invariant vector. Then for any  $K$ -types  $\tau_1 = (t_1, t_2)$  and  $\tau_2 = (t_3, t_4)$  of  $\mathcal{U}(v, s)$ , if  $s < d - 1$ ,*

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} = \frac{\Gamma(d - s + t_3)}{\Gamma(s + t_3)} \cdot \frac{\Gamma(d - s + t_4 - 1)}{\Gamma(s + t_4 - 1)} \cdot \frac{\Gamma(s + t_1)}{\Gamma(d - s + t_1)} \cdot \frac{\Gamma(s + t_2 - 1)}{\Gamma(d - s + t_2 - 1)},$$

and if  $s \geq d - 1$ ,

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} = \frac{\Gamma(d - s + t_3)}{\Gamma(s + t_3)} \cdot \frac{\Gamma(s + t_1)}{\Gamma(d - s + t_1)}.$$

*Proof.* The claimed formula follows from a recursion formula similar to that for the Harish-Chandra  $c$ -function given in [10]. We start by letting  $H \in \mathfrak{a}$  be such that  $a_t = \exp(tH)$ , and defining an inner product  $\langle \cdot, \cdot \rangle_{\mathfrak{g}}$  on  $\mathfrak{g}_{\mathbb{C}}$  by

$$\langle X, Y \rangle_{\mathfrak{g}} := c \cdot (-B(X, \theta Y)),$$

where  $B(\cdot, \cdot)$  denotes the Killing form on  $\mathfrak{g}$ ,  $\theta$  is the Cartan involution defining  $K$ , and  $c \in \mathbb{R}_{>0}$  is chosen so that  $\langle H, H \rangle_{\mathfrak{g}} = 1$ . Denote the  $-1$ -eigenspace of  $\theta$  by  $\mathfrak{p}$ . Then for all  $X \in \mathfrak{p}_{\mathbb{C}}$ ,  $\mathbf{v} \in C^\infty(K)$ , and  $k \in K$ , we have

$$\begin{aligned} [dU^s(X)\mathbf{v}](k) &= (s - \frac{d}{2})\langle \mathrm{Ad}_{k^{-1}}X, H \rangle_{\mathfrak{g}}\mathbf{v}(k) \\ &\quad + \frac{1}{2} \left( d\rho(\Omega_K) \{ \langle \mathrm{Ad}_{k^{-1}}X, H \rangle_{\mathfrak{g}}\mathbf{v}(k) \} - \langle \mathrm{Ad}_{k^{-1}}X, H \rangle_{\mathfrak{g}} [d\rho(\Omega_K)\mathbf{v}](k) \right), \end{aligned} \quad (4.6)$$

where  $\rho$  denotes right-translation (cf. [42, Lemma 1] and [10, Lemma 3.2]). Let  $\tau = (r_1, r_2)$  be an arbitrary  $K$ -type of  $\mathcal{U}(v, s)$ . By e.g. [42, Lemma 2] (cf. also [21, Proposition 5.28]), for any vector  $\mathbf{v} \in \mathcal{U}(v, s)_\tau$ ,

$$d\rho(\Omega_K)\mathbf{v} = dU^s(\Omega_K)\mathbf{v} = (r_1^2 + r_2^2 + (d - 1)r_1 + (d - 3)r_2)\mathbf{v}. \quad (4.7)$$

We now denote the orthogonal projection onto  $\tau$  by  $\mathbf{P}_{(r_1, r_2)}$ . Combining the above expression for  $d\rho(\Omega_K)\mathbf{v}$  with (4.6) gives

$$\begin{aligned} & [\mathbf{P}_{(r_1+1, r_2)} dU^s(X)\mathbf{v}](k) \\ &= \left( s - \frac{d}{2} + \frac{(r_1+1)^2 + r_2^2 + (d-1)(r_1+1) + (d-3)r_2}{2} - \frac{r_1^2 + r_2^2 + (d-1)r_1 + (d-3)r_2}{2} \right) \\ & \quad \times \mathbf{P}_{(r_1+1, r_2)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}(k) \} \\ &= (s + r_1) \mathbf{P}_{(r_1+1, r_2)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}(k) \}, \end{aligned} \quad (4.8)$$

and similarly

$$[\mathbf{P}_{(r_1, r_2+1)} dU^s(X)\mathbf{v}](k) = (s + r_2 - 1) \mathbf{P}_{(r_1, r_2+1)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}(k) \}. \quad (4.9)$$

From (4.4), we obtain

$$[\mathbf{P}_{(r_1, r_2)} dU^s(X)\mathcal{A}(v, s)\mathbf{v}](k) = [\mathbf{P}_{(r_1, r_2)} \mathcal{A}(v, s) dU^{d-s}(X)\mathbf{v}](k).$$

Combining this with (4.5) and (4.8) and (4.9), respectively, gives

$$\begin{aligned} & a(v, s, r_1, r_2)(s + r_1) \mathbf{P}_{(r_1+1, r_2)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}(k) \} \\ &= a(v, s, r_1 + 1, r_2)(d - s + r_1) \mathbf{P}_{(r_1+1, r_2)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}(k) \}, \end{aligned}$$

and

$$\begin{aligned} & a(v, s, r_1, r_2)(s + r_2 - 1) \mathbf{P}_{(r_1, r_2+1)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}(k) \} \\ &= a(v, s, r_1, r_2 + 1)(d - s + r_2 - 1) \mathbf{P}_{(r_1, r_2+1)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}(k) \}. \end{aligned}$$

The decomposition of  $(\text{Ad}, \mathfrak{p}_{\mathbb{C}}) \otimes (U^s, \mathcal{U}(v, s)_{\tau})$  into irreducible representations of  $K$  ensures the existence of  $X, Y \in \mathfrak{p}_{\mathbb{C}}$  and  $\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{U}(v, s)_{\tau}$  such that  $\mathbf{P}_{(r_1+1, r_2)} \{ \langle \text{Ad}_{k^{-1}} X, H \rangle_{\mathfrak{g}} \mathbf{v}_1(k) \} \neq 0$ , and  $\mathbf{P}_{(r_1, r_2+1)} \{ \langle \text{Ad}_{k^{-1}} Y, H \rangle_{\mathfrak{g}} \mathbf{v}_2(k) \} \neq 0$  if  $(r_1 + 1, r_2)$ , and  $(r_1, r_2 + 1)$  are  $K$ -types of  $\mathcal{U}(v, s)$ , cf., e.g., [42, Lemma 3] and [44, Theorem 3.4.12]. This gives

$$(s + r_1) \cdot a(v, s, r_1 + 1, r_2) = (d - s + r_1) \cdot a(v, s, r_1, r_2)$$

and

$$(s + r_2 - 1) \cdot a(v, s, r_1, r_2 + 1) = (d - s + r_2 - 1) \cdot a(v, s, r_1, r_2)$$

(cf. [10, (6.1) and (6.7)]). Note that if one of the four factors that are multiplied with  $a(v, s, \dots)$ 's is zero, then the representation given by that choice of  $v$  and  $s$  is not in the unitary dual of  $G$ . In particular, if  $s \geq d - 1$ , then  $\mathcal{U}(v, s)$  is spherical, hence  $v = r_2 = 0$ . These two recursion formulas imply that if  $v \neq 0$ ,

$$a(v, s, r_1, r_2) = \frac{\Gamma(d-s+r_1)}{\Gamma(s+r_1)} \cdot \frac{\Gamma(d-s+r_2-1)}{\Gamma(s+r_2-1)} \cdot \frac{\Gamma(s+v)}{\Gamma(d-s+v)} \cdot \frac{\Gamma(s+v-1)}{\Gamma(d-s+v-1)} \cdot a(v, s, v, v),$$

and if  $v = 0$ ,

$$a(0, s, r_1, 0) = \frac{\Gamma(d-s+r_1)}{\Gamma(s+r_1)} \cdot \frac{\Gamma(s)}{\Gamma(d-s)} \cdot a(0, s, 0, 0).$$

The formulas claimed in the proposition then follow.  $\square$

By Schur's lemma,  $\Omega_K$  acts on any realization of a  $K$ -type  $\tau$  by the same scalar, which we denote  $\tau(\Omega_K)$ .

