

THE HARMONIC ANALYSIS OF LATTICE COUNTING  
ON REAL SPHERICAL SPACES

BERNHARD KRÖTZ

*Universität Paderborn, Institut für Mathematik  
Warburger Straße 100, 33098 Paderborn, Deutschland*

EITAN SAYAG

*Department of Mathematics, Ben Gurion University of the Negev  
P.O.B. 653, Be'er Sheva 84105, Israel*

HENRIK SCHLICHTKRULL

*University of Copenhagen, Department of Mathematics  
Universitetsparken 5, DK-2100 Copenhagen Ø, Denmark*

ABSTRACT. By the collective name of *lattice counting* we refer to a setup introduced in [10] that aim to establish a relationship between arithmetic and randomness in the context of affine symmetric spaces. In this paper we extend the geometric setup from symmetric to real spherical spaces and continue to develop the approach with harmonic analysis which was initiated in [10].

---

*E-mail addresses:* bkroetz@math.uni-paderborn.de,  
eitan.sayag@gmail.com, schlicht@math.ku.dk.

*Date:* August 14, 2014.

The first author was supported by ERC Advanced Investigators Grant HARG 268105. The second author was partially supported by ISF 1138/10 and ERC 291612.

## 1. INTRODUCTION

**1.1. Lattice counting.** Let us recall from Duke, Rudnick and Sarnak [10] the setup of lattice counting on a homogeneous space  $Z = G/H$ . Here  $G$  is an algebraic real reductive group and  $H < G$  an algebraic subgroup such that  $Z$  carries an invariant measure. Further we are given a lattice  $\Gamma < G$  such that its trace  $\Gamma_H := \Gamma \cap H$  in  $H$  is a lattice in  $H$ .

Attached to invariant measures  $dh$  and  $dg$  on  $H$  and  $G$  we obtain an invariant measure  $d(gH)$  on  $Z$  via Weil-integration:

$$\int_Z \left( \int_H f(gh)dh \right) d(gH) = \int_G f(g) dg \quad (f \in C_c(G)).$$

Likewise the measures  $dg$  and  $dh$  give invariant measures  $d(g\Gamma)$  and  $d(h\Gamma_H)$  on  $Y := G/\Gamma$  and  $Y_H := H/\Gamma_H$ . We pin down the measures  $dg$  and  $dh$  and hence  $d(gH)$  by the request that  $Y$  and  $Y_H$  have volume one.

Further we are given a family  $\mathcal{B}$  of “balls”  $B_R \subset Z$  depending on a parameter  $R \geq 0$ . At this point we are rather imprecise about the structure of these balls and content us with the property that they constitute an exhausting family of compact sets as  $R \rightarrow \infty$ .

Let  $z_0 = H \in Z$  be the standard base point. The *lattice counting problem* for  $\mathcal{B}$  consists of the determination of the asymptotic behavior of the density of  $\Gamma \cdot z_0$  in balls  $B_R \subset Z$ , as the radius  $R \rightarrow \infty$ . By *main term counting* for  $\mathcal{B}$  we understand the statement that the asymptotic density is 1. More precisely, with

$$N_R(\Gamma, Z) := \#\{\gamma \in \Gamma/\Gamma_H \mid \gamma \cdot z_0 \in B_R\}$$

and  $|B_R| := \text{vol}_Z(B_R)$  we say that main term counting holds if

$$(1.1) \quad N_R(\Gamma, Z) \sim |B_R| \quad (R \rightarrow \infty).$$

**1.2. Relevant previous works.** The main term counting was established in [10] for symmetric spaces  $G/H$  and certain families of balls, for lattices with  $Y_H$  compact. In subsequent work Eskin and McMullen [11] removed the obstruction that  $Y_H$  is compact and presented an ergodic approach. Later Eskin, Mozes and Shah [12] refined the ergodic methods and discovered that main term counting holds for a wider class of reductive spaces: For reductive algebraic groups  $G, H$  defined over  $\mathbb{Q}$  and arithmetic lattices  $\Gamma < G(\mathbb{Q})$  it is enough to request that  $H$  is not contained in a proper parabolic subgroup of  $G$  which is defined over  $\mathbb{Q}$ . In particular all maximal reductive subgroups have this property.

In these works the balls  $B_R$  are constructed as follows. All spaces considered are affine in the sense that there exists a  $G$ -equivariant

embedding of  $Z$  into the representation module  $V$  of a rational representation of  $G$ . For any such embedding and any norm on the vector space  $V$ , one then obtains a family of balls  $B_R$  on  $Z$  by intersection with the metric balls in  $V$ . For symmetric spaces all families of balls produced this way are suitable for the lattice counting, but on a more general reductive homogeneous space one needs to be more careful. In [12] a technical condition called “focusing balls” is requested.

**1.3. Real spherical spaces.** In this paper we investigate the lattice counting for a real spherical space  $Z$ , that is, it is requested that the action of a minimal parabolic subgroups  $P < G$  on  $Z$  admits an open orbit. In addition we assume that  $H$  is reductive and remark that for spherical spaces this is automatically satisfied if the Lie algebra  $\mathfrak{h}$  of  $H$  is self-normalizing.

Our approach is based on spectral theory and is a natural continuation to [10]. We consider a particular type of balls which are intrinsically defined by the geometry of  $Z$  (and thus not related to a particular representation  $V$  as before).