**Proposition 4.5.** *Assume that  $\mathcal{U}(v, s)$  has a non-trivial  $M$ -invariant vector. Then for any  $K$ -types  $\tau_1, \tau_2$  of  $\mathcal{U}(v, s)$ ,*

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} \ll_s (1 + \tau_1(\Omega_K)^d)(1 + \tau_2(\Omega_K)^d),$$

and the implied constant is uniformly bounded over  $s$  in compact subsets of  $\mathcal{I}_v$ .

*Proof.* The key result needed in the proof is a consequence of [17, Theorem 1]: for all  $t \geq 0$ ,

$$\frac{\Gamma(s+t)}{\Gamma(d-s+t)} \asymp_s 1 + t^{2s-d}, \quad (4.10)$$

where the implied constant is uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ . We now write  $\tau_1 = (t_1, t_2)$ ,  $\tau_2 = (t_3, t_4)$ .

Assuming first that  $d > 3$  and  $v > 0$ , we then have  $t_i \geq 0$  for all  $i$ . Since  $v > 0$ ,  $s < d - 1$ , and so rewriting Lemma 4.4 gives

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} = \frac{\Gamma(d-s+t_3)}{\Gamma(s+t_3)} \cdot \frac{\Gamma(d-s+t_4)}{\Gamma(s+t_4)} \cdot \frac{\Gamma(s+t_1)}{\Gamma(d-s+t_1)} \cdot \frac{\Gamma(s+t_2)}{\Gamma(d-s+t_2)} \cdot \frac{s+t_4-1}{d-s+t_4-1} \cdot \frac{d-s+t_2-1}{s+t_2-1},$$

and so (after recalling the formula (4.7) for  $\tau_1(\Omega_K)$  and  $\tau_2(\Omega_K)$ ),

$$\begin{aligned} \frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} &\asymp_s (1 + t_1^{2s-d})(1 + t_2^{2s-d})(1 + t_3^{2s-d})(1 + t_4^{2s-d}) \\ &\ll (1 + \tau_1(\Omega_K)^d)(1 + \tau_2(\Omega_K)^d), \end{aligned}$$

with the implied constant uniformly bounded over  $(\frac{d}{2}, d - 1) = \mathcal{I}_v$ .

In the case  $d \geq 2$  and  $v = 0$  (and so  $t_2 = t_4 = 0$ ), Lemma 4.4 gives

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} = \frac{\Gamma(d-s+t_3)}{\Gamma(s+t_3)} \cdot \frac{\Gamma(s+t_1)}{\Gamma(d-s+t_1)},$$

so by a direct application of (4.7) and (4.10),

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} \ll_s (1 + t_1^{2s-d})(1 + t_3^{2s-d}) \ll (1 + \tau_1(\Omega_K)^{\frac{d}{2}})(1 + \tau_2(\Omega_K)^{\frac{d}{2}}), \quad (4.11)$$

with the implied constant uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ .

For the remaining cases,  $d = 1$  and  $d = 3$  with  $v > 0$ , negative values of the  $t_i$  can appear in the formulas given in Lemma 4.4. It thus remains to bound quotients of the form  $\frac{\Gamma(s+t)}{\Gamma(d-s+t)}$ , where  $t$  is a negative integer and  $s \in (\frac{1}{2}, 1)$  if  $d = 1$ , and  $s \in (\frac{3}{2}, 2)$  if  $d = 3$ . By the reflection formula, in both cases,

$$\frac{\Gamma(s+t)}{\Gamma(d-s+t)} = \frac{\Gamma(s-|t|)}{\Gamma(d-s-|t|)} = \frac{\Gamma(|t|+1+s-d)}{\Gamma(|t|+1-s)}. \quad (4.12)$$

If  $d = 1$ , we then have

$$\frac{\Gamma(s+t)}{\Gamma(1-s+t)} = \frac{\Gamma(|t|+s)}{\Gamma(|t|+1-s)},$$



and so (4.10) gives

$$\begin{aligned} \frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} &= \frac{\Gamma(1-s+|t_3|)}{\Gamma(s+|t_3|)} \cdot \frac{\Gamma(s+|t_1|)}{\Gamma(1-s+|t_1|)} \\ &\ll_s (1 + |t_1|^{2s-1})(1 + |t_2|^{2s-1}) \ll (1 + \tau_1(\Omega_K)^{\frac{1}{2}})(1 + \tau_2(\Omega_K)^{\frac{1}{2}}). \end{aligned} \quad (4.13)$$

If  $d = 3$ , using (4.12), we have

$$\frac{\Gamma(s+t)}{\Gamma(3-s+t)} = \frac{\Gamma(|t+s-2|)}{\Gamma(|t+2-s|)} = \frac{(|t+1-s|)(|t+2-s|)}{(s+|t|-1)(s+|t|-2)} \cdot \frac{\Gamma(s+|t|)}{\Gamma(3-s+|t|)}.$$

Now using  $\frac{(|t+1-s|)(|t+2-s|)}{(s+|t|-1)(s+|t|-2)} \asymp_s 1$ , where the implied constants are uniformly bounded over  $s$  in compact subsets of  $(\frac{3}{2}, 2)$ , together with (4.10) as previously completes the proof.  $\square$

For  $K$ -types that contain  $M$ -invariant vectors, we have the following strengthening of the bound in Proposition 4.5:

**Proposition 4.6.** *Let  $\tau_1, \tau_2$  be  $K$ -types of  $\mathcal{U}(v, s)$  containing non-trivial  $M$ -invariant vectors. Then*

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} \ll_s (1 + \tau_1(\Omega_K)^{\frac{d}{2}})(1 + \tau_2(\Omega_K)^{\frac{d}{2}}),$$

and the implied constant is uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ .

*Proof.* For  $d \geq 3$ , we observe that if a  $K$ -type  $\tau = (r_1, r_2)$  contains an  $M$ -invariant vector, the interlacing relation implies  $r_2 = 0$ , hence for  $t_1, t_2$  as in the statement of the lemma, we have  $\tau_1 = (t_1, 0)$  and  $\tau_2 = (t_2, 0)$ . Lemma 4.4 in this case reads

$$\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)} = \frac{\Gamma(d-s+t_2)}{\Gamma(s+t_2)} \cdot \frac{\Gamma(s+t_1)}{\Gamma(d-s+t_1)}.$$

The lemma then follows from (4.11) for  $d \geq 2$  and (4.13) for  $d = 1$ .  $\square$

**Remark 4.14.** The main point of Proposition 4.6 is that the implied constant remains bounded even as  $s$  approaches  $d - 1$  in the case  $v > 0$ . This allows us to use the bound uniformly over complementary series representations appearing in the direct integral decomposition of  $L^2(\Gamma \backslash G)$ .

The bounds from Propositions 4.5 and 4.6 allow us to restate Theorem 4.3 in terms of  $\langle \cdot, \cdot \rangle_{\mathcal{U}(v, s)}$ . This can then be applied to the matrix coefficients of irreducible unitary representations weakly contained in  $L^2(\Gamma \backslash G)$ . Retaining the notation of Theorem 4.3, we have:

**Proposition 4.7.** *There exists  $m \in \mathbb{N}$  such that for any  $\mathcal{U}(v, s)$  with a non-trivial  $M$ -invariant vector, for all  $\mathbf{u} \in \mathcal{U}(v, s)_{\tau_1}$ ,  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_2}$ , and  $t \geq 0$ ,*

$$\begin{aligned} \langle U^s(a_t) \mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)} &= e^{(s-d)t} \langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)} \\ &\quad + O_s(e^{(s-d-\eta_s)t} \|\mathbf{u}\|_{\mathcal{S}^m(v, s)} \|\mathbf{v}\|_{\mathcal{S}^m(v, s)}), \end{aligned}$$

where the implied constant is uniformly bounded over  $s$  in compact subsets of  $\mathcal{I}_v$ . Furthermore, if the  $K$ -types  $\tau_1$  and  $\tau_2$  both contain non-trivial  $M$ -invariant vectors, then the implied constant is uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ .

*Proof.* Since  $\mathbf{v} \in \mathcal{U}(v, s)_{\tau_2}$ , using the expression for  $\langle \cdot, \cdot \rangle_{\mathcal{U}(v, s)}|_{\mathcal{U}(v, s)_{\tau_2}}$ , Theorem 4.3 gives

$$\begin{aligned} \langle U^s(a_t)\mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)} &= a(v, s, \tau_2) \langle U^s(a_t)\mathbf{u}, \mathbf{v} \rangle_K \\ &= e^{(s-d)t} a(v, s, \tau_2) \langle T_{\tau_1}^{\tau_2} C_+(s)\mathbf{u}, \mathbf{v} \rangle_K \\ &\quad + O_s \left( e^{(s-d-\eta_s)t} a(v, s, \tau_2) \|T_{\tau_1}^{\tau_2}\|_K \|\mathbf{u}\|_K \|\mathbf{v}\|_{\mathcal{S}^1(K)} \right) \\ &= e^{(s-d)t} \langle T_{\tau_1}^{\tau_2} C_+(s)\mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)} \\ &\quad + O_s \left( e^{(s-d-\eta_s)t} \sqrt{\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)}} \|T_{\tau_1}^{\tau_2}\|_K \|\mathbf{u}\|_{\mathcal{U}(v, s)} \|\mathbf{v}\|_{\mathcal{S}^1(v, s)} \right). \end{aligned}$$

By Proposition 4.5, or Proposition 4.6 if  $\tau_1$  and  $\tau_2$  both have  $M$ -invariant vectors,

$$\begin{aligned} &\sqrt{\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)}} \|\mathbf{u}\|_{\mathcal{U}(v, s)} \|\mathbf{v}\|_{\mathcal{S}^1(v, s)} \\ &\ll_s \sqrt{(1 + \tau_1(\Omega_K)^d)(1 + \tau_2(\Omega_K)^d)} \|\mathbf{u}\|_{\mathcal{U}(v, s)} \|\mathbf{v}\|_{\mathcal{S}^1(v, s)} \\ &\ll \|(1 + dU^s(\Omega_K^{\lceil d/2 \rceil}))\mathbf{u}\|_{\mathcal{U}(v, s)} \|(1 + dU^s(\Omega_K^{\lceil d/2 \rceil}))\mathbf{v}\|_{\mathcal{S}^1(v, s)} \\ &\ll \|\mathbf{u}\|_{\mathcal{S}^{d+1}(v, s)} \|\mathbf{v}\|_{\mathcal{S}^{d+2}(v, s)}, \end{aligned}$$

with the implied constant uniformly bounded over  $s$  in compact subsets of  $\mathcal{I}_v$ , or  $(\frac{d}{2}, d)$ , respectively. Now Corollary 3.8 and Lemma 3.1 give

$$\begin{aligned} \|T_{\tau_1}^{\tau_2}\|_K \|\mathbf{u}\|_{\mathcal{S}^{d+1}(v, s)} \|\mathbf{v}\|_{\mathcal{S}^{d+2}(v, s)} &\leq \sqrt{\dim(\tau_1) \dim(\tau_2)} \|\mathbf{u}\|_{\mathcal{S}^{d+1}(v, s)} \|\mathbf{v}\|_{\mathcal{S}^{d+2}(v, s)} \\ &\ll \|\mathbf{u}\|_{\mathcal{S}^m(v, s)} \|\mathbf{v}\|_{\mathcal{S}^m(v, s)} \end{aligned}$$

for some  $m \in \mathbb{N}$ , completing the proof.  $\square$

**Remark 4.15.** Both Proposition 4.5 and Proposition 4.7 are expected to hold for all complementary series  $\mathcal{U}(v, s)$ . We have proved them only for those representations with  $M$ -invariant vectors since in this case slight simplifications occur in the proof of Lemma 4.4. Note that if Proposition 4.5 were proved for all complementary series representations, Proposition 4.7 would also follow automatically for all complementary series.

**Theorem 4.8.** *There exists  $m \in \mathbb{N}$  such that for any complementary series representation  $\mathcal{U}(v, s)$  containing a non-trivial  $M$ -invariant vector, for all  $\mathbf{u}, \mathbf{v} \in \mathcal{S}^m(v, s)$  and  $t \geq 0$ ,*

$$\begin{aligned} \langle U^s(a_t)\mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v, s)} &= e^{(s-d)t} \left( \sum_{\tau_1, \tau_2 \in \hat{K}} \langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{P}_{\tau_1} \mathbf{u}, \mathbf{P}_{\tau_2} \mathbf{v} \rangle_{\mathcal{U}(v, s)} \right) \\ &\quad + O_s \left( e^{(s-d-\eta_s)t} \|\mathbf{u}\|_{\mathcal{S}^m(v, s)} \|\mathbf{v}\|_{\mathcal{S}^m(v, s)} \right), \end{aligned}$$

and the sum

$$\sum_{\tau_1, \tau_2 \in \hat{K}} \langle T_{\tau_1}^{\tau_2} C_+(s) P_{\tau_1} \mathbf{u}, P_{\tau_2} \mathbf{v} \rangle_{\mathcal{U}(v,s)} \quad (4.16)$$

converges absolutely.