1.3.1. *Factorization of spherical spaces.* In the spectral approach it is of relevance to get a control over intermediate subgroups  $H < H^* < G$  which arise in the following way: Given a unitary representation  $(\pi, \mathcal{H})$  one looks at the smooth vectors  $\mathcal{H}^\infty$  and its continuous dual  $\mathcal{H}^{-\infty}$ , the distribution vectors. The space  $(\mathcal{H}^{-\infty})^H$  of  $H$ -invariant distribution vectors is of fundamental importance. For all pairs  $(v, \eta) \in \mathcal{H}^\infty \times (\mathcal{H}^{-\infty})^H$  one obtains a smooth function on  $Z$ , a *generalized matrix-coefficient*, via

$$(1.2) \quad m_{v,\eta}(z) = \eta(g^{-1} \cdot v) \quad (z = gH \in Z).$$

The functions (1.2) are the building blocks for the harmonic analysis on  $Z$ . The stabilizer  $H_\eta$  in  $G$  of  $\eta \in (\mathcal{H}^{-\infty})^H$  is a closed subgroup which contains  $H$ , but in general it can be larger than  $H$  even if  $\pi$  is non-trivial.

Let us call  $Z^* = G/H^*$  a *factorization* of  $Z$  if  $H < H^*$  and  $Z^*$  is unimodular. For a general real spherical space  $Z$  the homogeneous spaces  $Z_\eta = G/H_\eta$  can happen to be non-unimodular (see [18] for  $H$  the Iwasawa  $N$ -subgroup). However there is a large subclass of real spherical spaces which behave well under factorization. Let us call a factorization co-compact if  $H^*/H$  is compact and *basic* if  $H^*$  is of the form  $H_I := HI$  for a normal subgroup  $I \triangleleft G$ . Finally we call a factorization *weakly basic* if it is obtained by a composition of a basic and a co-compact factorization.

1.3.2. *Wavefront spherical spaces.* A real spherical space is called *wavefront* if the attached compression cone is a quotient of a closed Weyl-chamber – see [17]. Many real spherical spaces are wavefront: all symmetric spaces and all Gross-Prasad type spaces  $G \times H/H$  are wavefront<sup>1</sup>. The terminology *wavefront* originates from [23] because wavefront real spherical spaces satisfy the “wavefront lemma” of Eskin-McMullen (see [11], [17]) which is fundamental in the approach of [11] to lattice counting.

On the geometric side wavefront real spherical spaces enjoy the following property from [18]: All  $Z_\eta$  are unimodular and the factorizations of the type  $Z_\eta$  are precisely the weakly basic factorizations of  $Z$ .

On the spectral level wavefront real spherical spaces are distinguished by the following integrability property, also from [18]: The generalized matrix coefficients  $m_{v,\eta}$  of (1.2) belong to  $L^p(Z_\eta)$  for some  $1 \leq p < \infty$  only depending on  $\pi$  and  $\eta$ .

1.3.3. *Main term counting.* In the theorem below we assume that  $Z$  is a wavefront real spherical space of reductive type, for which all factorizations are basic. For simplicity we also assume that all compact normal subgroups of  $G$  are finite.

Using soft techniques from harmonic analysis and a general property of decay from [20], our first result (see Section 5) is:

**Theorem A.** Let  $Z = G/H$  be as above, and assume that  $Y = G/\Gamma$  is compact. Then main term counting (1.1) holds.

Since wavefront real spherical spaces satisfy the wavefront lemma by [17], Section 6, this theorem could also be derived with the ergodic method of [11]. In the current context the main point is thus the proof by harmonic analysis.

To remove the assumption that  $Y$  is compact and to obtain error term bounds for the lattice counting problem we need to apply more sophisticated tools from harmonic analysis. This will be discussed in the next paragraph with some extra assumptions on  $G/H$ .

1.4. **Error Terms.** The problem of determining the error term in counting problems is notoriously difficult and in many cases relies on deep arithmetic information. Sometimes, like in the Gauss circle problem, some error term is easy to establish but getting an optimal error term is a very difficult problem.

---

<sup>1</sup>Also, if  $Z$  is complex, then of the 78 cases in the list of [4], the non-wavefront cases are (11), (24), (25), (27), (39-50), (60), (61)

We restrict ourselves to the cases where the cycle  $H/\Gamma_H$  is compact.<sup>2</sup> To simplify the exposition here we assume in addition that  $\Gamma < G$  is irreducible, i.e. there do not exist non-trivial normal subgroups  $G_1, G_2$  of  $G$  and lattices  $\Gamma_i < G_i$  such that  $\Gamma_1\Gamma_2$  has finite index in  $\Gamma$ .

The error we study is measure theoretic in nature, and will be denoted here as  $\text{err}(R, \Gamma)$ . Thus,  $\text{err}(R, \Gamma)$  measures the deviation of two measures on  $Y = \Gamma \backslash G$ , the counting measure arising from lattice points in a ball of radius  $R$ , and the invariant measure on  $Y$ . It is easy to compare this error term with the pointwise error  $\text{err}_{pt}(R, \Gamma) = |N_R(\Gamma, Z) - |B_R||$ , see Remark 7.2.

To formulate our result we introduce the exponent  $p_H(\Gamma)$ , which measures the worst  $L^p$ -behavior of any generalized matrix coefficient associated with a spherical unitary representation  $\pi$ , which is  $H$ -distinguished and occurs in the automorphic spectrum of  $L^2(\Gamma \backslash G)$ . We first state our result for the non-symmetric case of triple product spaces, which is Theorem 8.2 from the body of the paper.

**Theorem B.** Let  $Z = G_0^3/\text{diag}(G_0)$  for  $G_0 = \text{SO}_e(1, n)$  and assume that  $H/\Gamma_H$  is compact. For all  $p > p_H(\Gamma)$  there exists a  $C = C(p) > 0$  such that

$$\text{err}(R, \Gamma) \leq C|B_R|^{-\frac{1}{(6n+3)p}}$$

for all  $R \geq 1$ . (In particular, main term counting holds in this case).