*Proof.* Since smooth vectors are dense in  $\mathcal{S}^m(v, s)$  for all  $m \in \mathbb{N}$ , and both sides of the inequality are continuous with respect to  $\|\cdot\|_{\mathcal{S}^m(v,s)}$ , we start by assuming that  $\mathbf{u}$  and  $\mathbf{v}$  are smooth, and decompose them according to the  $K$ -types of  $\mathcal{U}(v, s)$ :

$$\mathbf{u} = \sum_{\tau_1 \subset \mathcal{U}(v,s)} \mathbf{u}_{\tau_1}, \quad \mathbf{v} = \sum_{\tau_2 \subset \mathcal{U}(v,s)} \mathbf{v}_{\tau_2},$$

where  $\mathbf{u}_{\tau_1} = P_{\tau_1} \mathbf{u}$  and  $\mathbf{u}_{\tau_2} = P_{\tau_2} \mathbf{v}$ . By [46, Theorem 4.4.2.1],

$$\langle U^s(a_t) \mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v,s)} = \sum_{\tau_1, \tau_2} \langle U^s(a_t) \mathbf{u}_{\tau_1}, \mathbf{v}_{\tau_2} \rangle_{\mathcal{U}(v,s)},$$

with the sum converging absolutely. Applying Proposition 4.7 gives

$$\begin{aligned} \langle U^s(a_t) \mathbf{u}_{\tau_1}, \mathbf{v}_{\tau_2} \rangle_{\mathcal{U}(v,s)} &= e^{(s-d)t} \langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{u}_{\tau_1}, \mathbf{v}_{\tau_2} \rangle_{\mathcal{U}(v,s)} \\ &\quad + O_s(e^{(s-d-\eta_s)t} \|\mathbf{u}_{\tau_1}\|_{\mathcal{S}^{m'}(v,s)} \|\mathbf{v}_{\tau_2}\|_{\mathcal{S}^{m'}(v,s)}) \end{aligned} \quad (4.17)$$

for some  $m' \in \mathbb{N}$ . Hence

$$\begin{aligned} \langle U^s(a_t) \mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v,s)} &= e^{(s-d)t} \left( \sum_{\tau_1, \tau_2} \langle T_{\tau_1}^{\tau_2} C_+(s) P_{\tau_1} \mathbf{u}, P_{\tau_2} \mathbf{v} \rangle_{\mathcal{U}(v,s)} \right) \\ &\quad + O_s \left( e^{(s-d-\eta_s)t} \left( \sum_{\tau} \|\mathbf{u}_{\tau}\|_{\mathcal{S}^{m'}(v,s)} \right) \left( \sum_{\tau} \|\mathbf{v}_{\tau}\|_{\mathcal{S}^{m'}(v,s)} \right) \right). \end{aligned}$$

By Lemma 3.1, there exists  $m \geq m'$  (depending only on  $K$ ) large enough such that

$$\sum_{\tau \in \mathcal{U}(v,s)} \|\mathbf{u}_{\tau}\|_{\mathcal{S}^{m'}(v,s)} \ll \|\mathbf{u}\|_{\mathcal{S}^m(v,s)} \quad \text{and} \quad \sum_{\tau \in \mathcal{U}(v,s)} \|\mathbf{v}_{\tau}\|_{\mathcal{S}^{m'}(v,s)} \ll \|\mathbf{v}\|_{\mathcal{S}^m(v,s)}$$

where the implied constants depend only on  $K$ .

It remains to prove that the sum  $\sum_{\tau_1, \tau_2} \langle T_{\tau_1}^{\tau_2} C_+(s) P_{\tau_1} \mathbf{u}, P_{\tau_2} \mathbf{v} \rangle_{\mathcal{U}(v,s)}$  converges absolutely. Looking at an individual summand, we have

$$\begin{aligned} |\langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{u}_{\tau_1}, \mathbf{v}_{\tau_2} \rangle_{\mathcal{U}(v,s)}| &= a(v, s, \tau_2) |\langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{u}_{\tau_1}, \mathbf{v}_{\tau_2} \rangle_K| \\ &\leq a(v, s, \tau_2) \|T_{\tau_1}^{\tau_2}\|_K \cdot \|C_+(s) \mathbf{u}_{\tau_1}\|_K \cdot \|\mathbf{v}_{\tau_2}\|_K \end{aligned}$$

Since  $U^s|_K$  is unitary on  $L^2(K)$ , from the definition of  $C_+(s)$  (see Theorem 4.3), we get

$$\|C_+(s) \mathbf{u}_{\tau_1}\|_K \ll_s \|\mathbf{u}_{\tau_1}\|_K,$$

giving

$$\begin{aligned} |\langle T_{\tau_1}^{\tau_2} C_+(s) \mathbf{u}_{\tau_1}, \mathbf{v}_{\tau_2} \rangle_{\mathcal{U}(v,s)}| &\ll_s a(v, s, \tau_2) \|T_{\tau_1}^{\tau_2}\|_K \cdot \|\mathbf{u}_{\tau_1}\|_K \cdot \|\mathbf{v}_{\tau_2}\|_K \\ &\leq \sqrt{\frac{a(v, s, \tau_2)}{a(v, s, \tau_1)}} \|T_{\tau_1}^{\tau_2}\|_K \cdot \|\mathbf{u}_{\tau_1}\|_{\mathcal{U}(v,s)} \cdot \|\mathbf{v}_{\tau_2}\|_{\mathcal{U}(v,s)}. \end{aligned}$$

This expression is now bounded using Proposition 4.5, Corollary 3.8, and Lemma 3.1 as in the proof of Proposition 4.7. Lemma 3.1 then gives the desired convergence of the sum.  $\square$

**Theorem 4.9.** *There exists  $m \in \mathbb{N}$  such that for any non-spherical complementary series representation  $\mathcal{U}(v, s)$ , for all  $M$ -invariant vectors  $\mathbf{u}, \mathbf{v} \in \mathcal{S}^m(v, s)$ , and  $t \geq 0$ , we have*

$$|\langle U^s(a_t) \mathbf{u}, \mathbf{v} \rangle_{\mathcal{U}(v,s)}| \ll_s e^{(s-d-\eta_s)t} \|\mathbf{u}\|_{\mathcal{S}^m(v,s)} \|\mathbf{v}\|_{\mathcal{S}^m(v,s)},$$

where the implied constant is uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ .

*Proof.* Note that since  $\mathbf{u}, \mathbf{v}$  are both  $M$ -invariant and  $M \subset K$ ,  $P_{\tau} \mathbf{u}$  and  $P_{\tau} \mathbf{v}$  are as well for any  $K$ -type  $\tau$  of  $\mathcal{U}(v, s)$ . Proposition 4.7 then gives that the implied constant in Theorem 4.8 is uniformly bounded over  $s$  in compact subsets of  $(\frac{d}{2}, d)$ . Thus, in order to prove the theorem, it suffices to show that for any  $K$ -types  $\tau_1, \tau_2$  of  $\mathcal{U}(v, s)$  and an arbitrary  $M$ -invariant vector  $\mathbf{w}$  of  $\mathcal{U}(v, s)$ ,

$$T_{\tau_1}^{\tau_2} C_+(s) P_{\tau_1} \mathbf{w} = 0. \quad (4.18)$$

Since  $C_+(s)$  preserves the  $M$ -types of each  $K$ -type by Lemma 4.2 and  $\mathbf{w}$  is  $M$ -invariant (hence  $P_{\tau_1} \mathbf{w}$  is as well),  $C_+(s) P_{\tau_1} \mathbf{w}$  is  $M$ -invariant. Therefore (4.18) follows from Corollary 3.6.  $\square$

## 5. LEADING TERM FOR SOBOLEV FUNCTIONS

In this section, we will extend Roblin's result on mixing of the geodesic flow for continuous functions with compact support to general functions in a Sobolev space of sufficiently high order.

Recall the following theorem of Roblin [33, Theorem 3.4]:

**Theorem 5.1.** *For all  $M$ -invariant  $f_1, f_2 \in C_c(\Gamma \backslash G)$ ,*

$$\lim_{t \rightarrow +\infty} e^{(d-\delta)t} \langle \rho(a_t) f_1, f_2 \rangle = m^{\text{BR}}(f_1) m^{\text{BR}*}(f_2).$$

We note that  $m^{\text{BR}}(f)$  may be infinite for a general function  $f \in L^2(\Gamma \backslash G)$ , and hence Theorem 5.1 does not generalize to arbitrary  $L^2$ -functions.

In order to extend this theorem to Sobolev functions (which are not necessarily compactly supported), we first prove the following bound on matrix coefficients of  $L^2(\Gamma \backslash G)$ :

**Lemma 5.2.** *There exists  $m \in \mathbb{N}$  such that for all  $f_1, f_2 \in \mathcal{S}^m(\Gamma \backslash G)$  and  $t \geq 0$ ,*

$$|\langle \rho(a_t) f_1, f_2 \rangle| \ll_{\Gamma} e^{(\delta-d)t} \mathcal{S}^m(f_1) \mathcal{S}^m(f_2).$$

*Proof.* We start by assuming that  $f_1$  and  $f_2$  are  $\rho(K)$ -invariant. Using the notation of Section 2.4, the decomposition of the functions according to  $L^2(\Gamma \backslash G)_{\text{sph}} = \mathcal{B}_\delta \oplus \mathcal{W}$  reads

$$f_i = \langle f_i, \phi_0 \rangle \phi_0 + f'_i, \quad i = 1, 2,$$

where  $f'_i \in \mathcal{W}$  are  $K$ -invariant, and  $\phi_0 \in \mathcal{B}_\delta$  is the base-eigenfunction in  $L^2(\Gamma \backslash \mathbb{H}^{d+1}) = L^2(\Gamma \backslash G)^K$  of unit norm.

Since  $\mathcal{B}_\delta$  is a complementary series representation  $\mathcal{U}(1, \delta)$  (cf. Sections 2.4 and 3.2), it follows from Theorem 4.8 that for all  $t \geq 0$ ,

$$|\langle \rho(a_t) \phi_0, \phi_0 \rangle| \ll e^{(\delta-d)t}.$$

As a consequence of Theorem 2.1,  $f'_1$  and  $f'_2$  are orthogonal to any complementary series representation  $\mathcal{U}(v, s)$  with  $s_1 < s \leq \delta$ . It follows that for any  $\varepsilon > 0$ , for any  $K$ -invariant  $f_1, f_2 \in L^2(\Gamma \backslash G)$  and  $t > 0$ ,

$$|\langle \rho(a_t) f'_1, f'_2 \rangle| \ll_\varepsilon e^{(s_1-d+\varepsilon)t} \|f'_1\| \|f'_2\|$$

(see [37, Theorem 2.1, 2], [23, Proposition 5.3]).

Therefore, choosing  $0 < \varepsilon < \delta - s_1$  gives

$$|\langle \rho(a_t) f_1, f_2 \rangle| \ll e^{(\delta-d)t} (|\langle f_1, \phi_0 \rangle| |\langle f_2, \phi_0 \rangle| + \|f'_1\| \|f'_2\|) \leq e^{(\delta-d)t} \|f_1\| \|f_2\|.$$

In view of [37, Proposition 2.5], this bound extends for all  $K$ -finite functions  $f_1, f_2$  in  $L^2(\Gamma \backslash G)$ :

$$|\langle \rho(a_t) f_1, f_2 \rangle| \ll e^{(\delta-d)t} \|f_1\| \|f_2\| \sqrt{\dim(\rho(K)f_1) \dim(\rho(K)f_2)}.$$

To pass to Sobolev functions, we first observe that

$$\dim(\rho(K)P_\tau f) \leq \dim(\tau)^2$$

for all  $f \in L^2(\Gamma \backslash G)$  (cf. [20, p. 206 (1)]). Then for all  $f_1, f_2 \in \mathcal{S}^m(\Gamma \backslash G)$ ,

$$\begin{aligned} |\langle \rho(a_t) f_1, f_2 \rangle| &\leq e^{(\delta-d)t} \sum_{\tau_1, \tau_2 \in \hat{K}} \dim(\tau_1) \dim(\tau_2) \|P_{\tau_1} f_1\| \|P_{\tau_2} f_2\| \\ &\ll e^{(\delta-d)t} \|f_1\|_{\mathcal{S}^m(\Gamma \backslash G)} \|f_2\|_{\mathcal{S}^m(\Gamma \backslash G)} \end{aligned}$$

for  $m \in \mathbb{N}$  large enough, by Lemma 3.1.  $\square$

**Proposition 5.3.** *There exists  $m \in \mathbb{N}$  such that for all  $M$ -invariant  $f_1, f_2 \in \mathcal{S}^m(\Gamma \backslash G)$ ,*

$$\lim_{t \rightarrow +\infty} e^{(d-\delta)t} \langle \rho(a_t) f_1, f_2 \rangle = m^{\text{BR}}(f_1) m^{\text{BR}*}(f_2).$$

*Proof.* Let  $m > d(d+1)/4$  satisfy the conclusion of Lemma 5.2. For simplicity, we write  $\|f\|_{\mathcal{S}^m} = \|f\|_{\mathcal{S}^m(\Gamma \backslash G)}$ . Using the density of  $C_c^\infty(\Gamma \backslash G)^M$  in  $\mathcal{S}^m(\Gamma \backslash G)^M$ , given  $\varepsilon > 0$ , there exist  $f_1^\varepsilon, f_2^\varepsilon \in C_c^\infty(\Gamma \backslash G)$  such that

$$\|f_i - f_i^\varepsilon\|_{\mathcal{S}^m} \leq \varepsilon \quad i = 1, 2.$$

We then write, using Lemma 5.2,

$$\begin{aligned} & e^{(d-\delta)t} \langle \rho(a_t) f_1, f_2 \rangle \\ &= e^{(d-\delta)t} \left( \langle \rho(a_t) f_1^\varepsilon, f_2^\varepsilon \rangle + \langle \rho(a_t) (f_1 - f_1^\varepsilon), f_2 \rangle + \langle \rho(a_t) f_1^\varepsilon, f_2 - f_2^\varepsilon \rangle \right) \\ &= e^{(d-\delta)t} \langle \rho(a_t) f_1^\varepsilon, f_2^\varepsilon \rangle + O(\varepsilon(\|f_1\|_{\mathcal{S}^m} + \|f_2\|_{\mathcal{S}^m})). \end{aligned}$$

By Theorem 5.1, we have

$$\lim_{t \rightarrow \infty} e^{(d-\delta)t} \langle \rho(a_t) f_1^\varepsilon, f_2^\varepsilon \rangle = m^{\text{BR}}(f_1^\varepsilon) m^{\text{BR}*}(f_2^\varepsilon).$$

So

$$\limsup_{t \rightarrow \infty} e^{(d-\delta)t} \langle \rho(a_t) f_1, f_2 \rangle = m^{\text{BR}}(f_1^\varepsilon) m^{\text{BR}*}(f_2^\varepsilon) + O(\varepsilon(\|f_1\|_{\mathcal{S}^m} + \|f_2\|_{\mathcal{S}^m})).$$

Since  $m > d(d+1)/4$ , we may apply Lemma 2.2 to get

$$m^{\text{BR}}(f_1^\varepsilon) m^{\text{BR}*}(f_2^\varepsilon) = m^{\text{BR}}(f_1) m^{\text{BR}*}(f_2) + O(\varepsilon(\|f_1\|_{\mathcal{S}^m} + \|f_2\|_{\mathcal{S}^m})).$$