To the best of our knowledge this is the first error term obtained for a non-symmetric space. The crux of the proof is locally uniform comparison between  $L^p$  and  $L^\infty$  norms of generalized matrix coefficients  $m_{v, \eta}$  which is achieved by applying the model of [3] for the triple product functional  $\eta$  in spherical principal series.

It is possible to obtain error term bounds under a certain technical hypothesis introduced in Section 6 and referred to as Hypothesis A. This hypothesis in turn is implied by a conjecture on the analytic structure of families of Harish-Chandra modules which we explain in Section 9.1. The conjecture and hence the hypothesis appear to be true for symmetric spaces but requires quite a technical tour de force. In general, the techniques currently available do not allow for an elegant and efficient solution. Under this hypothesis we show that:

**Theorem C.** Let  $Z$  be wavefront real spherical space for which Hypothesis A is valid. Assume

- $\Gamma_H = H \cap \Gamma$  is co-compact in  $H$ .
- $p > p_H(\Gamma)$

---

<sup>2</sup>After a theory for regularization of  $H$ -periods of Eisenstein series is developed, one can drop this assumption.

- $k > \frac{\text{rank}(G/K)+1}{2} \dim(G/K) + 1$

Then, there exists a constant  $C = C(p, k) > 0$  such that

$$\text{err}(R, \Gamma) \leq C |B_R|^{-\frac{1}{(2k+1)p}}$$

for all  $R \geq 1$ . Moreover, if  $Y = \Gamma \backslash G$  is compact one can replace the third condition by  $k > \dim(G/K) + 1$ .

The existence of a non-quantitative error term for symmetric spaces was established in [1] and improved in [13].

We note that in case of the hyperbolic plane our error term is still far from the quality of the bound of A. Selberg. This is because we only use a weak version of the trace formula, namely Weyl's law, and use simple soft Sobolev bounds between eigenfunctions on  $Y$ .

## 2. REDUCTIVE HOMOGENEOUS SPACES

In this section we review a few facts on reductive homogeneous spaces: the Mostow decomposition, the associated geometric balls and their factorizations.

We use the convention that real Lie groups are denoted by upper case Latin letters, e.g  $A, B, C$ , and their Lie algebras by the corresponding lower case German letter  $\mathfrak{a}, \mathfrak{b}, \mathfrak{c}$ .

Throughout this paper  $G$  will denote an algebraic real reductive group and  $H < G$  is an algebraic subgroup. We form the homogeneous space  $Z = G/H$  and write  $z_0 = H$  for the standard base point.

Furthermore, unless otherwise mentioned we assume that  $H$  is reductive in  $G$ , that is, the adjoint representation of  $H$  on  $\mathfrak{g}$  is completely reducible. In this case we say that  $G/H$  is of *reductive type*.

Let us fix a maximal compact subgroup  $K < G$ . It is no loss of generality to request that  $H \cap K$  is maximal compact in  $H$ . Attached to the choice of  $K$  is the infinitesimal Cartan decomposition  $\mathfrak{g} = \mathfrak{k} + \mathfrak{s}$  where  $\mathfrak{s} = \mathfrak{k}^\perp$  is the orthogonal complement with respect to a non-degenerate invariant bilinear form  $\kappa$  on  $\mathfrak{g}$  which is positive definite on  $\mathfrak{s}$  (if  $\mathfrak{g}$  is semi-simple, then we can take for  $\kappa$  the Cartan-Killing form). Further we set  $\mathfrak{q} := \mathfrak{h}^\perp$ .

**2.1. Mostow decomposition.** We recall Mostow's polar decomposition:

$$(2.1) \quad K \times_{H \cap K} \mathfrak{q} \cap \mathfrak{s} \rightarrow Z, \quad [k, X] \mapsto k \exp(X) \cdot z_0$$

which is a homeomorphism. With that we define

$$\|k \exp(X) \cdot z_0\|_Z = \|X\| := \kappa(X, X)^{\frac{1}{2}}$$

for  $k \in K$  and  $X \in \mathfrak{q} \cap \mathfrak{s}$ .

**2.2. Geometric balls.** The problem of lattice counting in  $Z$  leads to a question of exhibiting natural exhausting families of compact subsets. We use balls which are intrinsically defined by the geometry of  $Z$ .

We define the *intrinsic ball* of radius  $R > 0$  on  $Z$  by

$$B_R := \{z \in Z \mid \|z\|_Z < R\}.$$

Write  $B_R^G$  for the intrinsic ball of  $Z = G$ , that is, if  $g = k \exp(X)$  with  $k \in K$  and  $X \in \mathfrak{s}$ , then we put  $\|g\|_G = \|X\|$  and define  $B_R^G$  accordingly.

Our first interest is the growth of the volume  $|B_R|$  for  $R \rightarrow \infty$ . We have the following upper bound.

**Lemma 2.1.** *There exists a constant  $c > 0$  such that:*

$$|B_{R+r}| \leq e^{cr} |B_R|$$

for all  $R \geq 1, r \geq 0$ .

*Proof.* Recall the integral formula

$$(2.2) \quad \int_Z f(z) dz = \int_K \int_{\mathfrak{q} \cap \mathfrak{p}} f(k \exp(X).z_0) \delta(X) dX dk,$$

for  $f \in C_c(Z)$ , where  $\delta(Y)$  is the Jacobian at  $(k, Y)$  of the map (2.1). It is independent of  $k$  because  $dz$  is invariant. Then

$$|B_R| = \int_{X \in \mathfrak{q} \cap \mathfrak{s}, \|X\| < R} \delta(X) dX.$$

Hence it suffices to prove that there exists  $c > 0$  such that

$$\int_0^{R+r} \delta(tX) t^{l-1} dt \leq e^{cr} \int_0^R \delta(tX) t^{l-1} dt$$

for all  $X \in \mathfrak{q} \cap \mathfrak{s}$  with  $\|X\| = 1$ . Here  $l = \dim \mathfrak{q} \cap \mathfrak{s}$ . Equivalently, the function

$$R \mapsto e^{-cR} \int_0^R \delta(tX) t^{l-1} dt$$

is decreasing, or by differentiation,

$$\delta(RX) R^{l-1} \leq c \int_0^R \delta(tX) t^{l-1} dt$$

for all  $R$ . The latter inequality is established in [12, Lemma A.3] with  $c$  independent of  $X$ .  $\square$

Further we are interested how the volume behaves under distortion by elements from  $G$ .