Therefore

$$\limsup_{t \rightarrow \infty} e^{(d-\delta)t} \langle \rho(a_t) f_1, f_2 \rangle = m^{\text{BR}}(f_1) m^{\text{BR}*}(f_2) + O(\varepsilon(\|f_1\|_{\mathcal{S}^m} + \|f_2\|_{\mathcal{S}^m})).$$

Since  $\varepsilon > 0$  was arbitrary, we in fact have

$$\limsup_{t \rightarrow \infty} e^{(d-\delta)t} \langle \rho(a_t) f_1, f_2 \rangle = m^{\text{BR}}(f_1) m^{\text{BR}*}(f_2).$$

An analogous calculation with  $\liminf$  in place of  $\limsup$  proves the proposition.  $\square$

## 6. PROOF OF THEOREM 1.1

In this final section, we prove Theorem 1.1. Recall the isomorphism

$$(\rho, L^2(\Gamma \backslash G)) \cong \int_{\mathbf{Z}}^{\oplus} (\pi_\zeta, \mathcal{H}_\zeta) d\mu_{\mathbf{Z}}(\zeta).$$

Fix arbitrary  $0 < r < \delta - s_1$ . We partition  $\mathbf{Z}$  as

$$\mathbf{Z} = \mathbf{Z}_r^- \cup \mathbf{Z}_r^+,$$

where

$$\mathbf{Z}_r^- = \left\{ \zeta \in \mathbf{Z} : (\pi_\zeta, \mathcal{H}_\zeta) \cong \mathcal{U}(v, s), \text{ where } v \in \hat{M} \text{ and } s \in [s_1 + r, \delta] \right\}$$

and  $\mathbf{Z}_r^+ = \mathbf{Z} \setminus \mathbf{Z}_r^-$  (cf. (2.2)). This partition is then used to decompose  $(\rho, L^2(\Gamma \backslash G))$  as

$$(\rho, L^2(\Gamma \backslash G)) = (\rho, L^2(\Gamma \backslash G)^-) \oplus (\rho, L^2(\Gamma \backslash G)^+),$$

where

$$(\rho, L^2(\Gamma \backslash G)^\pm) \cong \int_{\mathbf{Z}^\pm}^{\oplus} (\pi_\zeta, \mathcal{H}_\zeta) d\mu_{\mathbf{Z}}(\zeta)$$

with the dependency on  $r$  being slightly suppressed. Recall that  $\mathcal{B}_\delta$  occurs as a subrepresentation of  $(\rho, L^2(\Gamma \backslash G)^-)$  and  $\mathcal{B}_\delta \cong \mathcal{U}(1, \delta)$  (cf. Sections 2.4 and

3.2). We may further decompose  $(\rho, L^2(\Gamma \backslash G)^-)$  accordingly: let  $L^2(\Gamma \backslash G)_0^-$  be the orthogonal complement of  $\mathcal{B}_\delta$  in  $L^2(\Gamma \backslash G)^-$ . We thus have

$$(\rho, L^2(\Gamma \backslash G)) = (\rho, \mathcal{B}_\delta) \oplus (\rho, L^2(\Gamma \backslash G)_0^-) \oplus (\rho, L^2(\Gamma \backslash G)^+).$$

Note that Theorem 2.1 and the duality between eigenfunctions of the Laplacian on  $\mathcal{M} = \Gamma \backslash G / K$  and spherical representations in the decomposition of  $L^2(\Gamma \backslash G)$  imply that no spherical representation is weakly contained in  $(\rho, L^2(\Gamma \backslash G)_0^-)$ . At the level of functions, we write

$$f_i = f_i^0 + f_i^- + f_i^+, \quad i = 1, 2$$

where  $f_i^0 \in \mathcal{B}_\delta$ ,  $f_i^- \in L^2(\Gamma \backslash G)_0^-$ , and  $f_i^+ \in L^2(\Gamma \backslash G)^+$ .

The matrix coefficients we are interested in now decompose as

$$\langle \rho(a_t) f_1, f_2 \rangle = \langle \rho(a_t) f_1^0, f_2^0 \rangle + \langle \rho(a_t) f_1^-, f_2^- \rangle + \langle \rho(a_t) f_1^+, f_2^+ \rangle. \quad (6.1)$$

We deal with the three summands in turn: by construction,  $(\rho, L^2(\Gamma \backslash G)^+)$  does not weakly contain any  $\mathcal{U}(v, s)$  with  $s > s_1 + r$  (this also uses the fact  $(\rho, L^2(\Gamma \backslash G))$  does not weakly contain any  $\mathcal{U}(v, s)$  with  $s > \delta$ ; cf. [27, Proposition 3.23]).

We assume that  $m$  is large enough so that all the results of the previous sections hold for  $f_1, f_2 \in \mathcal{S}^m(\Gamma \backslash G)$ . By [23, Proposition 5.3] (cf. also [27, Proposition 3.29] or [37, Theorem 2.1] combined with the argument from the proof of Lemma 5.2), we have that for all  $\xi > 0$ ,

$$\begin{aligned} |\langle \rho(a_t) f_1^+, f_2^+ \rangle| &\ll_\xi e^{(s_1 - d + r + \xi)t} \|f_1^+\|_{\mathcal{S}^m(\Gamma \backslash G)} \|f_2^+\|_{\mathcal{S}^m(\Gamma \backslash G)} \\ &\leq e^{(s_1 - d + r + \xi)t} \|f_1\|_{\mathcal{S}^m(\Gamma \backslash G)} \|f_2\|_{\mathcal{S}^m(\Gamma \backslash G)}. \end{aligned} \quad (6.2)$$

We will now use Theorem 4.9 to bound  $\langle \rho(a_t) f_1^-, f_2^- \rangle$ . Let  $\tilde{Z}_r^- \subset Z_r^-$  be such that  $(\rho, L^2(\Gamma \backslash G)_0^-) \cong \int_{\tilde{Z}_r^-}^\oplus (\pi_\zeta, \mathcal{H}_\zeta) d\mu_Z(\zeta)$ . The corresponding decomposition of the functions  $f_1^-, f_2^-$  reads  $f_i^- = \int_{\tilde{Z}_\varepsilon^-} (f_i^-)_\zeta d\mu_Z(\zeta)$  ( $i = 1, 2$ ). Since  $f_i$  is  $M$ -invariant, so is  $f_i^-$  and  $\mu_Z$ -a. e.  $(f_i^-)_\zeta$ . The matrix coefficient  $\langle \rho(a_t) f_1^-, f_2^- \rangle$  is now written as

$$\langle \rho(a_t) f_1^-, f_2^- \rangle = \int_{\tilde{Z}_\varepsilon^-} \langle \pi_\zeta(a_t) (f_1^-)_\zeta, (f_2^-)_\zeta \rangle_{\mathcal{H}_\zeta} d\mu_Z(\zeta).$$