**Lemma 2.2.** *For all  $r, R > 0$  one has  $B_r^G B_R \subset B_{R+r}$ .*

To prove the lemma we first record that:

**Lemma 2.3.** *Let  $z = gH \in Z$ . Then  $\|z\|_Z = \inf_{h \in H} \|gh\|_G$ .*

*Proof.* It suffices to prove that  $\|\exp(X)h\|_G \geq \|X\|$  for  $X \in \mathfrak{q} \cap \mathfrak{s}$ ,  $h \in H$ , and by Cartan decomposition of  $H$ , we may assume  $h = \exp(T)$  with  $T \in \mathfrak{h} \cap \mathfrak{s}$ . Thus we have reduced to the statement that

$$\|\exp(X)\exp(T)\|_G \geq \|\exp(X)\|_G$$

for  $X \perp T$  in  $\mathfrak{s}$ . This follows from the fact that the sectional curvatures of  $K \backslash G$  are  $\leq 0$ .  $\square$

In particular, it follows that

$$(2.3) \quad \|gz\|_Z \leq \|z\|_Z + \|g\|_G \quad (z \in Z, g \in G)$$

and Lemma 2.2 follows.

**2.3. Factorization.** By a (reductive) factorization of  $Z = G/H$  we understand a homogeneous space  $Z^* = G/H^*$  with  $H^*$  an algebraic subgroup of  $G$  such that

- $H^*$  is reductive.
- $H \subset H^*$ .

A factorization is called *compact* if  $Z^*$  is compact, and *co-compact* if the fiber space  $\mathcal{F} := H^*/H$  is compact. It is called *proper* if  $\dim H < \dim H^* < \dim G$ .

Let  $\mathcal{F} \rightarrow Z \rightarrow Z^*$  be a factorization of  $Z$ . We write  $B_R^*$  and  $\mathcal{B}_R^{\mathcal{F}}$  for the intrinsic balls in  $Z^*$  and  $\mathcal{F}$ , respectively.

**Lemma 2.4.** *We have  $B_R^* = B_R H^*/H^*$  and  $\mathcal{B}_R^{\mathcal{F}} = B_R \cap \mathcal{F}$ .*

*Proof.* Follows from Lemma 2.3.  $\square$

For a compactly supported bounded measurable function  $\phi$  on  $Z$  we define the fiberwise integral

$$\phi^{\mathcal{F}}(gH^*) := \int_{H^*/H} \phi(gh^*) d(h^*H)$$

and recall the integration formula

$$(2.4) \quad \int_Z \phi(gH) d(gH) = \int_{Z^*} \phi^{\mathcal{F}}(gH^*) d(gH^*)$$

under appropriate normalization of measures. Consider the characteristic function  $\mathbf{1}_R$  of  $B_R$  and note that its fiber average  $\mathbf{1}_R^{\mathcal{F}}$  is supported in the compact ball  $B_R^*$ . We say that the family of balls  $(B_R)_{R>0}$  *factorizes well to  $Z^*$*  provided for all compact subsets  $Q \subset G$

$$(2.5) \quad \lim_{R \rightarrow \infty} \frac{\sup_{g \in Q} \mathbf{1}_R^{\mathcal{F}}(gH^*)}{|B_R|} = 0.$$



Observe that for all compact subsets  $Q$  there exists an  $R_0 = R_0(Q) > 0$  such that

$$\sup_{g \in Q} \mathbf{1}_R^{\mathcal{F}}(gH^*) \leq |B_{R+R_0}^{\mathcal{F}}|$$

by Lemma 2.2. Thus the balls  $B_R$  factorize well provided

$$(2.6) \quad \lim_{R \rightarrow \infty} \frac{|B_{R+R_0}^{\mathcal{F}}|}{|B_R|} = 0.$$

for all  $R_0 > 0$ .

**2.4. Basic factorizations.** There is a special class of factorizations with which we are dealing with in the sequel. From now on we assume that  $\mathfrak{g}$  is semi-simple and write

$$\mathfrak{g} = \mathfrak{g}_1 \oplus \dots \oplus \mathfrak{g}_m$$

for the decomposition into simple ideals. For a reductive sub algebra  $\mathfrak{h} < \mathfrak{g}$  and a subset  $I \subset \{1, \dots, m\}$  we define the reductive subalgebra

$$(2.7) \quad \mathfrak{h}_I := \mathfrak{h} + \bigoplus_{i \in I} \mathfrak{g}_i.$$

We say that the factorization is *basic* provided that  $\mathfrak{h}^* = \mathfrak{h}_I$  for some  $I$ . Finally we call a factorization *weakly basic* if it is built from a sequence of basic and co-compact factorizations. To be more explicit:

$$\mathfrak{h}^* = \mathfrak{h}_k \supset \dots \supset \mathfrak{h}_0 = \mathfrak{h}$$

such that for each  $i$  we have  $\mathfrak{h}_i = (\mathfrak{h}_{i-1})_I$  for some  $I$  or  $\mathfrak{h}_i/\mathfrak{h}_{i-1}$  compact.

### 3. WAVEFRONT REAL SPHERICAL SPACES

We assume that  $Z$  is real spherical, i.e. a minimal parabolic subgroup  $P < G$  has an open orbit on  $Z$ . It is no loss of generality to assume that  $PH \subset G$  is open, or equivalently that  $\mathfrak{g} = \mathfrak{h} + \mathfrak{p}$ .

If  $L$  is a real algebraic group, then we write  $L_n$  for the normal subgroup of  $L$  which is generated by all unipotent element. In case  $L$  is reductive we observe that  $\mathfrak{l}_n$  is the sum of all non-compact simple ideals of  $\mathfrak{l}$ .

According to [19] there is a unique parabolic subgroup  $Q \supset P$  with the following two properties:

- $QH = PH$ .
- There is a Levi decomposition  $Q = LU$  with  $L_n \subset Q \cap H \subset L$ .

Following [19] we call  $Q$  a  $Z$ -adapted parabolic subgroup.

Having fixed  $L$  we let  $L = K_L A_L N_L$  be an Iwasawa decomposition of  $L$ . We choose an Iwasawa decomposition  $G = KAN$  which inflates the one of  $L$ , i.e.  $K_L < K, A_L = A$  and  $N_L < N$ . Further we may assume that  $N$  is the unipotent radical of the minimal parabolic  $P$ .

Set  $A_H := A \cap H$  and put  $A_Z = A/A_H$ . We recall that  $\dim A_Z$  is an invariant of the real spherical space, called the real rank (see [19]).

In [17], Section 6, we defined the notion of *wavefront* for a real spherical space, which we quickly recall. Attached to  $Z$  is a geometric invariant, the so-called compression cone which is a closed and convex subcone  $\mathfrak{a}_Z^-$  of  $\mathfrak{a}_Z$ . If one denotes by  $\mathfrak{a}^- \subset \mathfrak{a}$  the closure of the negative Weyl-chamber, then  $Z$  being wavefront means that

$$A^- A_H / A_H = A_Z^-.$$

Let us mention that many real spherical spaces are wavefront; for example all symmetric spaces and all Gross-Prasad type spaces  $Z = G \times H/H$  have this property. We recall from [17] the polar decomposition for real spherical spaces

$$(3.1) \quad Z = \Omega A_Z^- F \cdot z_0$$

where

- $\Omega$  is a compact set of the type  $F'K$  with  $F' \subset G$  a finite set.
- $F \subset G$  is a finite set with the property that  $F \cdot z_0 = T \cdot z_0 \cap Z$  where  $T = \exp(i\mathfrak{a})$  and the intersection is taken in  $Z_{\mathbb{C}} = G_{\mathbb{C}}/H_{\mathbb{C}}$ .

**Remark 3.1.** With regard to lattice counting one needs that  $Z = G/H$  carries an invariant measure. If we assume in addition that  $N_G(H) = H$ , then it follows from [16] that  $H$  is reductive.

**3.1. Volume growth.** Define  $\rho_Q \in \mathfrak{a}^*$  by  $\rho_Q(X) = \frac{1}{2} \text{tr}(\text{ad}_{\mathfrak{u}} X)$ ,  $X \in \mathfrak{a}$ . It follows from the unimodularity of  $Z$  and the local structure theorem that  $\rho_Q|_{\mathfrak{a}_H} = 0$ , i.e.  $\rho_Q \in \mathfrak{a}_Z^* = \mathfrak{a}_H^\perp$ .

**Lemma 3.2.** *Let  $Z = G/H$  be a wavefront real spherical space. Then*

$$(3.2) \quad |B_R| \asymp \sup_{\substack{X \in \mathfrak{a} \\ \|X\| \leq R}} e^{2\rho_Q(X)} = \sup_{\substack{X \in \mathfrak{a}_Z^- \\ \|X\| \leq R}} e^{-2\rho_Q(X)}.$$

*Proof.* First note that the equality in (3.2) is immediate from the wavefront assumption.

Let us first show the lower bound, i.e. there exists a  $C > 0$  such that for all  $R > 0$  one has

$$|B_R| \geq C \sup_{\substack{X \in \mathfrak{a} \\ \|X\| \leq R}} e^{2\rho_Q(X)}.$$

For that we recall the volume bound from [18], Prop. 4.2: for all compact subsets  $B \subset G$  with non-empty interior there exists a constant  $C > 0$  such that  $\text{vol}_Z(Ba \cdot z_0) \geq Ca^{2\rho_Q}$  for all  $a \in A_Z^-$ . Together with the polar decomposition (3.1) this gives us the lower bound.

As for the upper bound let

$$\mathfrak{a}_R^- := \{X \in \mathfrak{a}^- \mid \|X\| \leq R\}.$$

Observe that  $B_R \subset B'_R := KA_R^-K \cdot z_0$ . In the sequel it is convenient to realize  $A_Z$  as a subgroup of  $A$  (and not as quotient): we identify  $A_Z$  with  $A_H^\perp \subset A$ . The upper bound will follow if we can show that

$$|B'_R| \leq C \sup_{\substack{X \in \mathfrak{a} \\ \|X\| \leq R}} e^{2\rho_Q(X)} \quad (R > 0).$$

for some constant  $C > 0$ . This in turn will follow from the argument for the upper bound in the proof of Prop. 4.2 in [18]: in this proof we considered for  $a \in A_Z^-$  the map

$$\Phi_a : K \times \Omega_A \times \Xi \rightarrow G, \quad (k, b, X) \mapsto kb \exp(\text{Ad}(a)X)$$

where  $\Omega_A \subset A$  is a compact neighborhood of  $\mathbf{1}$  and  $\Xi \subset \mathfrak{h}$  is a compact neighborhood of  $0$ . It was shown that the Jacobian of  $\Phi_a$ , that is  $\sqrt{\det(d\Phi_a d\Phi_a^t)}$ , is bounded by  $Ca^{-2\rho_Q}$ . Now this bounds holds as well for the right  $K$ -distorted map

$$\Psi_a : K \times \Omega_A \times K \times \Xi \rightarrow G, \quad (k, b, k', X) \mapsto kb \exp(\text{Ad}(ak')X).$$