Each  $(\pi_\zeta, \mathcal{H}_\zeta)$  is isomorphic to some  $\mathcal{U}(v, s)$  with  $v$  *non-trivial* and  $s \in [s_1 + r, \delta]$ . Setting

$$\lambda(\delta, r) = \max_{s \in [s_1 + r, \delta]} s - d - \eta_s$$

where  $\eta_s = \min(2s-d, 1)$ , we apply Theorem 4.9 to each  $\langle \pi_\zeta(a_t)(f_1^-)_\zeta, (f_2^-)_\zeta \rangle_{\mathcal{H}_\zeta}$  and obtain

$$\begin{aligned}
& |\langle \rho(a_t)f_1^-, f_2^- \rangle| \ll_{r,\delta} e^{\lambda(\delta,r)t} \int_{\tilde{\mathcal{Z}}_r^-} \|(f_1^-)_\zeta\|_{\mathcal{S}^m(\mathcal{H}_\zeta)} \|(f_2^-)_\zeta\|_{\mathcal{S}^m(\mathcal{H}_\zeta)} d\mu_{\mathbf{Z}}(\zeta) \\
& \leq e^{\lambda(\delta,r)t} \sqrt{\int_{\tilde{\mathcal{Z}}_r^-} \|(f_1^-)_\zeta\|_{\mathcal{S}^m(\mathcal{H}_\zeta)}^2 d\mu_{\mathbf{Z}}(\zeta)} \sqrt{\int_{\tilde{\mathcal{Z}}_r^-} \|(f_2^-)_\zeta\|_{\mathcal{S}^m(\mathcal{H}_\zeta)}^2 d\mu_{\mathbf{Z}}(\zeta)} \\
& = e^{\lambda(\delta,r)t} \|f_1^-\|_{\mathcal{S}^m(\Gamma \setminus G)} \|f_2^-\|_{\mathcal{S}^m(\Gamma \setminus G)} \\
& \leq e^{\lambda(\delta,r)t} \|f_1\|_{\mathcal{S}^m(\Gamma \setminus G)} \|f_2\|_{\mathcal{S}^m(\Gamma \setminus G)}, \tag{6.3}
\end{aligned}$$

where the implied constant remains uniformly bounded as  $r \rightarrow 0$ . The remaining term in the left-hand side is  $\langle \rho(a_t)f_1^0, f_2^0 \rangle$ . Using the fact that  $(\rho, \mathcal{B}_\delta)$  is isomorphic to  $\mathcal{U}(1, \delta)$ , applying Proposition 4.8 gives

$$\begin{aligned}
\langle \rho(a_t)f_1^0, f_2^0 \rangle & = \Phi(f_1^0, f_2^0) e^{(\delta-d)t} + O_\delta \left( e^{(\delta-d-\eta_\delta)t} \|f_1^0\|_{\mathcal{S}^m(\Gamma \setminus G)} \|f_2^0\|_{\mathcal{S}^m(\Gamma \setminus G)} \right) \\
& = \Phi(f_1^0, f_2^0) e^{(\delta-d)t} + O_\delta \left( e^{(\delta-d-\eta_\delta)t} \|f_1\|_{\mathcal{S}^m(\Gamma \setminus G)} \|f_2\|_{\mathcal{S}^m(\Gamma \setminus G)} \right), \tag{6.4}
\end{aligned}$$

where  $\Phi(f_1^0, f_2^0)$  is given by the sum (4.16) under the aforementioned isomorphism.

Note that

$$\begin{aligned}
\beta & := \min \{ \eta_\delta, \delta - s_1 - r - \xi, \delta - d - \lambda(\delta, r) \} \\
& = \min \left\{ 1, 2\delta - d, \delta - s_1 - r - \xi, \min_{s \in [s_1+r, \delta]} (\delta - s + \min\{1, 2s - d\}) \right\} \\
& \geq \min \left\{ 1, 2\delta - d, \delta - s_1 - r, \min_{s \in [s_1+r, \delta]} (\delta - s + \min\{1, 2s - d\}) \right\} - \xi \\
& = \min \{ 1, 2\delta - d, \delta - s_1 - r, \delta + s_1 + r - d \} - \xi \\
& = \min \{ 1, \delta - s_1 - r, \delta + s_1 + r - d \} - \xi \\
& = \min \{ 1, \delta - s_1 - r \} - \xi \\
& \geq \min \{ 1, \delta - s_1 \} - (\xi + r).
\end{aligned}$$

Entering (6.2), (6.3), and (6.4) into (6.1) gives

$$e^{(d-\delta)t} \langle \rho(a_t)f_1, f_2 \rangle = \Phi(f_1^0, f_2^0) + O_{r,\xi} (e^{-\beta t} \|f_1\|_{\mathcal{S}^m(\Gamma \setminus G)} \|f_2\|_{\mathcal{S}^m(\Gamma \setminus G)}).$$

Writing  $\eta = \min(1, \delta - s_1)$  and  $\varepsilon = \xi + r$ , we thus have

$$e^{(d-\delta)t} \langle \rho(a_t)f_1, f_2 \rangle = \Phi(f_1^0, f_2^0) + O_\varepsilon (e^{-(\eta-\varepsilon)t} \|f_1\|_{\mathcal{S}^m(\Gamma \setminus G)} \|f_2\|_{\mathcal{S}^m(\Gamma \setminus G)}).$$

Choosing  $0 < \varepsilon < \eta$  gives

$$\lim_{t \rightarrow \infty} e^{(d-\delta)t} \langle \rho(a_t)f_1, f_2 \rangle = \Phi(f_1^0, f_2^0),$$

so  $\Phi(f_1^0, f_2^0) = m^{\text{BR}}(f_1) m^{\text{BR}*}(f_2)$  by Proposition 5.3, completing the proof of Theorem 1.1.



**Remark 6.5.** In the case when  $\Gamma$  is a lattice, i.e., when  $\delta = d$ , the above proof leads to the claim in Remark 1.1 (1) as follows. Firstly, we have  $\langle \rho(a_t)f_1^0, f_2^0 \rangle = \int f_1 dx \cdot \int f_2 dx$  for any  $t \in \mathbb{R}$ . As before, we use (6.2) to bound  $\langle \rho(a_t)f_1^+, f_2^+ \rangle$ . Now, since complementary series representations  $\mathcal{U}(v, s)$  with *non-trivial*  $v$  exist only for  $s \leq d - 1$ , the constant  $\lambda(d, r)$  appearing in the bound (6.3) for  $|\langle \rho(a_t)f_1^-, f_2^- \rangle|$  is at most  $\alpha := \max\{s - d - \eta_s : s \in [s_1 + r, d - 1]\}$ . We note that  $\alpha < -2 + r$ , and if, in addition,  $\{s > s_1 + 1 : \mathcal{U}(v, s) \text{ is weakly contained in } L^2(\Gamma \backslash G)\} = \emptyset$ , then one can take  $\lambda(d, r) < s_1 - d + r$ . Hence combining these bounds, we obtain Remark 1.1 (1).