The reason for that comes from an inspection of the proof; all what is needed is the following fact: let  $d := \dim \mathfrak{h}$  and consider the action of  $\text{Ad}(a)$  on  $V = \bigwedge^d \mathfrak{g}$ . Then for  $a \in A^-$  we have

$$a^{-2\rho} \geq \sup_{\substack{v \in V, \\ \|v\|=1}} \langle \text{Ad}(a)v, v \rangle.$$

We deduce an upper bound

$$(3.3) \quad \text{vol}_Z(K\Omega_A a K \cdot z_0) \leq Ca^{-2\rho}.$$

We need to improve that bound from  $\rho$  to  $\rho_Q$  on the right hand side of (3.3). For that let  $W_L$  be the Weyl group of the reductive pair  $(\mathfrak{l}, \mathfrak{a})$ . Note that  $\rho_Q = \frac{1}{|W_L|} \sum_{w \in W_L} w \cdot \rho$ . Further, the local structure theorem implies that  $L_n \subset H$  and hence  $W_L$  can be realized as a subgroup of

$W_{H \cap K} := N_{H \cap K}(\mathfrak{a})/Z_{H \cap K}(\mathfrak{a})$ . We choose  $\Omega_A$  to be invariant under  $N_{H \cap K}(\mathfrak{a})$  and observe that  $a \in A_Z$  is fixed under  $W_{H \cap K}$ . Thus using the  $N_{H \cap K}(\mathfrak{a})$ -symmetry in the  $a$ -variable we refine (3.3) to

$$\mathrm{vol}_Z(K\Omega_A a K \cdot z_0) \leq C a^{-2\rho_Q}.$$

The desired bound then follows.  $\square$

**Corollary 3.3.** *Let  $Z = G/H$  be a wavefront real spherical space of reductive type. Let  $Z \rightarrow Z^*$  be a basic factorization such that  $Z^*$  is not compact. Then the geometric balls  $B_R$  factorize well to  $Z^*$ .*

*Proof.* As  $Z \rightarrow Z^*$  is basic we may assume (ignoring connected components) that  $H^* = G_I H$  for some  $I$ . Note that  $\mathcal{F} = H^*/H \simeq G_I/G_I \cap H$  is real spherical.

Let  $Q$  be the  $Z$ -adapted parabolic subgroup attached to  $P$ . Let  $P_I = P \cap G_I$  and  $G_I \supset Q_I \supset P_I$  be the  $\mathcal{F}$ -adapted parabolic above  $P_I$  and note that  $Q_I = Q \cap G_I$ . With Lemma 3.2 we then get

$$|B_R^{\mathcal{F}}| \asymp \sup_{\substack{X \in \mathfrak{a}_I \\ \|X\| \leq R}} e^{2\rho_{Q_I}(X)},$$

which we are going to compare with (3.2).

Let  $\mathfrak{u}_I$  be the Lie algebra of the unipotent radical of  $Q_I$ . Note that  $\mathfrak{u}_I \subset \mathfrak{u}$  and that this inclusion is strict since  $G/H^*$  is not compact. The corollary now follows from (2.6).  $\square$

**3.2. Property I.** We briefly recall some results from [18].

Let  $(\pi, \mathcal{H}_\pi)$  be a unitary irreducible representation of  $G$ . We denote by  $\mathcal{H}_\pi^\infty$  the  $G$ -Fréchet module of smooth vectors and by  $\mathcal{H}_\pi^{-\infty}$  its strong dual. One calls  $\mathcal{H}_\pi^{-\infty}$  the  $G$ -module of distribution vectors; it is a DNF-space with continuous  $G$ -action.

Let  $\eta \in (\mathcal{H}_\pi^{-\infty})^H$  be an  $H$ -fixed element and  $H_\eta < G$  the stabilizer of  $\eta$ . Note that  $H < H_\eta$  and set  $Z_\eta := G/H_\eta$ . With regard to  $\eta$  and  $v \in \mathcal{H}^\infty$  we form the generalized matrix-coefficient

$$m_{v,\eta}(gH) := \eta(\pi(g^{-1})v) \quad (g \in G)$$

which is a smooth function on  $Z_\eta$ .

We recall the following fact from [18]:

**Proposition 3.4.** *Let  $Z$  be a wavefront real spherical space of reductive type. Then the following assertions hold:*

- (1) *Let  $H < H^* < G$  be a closed subgroup such that  $Z^*$  is unimodular. Then  $Z^*$  is a weakly basic factorization.*
- (2) *Let  $(\pi, \mathcal{H})$  be a unitary irreducible representation of  $G$  and let  $\eta \in (\mathcal{H}_\pi^{-\infty})^H$ . Then:*

- (a)  $Z \rightarrow Z_\eta$  is a weakly basic factorization.
- (b)  $Z_\eta$  is unimodular and there exists  $1 \leq p < \infty$  such that  $m_{v,\eta} \in L^p(Z_\eta)$  for all  $v \in \mathcal{H}_\pi^\infty$ .

The property of  $Z = G/H$  that (2b) is valid for all  $\pi$  and  $\eta$  as above is denoted *Property (I)* in [18]. Assuming this property we define  $p_H(\pi)$  as the smallest index  $\geq 1$  such that all  $K$ -finite generalized matrix coefficients  $m_{v,\eta}$  with  $\eta \in (\mathcal{H}_\pi^{-\infty})^H$  belong to  $L^p(Z_\eta)$  for any  $p > p_H(\pi)$ . It follows from finite dimensionality of  $(\mathcal{H}_\pi^{-\infty})^H$  (see [22]) that  $p_H(\pi) < \infty$ . We say that  $\pi$  is  $H$ -tempered if  $p_H(\pi) = 2$ .

The representation  $\pi$  is said to be  $H$ -distinguished if  $(\mathcal{H}_\pi^{-\infty})^H \neq \{0\}$ . Note that if  $\pi$  is not  $H$ -distinguished then  $p_H(\pi) = 1$ .

**Remark 3.5.** Let  $Z = G/H$  be a wavefront reductive homogeneous space. Then, with (1) above and a little bit of effort, one can show that every factorization  $Z \rightarrow Z^*$  (see Section 1.3.1) is of the type

$$Z \rightarrow Z_c^* \rightarrow Z^*$$

with  $Z \rightarrow Z_c^*$  co-compact and  $Z_c^* \rightarrow Z^*$  basic. Hence if we neglect compact symmetries of  $Z$ , which is natural in the context of lattice counting, then it is natural to assume that every factorization is basic.

#### 4. LATTICE POINT COUNTING: SETUP

Let  $G/H$  be a real algebraic homogeneous space. We further assume that we are given a lattice (a discrete subgroup with finite covolume)  $\Gamma \subset G$ , such that  $\Gamma_H := \Gamma \cap H$  is a lattice in  $H$ . We normalize Haar measures on  $G$  and  $H$  such that:

- $\text{vol}(G/\Gamma) = 1$ .
- $\text{vol}(H/\Gamma_H) = 1$ .

Our concern is with the double fibration

$$\begin{array}{ccc} & G/\Gamma_H & \\ \swarrow & & \searrow \\ Z := G/H & & Y := G/\Gamma \end{array}$$

Fibre-wise integration yields transfer maps from functions on  $Z$  to functions on  $Y$  and vice versa. In more precision,

$$(4.1) \quad L^\infty(Y) \rightarrow L^\infty(Z), \quad \phi \mapsto \phi^H; \quad \phi^H(gH) := \int_{H/\Gamma_H} \phi(gh\Gamma) d(h\Gamma_H)$$

and we record that this map is contractive, i.e

$$(4.2) \quad \|\phi^H\|_\infty \leq \|\phi\|_\infty \quad (\phi \in L^\infty(Y)).$$

Likewise we have

$$(4.3) \quad L^1(Z) \rightarrow L^1(Y), \quad f \mapsto f^\Gamma; \quad f^\Gamma(g\Gamma) := \sum_{\gamma \in \Gamma/\Gamma_H} f(g\gamma H),$$

which is contractive, i.e

$$(4.4) \quad \|f^\Gamma\|_1 \leq \|f\|_1 \quad (f \in L^1(Z)).$$

Unfolding with respect to the double fibration yields, in view of our normalization of measures, the following adjointness relation:

$$(4.5) \quad \langle f^\Gamma, \phi \rangle_{L^2(Y)} = \langle f, \phi^H \rangle_{L^2(Z)}$$

for all  $\phi \in L^\infty(Y)$  and  $f \in L^1(Z)$ . Let us note that (4.5) applied to  $|f|$  and  $\phi = \mathbf{1}_Y$  readily yields (4.4).

We write  $\mathbf{1}_R \in L^1(Z)$  for the characteristic function of  $B_R$  and deduce from the definitions and (4.5):

- $\mathbf{1}_R^\Gamma(e\Gamma) = N_R(\Gamma, Z) := \#\{\gamma \in \Gamma/\Gamma_H \mid \gamma \cdot z_0 \in B_R\}$ .
- $\|\mathbf{1}_R^\Gamma\|_{L^1(G/\Gamma)} = |B_R|$ .

**4.1. Weak asymptotics.** In the above setup,  $G/H$  need not be of reductive type, but we shall assume this again from now on. For spaces with property (I) and  $Y$  compact we prove analytically in the following section that

$$(MT) \quad N_R(\Gamma, Z) \sim |B_R| \quad (R \rightarrow \infty).$$

For that we will use the following result of [20]:

**Theorem 4.1.** *Let  $Z = G/H$  be of reductive type. The smooth vectors for the regular representation of  $G$  on  $L^p(Z)$  vanish at infinity, for all  $1 \leq p < \infty$ .*

With notation from (4.3) we set

$$F_R^\Gamma := \frac{1}{|B_R|} \mathbf{1}_R^\Gamma.$$

We shall concentrate on verifying the following limit of weak type:

$$(wMT) \quad \langle F_R^\Gamma, \phi \rangle_{L^2(Y)} \rightarrow \int_Y \bar{\phi} d\mu_Y \quad (R \rightarrow \infty), \quad (\forall \phi \in C_0(Y)).$$

Here  $C_0$  indicates functions vanishing at infinity.

**Lemma 4.2.** (wMT)  $\Rightarrow$  (MT).

*Proof.* As in [10] Lemma 2.3 this is deduced from Lemma 2.1 and Lemma 2.2.  $\square$

## 5. MAIN TERM COUNTING

In this section we will establish main term counting under the mandate of property (I) and  $Y$  being compact. Let us call a family of balls  $(B_R)_{R>0}$  *well factorizable* if it factorizes well to all proper factorizations of type  $Z \rightarrow Z_\eta$ .