## REFERENCES

- [1] T. Aubin. *Nonlinear Analysis on Manifolds*, Grundlehren Math. Wiss. 252, Springer (1982).
- [2] M. Babillot. *On the mixing property for hyperbolic systems*, Israel J. Math. 129 (2002), 61-76.
- [3] Y. Benoist and H. Oh. *Effective equidistribution of  $S$ -integral points on symmetric varieties*, Ann. Inst. Fourier 62 (2012), 1889-1942.
- [4] T. Bröcker and T. tom Dieck. *Representations of compact Lie groups*, GTM 98, Springer (1985).
- [5] J. Bourgain, A. Gamburd, and P. Sarnak. *Generalization of Selberg's 3/16 theorem and Affine sieve*, Acta Math. 20 (2011), 255-290.
- [6] J. Bourgain, A. Gamburd, and P. Sarnak. *Affine linear sieve, expanders, and sum-product*, Inventiones 179 (2010), 559-644.
- [7] J. Bourgain, A. Kontorovich, and P. Sarnak. *Sector estimates for hyperbolic isometries*, Geom. Funct. Anal. 20 (2010), 1175-1200.
- [8] N. Bergeron and L. Clozel. *Quelques conséquences des travaux d'Arthur pour le spectre et la topologie des variétés hyperboliques*. Invent. Math. 192 (2013), 505-532.
- [9] M. Burger and P. Sarnak. *Ramanujan duals II*. Invent. Math. 106 (1991), 1-11.
- [10] M. Eguchi, S. Koizumi, and M. Mamiuda. *The expressions of the Harish-Chandra  $C$ -functions of semisimple Lie groups  $Spin(n,1)$ ,  $SU(n,1)$* , J. Math. Soc. Japan 51 (1999), 955-985.
- [11] W. Duke, Z. Rudnick, and P. Sarnak. *Density of integer points on affine homogeneous varieties*, Duke Math. J. 71 (1993), 143-179.
- [12] A. Eskin and C. T. McMullen. *Mixing, counting, and equidistribution in Lie groups*, Duke Math. J. 71 (1993), 181-209.
- [13] A. Gamburd. *On the spectral gap for infinite index "congruence" subgroups of  $SL_2(\mathbb{Z})$* , Isr. J. Math. 123 (2002), 157-200.
- [14] R. Goodman and N. Wallach. *Symmetry, representations, and invariants*, GTM 255, Springer (2009).
- [15] W. He and N. de Saxcé. *Linear random walks on the torus*, preprint: arXiv:1910.1342.
- [16] T. Hirai. *On irreducible representations of the Lorentz group of  $n$ -th order*, Proc. Japan Acad. 38 (1962), 258-262.
- [17] J. Kečkić and P. Vasić. *Some inequalities for the gamma function*, Pub. Inst. Math. 11 (1971), 107-114.
- [18] D. Kelmer and H. Oh. *Exponential mixing and shrinking targets for geodesic flow on geometrically finite hyperbolic manifolds*, arXiv:1812.05251, to appear in JMD.
- [19] I. Kim. *Counting, Mixing and Equidistribution of horospheres in geometrically finite rank one locally symmetric manifolds*, J. Reine Angew. Math. 704 (2015), 85-133.

- [20] A. Knapp. *Representation Theory of semisimple Lie Groups*, Princeton University Press (1986).
- [21] A. Knapp. *Lie Groups Beyond an Introduction* Second Edition, Birkhäuser, 2002.
- [22] A. Knapp and E. Stein *Intertwining Operators for Semisimple Groups*, Ann. of Math., Second Series, Vol. 93 (1971), 489-578.
- [23] A. Kontorovich and H. Oh. *Apollonian circle packings and closed horospheres on hyperbolic 3-manifolds*, J. Amer. Math. Soc. 24 (2011), 603-648.
- [24] P. Lax and R. Phillips. *The asymptotic distribution of lattice points in Euclidean and non-Euclidean spaces*, J. Funct. Anal. 46 (1982), 280-350.
- [25] M. Magge. *Quantitative spectral gap for thin groups of hyperbolic isometries*, JEMS, 17 (2015), 151-187
- [26] G. Margulis, A. Mohammadi and H. Oh. *Closed geodesics and holonomies for Kleinian manifolds*, Geom. Funct. Anal. 24 (2014), 1608-1636.
- [27] A. Mohammadi and H. Oh. *Matrix coefficients, counting and primes for orbits of geometrically finite groups*, J. Eur. Math. Soc. (JEMS) 17 (2015), 837-897.
- [28] C. Moore. *Exponential decay of correlation coefficients for geodesic flows*, in: Group Representations, Ergodic Theory, Operator Algebras, and Mathematical Physics (eds C. C. Moore), MSRI Pub. Vol 6, Springer (1987).
- [29] H. Oh and D. Winter. *Uniform exponential mixing and resonance free regions for convex co-compact congruence subgroups of  $SL_2(\mathbb{Z})$* , J. Amer. Math. Soc. 29 (2016), 1069-1115.
- [30] S. J. Patterson. *The limit set of a Fuchsian group*, Acta Math. 136 (1976), 241-273.
- [31] J.-P. Otal and M. Peigné. *Principe variationnel et groupes Kleinien*, Duke Math. J. 125 (2004), 15-44.
- [32] M. Ratner. *The rate of mixing for geodesic and horocycle flows*, Ergod. Th. & Dynam. Sys. 7 (1987), 267-288.
- [33] T. Roblin. *Ergodicité et équidistribution en courbure négative*, Mém. Soc. Math. France (N.S.) 95 (2003), vi+96 pp.
- [34] P. Sarkar. *Generalization of Selberg's 3/16 theorem for convex cocompact thin subgroups of  $SO(n, 1)$* , preprint, arXiv:2006.07787.
- [35] P. Sarkar and D. Winter. *Exponential mixing of frame flows for convex cocompact hyperbolic manifolds*, preprint: arXiv:2004.14551.
- [36] P. Sarnak. *Notes on the Generalized Ramanujan Conjectures*. Clay Math. Proceedings 4 (2005).
- [37] Y. Shalom. *Rigidity, unitary representations of semisimple groups, and fundamental groups of manifolds with rank one transformation group*, Ann. of Math. 152 (2000), 113-182.
- [38] L. Stoyanov. *Spectra of Ruelle transfer operators for axiom A flows*, Nonlinearity 24 (2011), 1089-1120.
- [39] D. Sullivan. *Entropy, Hausdorff measures old and new, and limit sets of geometrically finite Kleinian groups*, Acta Math. 153 (1984), 259-277.
- [40] D. Sullivan. *The density at infinity of a discrete group of hyperbolic motions*, Inst. Hautes Études Sci. Publ. Math. 50 (1979), 171-202.
- [41] A. Salehi-Golsefidy and P. Varjú. *Expansion in perfect groups*, Geom. Funct. Anal. 22 (2012), 1832-1891.
- [42] E. Thieleker. *On the quasi-simple irreducible representations of the Lorentz groups*, Trans. Amer. Math. Soc. 179 (1973), 465-505.
- [43] I. Vinogradov. *Effective bisector estimate with applications to Apollonian circle packings*, Int. Math. Res. Not. IMRN 12 (2014), 3217-3262.
- [44] N. Wallach. *Real Reductive Groups I*, Academic Press (1988).
- [45] N. Wallach. *Real Reductive Groups II*, Academic Press (1992).
- [46] G. Warner. *Harmonic Analysis on Semi-Simple Lie Groups I*, Grundlehren Math. Wiss. 188, Springer (1972).

- [47] G. Warner. *Harmonic Analysis on Semi-Simple Lie Groups II*, Grundlehren Math. Wiss. 189, Springer (1972).
- [48] D. Želobenko. *Compact Lie groups and their representations*, Translations of Mathematical Monographs, Vol. 40, American Mathematical Society, Providence, R.I. (1973).

DEPARTMENT OF MATHEMATICS, YALE UNIVERSITY, NEW HAVEN, CT 06520  
*E-mail address:* `samuel.edwards@yale.edu`

DEPARTMENT OF MATHEMATICS, YALE UNIVERSITY, NEW HAVEN, CT 06520 AND  
KOREA INSTITUTE FOR ADVANCED STUDY, SEOUL, KOREA  
*E-mail address:* `hee.oh@yale.edu`