## 5.1. Main theorem on counting.

**Theorem 5.1.** *Let  $G$  be semi-simple and  $H$  a closed reductive subgroup. Suppose that  $Y$  is compact and  $Z$  admits (I). If  $(B_R)_{R>0}$  is well factorizable, then (wMT) and (MT) hold.*

**Remark 5.2.** In case  $Z = G/H$  is real spherical and wavefront, then  $Z$  has (I) by Proposition 3.4. If we assume in addition that  $G$  has no compact factors and that all proper factorizations are basic, then the family of geometric balls is well factorizable by Corollary 3.3. In particular, Theorem A of the introduction then follows from the above.

The proof is based on the following proposition. For a function space  $\mathcal{F}(Y)$  consisting of integrable functions on  $Y$  we denote by  $\mathcal{F}(Y)_{\text{van}}$  the subspace of functions with vanishing integral over  $Y$ .

**Proposition 5.3.** *Let  $Z = G/H$  be of reductive type. Assume that there exists a dense subspace  $\mathcal{A}(Y) \subset C_b(Y)_{\text{van}}^K$  such that*

$$(5.1) \quad \phi^H \in C_0(Z) \quad \text{for all } \phi \in \mathcal{A}(Y).$$

*Then (wMT) holds true.*

*Proof.* We will establish (wMT) for  $\phi \in C_b(Y)$ . As

$$C_b(Y) = C_b(Y)_{\text{van}} \oplus \mathbb{C}\mathbf{1}_Y,$$

and (wMT) is trivial for  $\phi$  a constant, it suffices to establish

$$(5.2) \quad \langle F_R^\Gamma, \phi \rangle_{L^2(Y)} \rightarrow 0 \quad (\phi \in C_b(Y)_{\text{van}}).$$

We will show (5.2) is valid for  $\phi \in \mathcal{A}(Y)$ . By density, as  $F_R^\Gamma$  is  $K$ -invariant and belongs to  $L^1(Y)$ , this will finish the proof.

Let  $\phi \in \mathcal{A}(Y)$  and let  $\epsilon > 0$ . By the unfolding identity (4.5) we have

$$(5.3) \quad \langle F_R^\Gamma, \phi \rangle_{L^2(Y)} = \frac{1}{|B_R|} \langle \mathbf{1}_R, \phi^H \rangle_{L^2(Z)}.$$

Using (5.1) we choose  $K_\epsilon \subset Z$  compact such that  $|\phi^H(z)| < \epsilon$  outside of  $K_\epsilon$ . Then

$$\frac{1}{|B_R|} \langle \mathbf{1}_R, \phi^H \rangle_{L^2(Z)} = \int_{K_\epsilon} + \int_{Z-K_\epsilon} \frac{\mathbf{1}_R(z)}{|B_R|} \phi^H(z) d\mu_Z(z).$$

By (4.2), the first term is bounded by  $\frac{|K_\epsilon| \|\phi\|_\infty}{|B_R|}$ , which is  $\leq \epsilon$  for  $R$  sufficiently large. As the second term is bounded by  $\epsilon$  for all  $R$ , we obtain (5.2). Hence (wMT) holds.  $\square$

**Remark 5.4.** It is possible to replace (5.1) by a weaker requirement: Suppose that an algebraic sum

$$(5.4) \quad \mathcal{A}(Y) = \sum_{j \in J} \mathcal{A}(Y)_j$$

is given together with a factorization  $Z_j^* = G/H_j^*$  for each  $j \in J$ . Suppose that the balls  $B_R$  all factorize well to  $Z_j^*$ ,  $j \in J$ . Suppose further that  $\phi^H$  factorizes to a function

$$(5.5) \quad \phi^{H_j^*} \in C_0(Z_j^*)$$

for all  $\phi \in \mathcal{A}(Y)_j$  and all  $j \in J$ . Then the conclusion in Proposition 5.3 is still valid. In fact, using (2.4) the last part of the proof modifies to:

$$\begin{aligned} \frac{1}{|B_R|} \langle \mathbf{1}_R, \phi^H \rangle_{L^2(Z)} &= \frac{1}{|B_R|} \langle \mathbf{1}_R^{\mathcal{F}}, \phi^{H_j^*} \rangle_{L^2(Z_j^*)} = \\ &= \int_{K_\epsilon^*} + \int_{Z_j^* - K_\epsilon^*} \frac{\mathbf{1}_R^{\mathcal{F}}(z)}{|B_R|} \phi^{H_j^*}(z) d\mu_{Z_j^*}(z) \end{aligned}$$

for  $\phi \in \mathcal{A}(Y)_j$ . As  $\|\mathbf{1}_R^{\mathcal{F}}\|_{L^1(Z_j^*)} = |B_R|$ , the second term is bounded by  $\epsilon$  for all  $R$ . As the balls factorize well to  $Z_j^*$  we get the first term as small as we wish with (2.5).

**5.2. The space  $\mathcal{A}(Y)$ .** We now construct a specific subspace  $\mathcal{A}(Y) \subset C_b(Y)_{\text{van}}^K$  and verify condition (5.5).

Denote by  $\widehat{G}_s \subset \widehat{G}$  the  $K$ -spherical unitary dual.

As  $Y$  is compact, the abstract Plancherel-theorem implies:

$$L^2(G/\Gamma)^K \simeq \bigoplus_{\pi \in \widehat{G}_s} (\mathcal{H}_\pi^{-\infty})^\Gamma.$$

If we denote the Fourier transform by  $f \mapsto f^\wedge$  then the corresponding inversion formula is given by

$$(5.6) \quad f = \sum_{\pi} m_{v_\pi, f^\wedge(\pi)}$$

with  $v_\pi \in \mathcal{H}_\pi$  normalized  $K$ -fixed and  $f^\wedge(\pi) \in (\mathcal{H}_\pi^{-\infty})^\Gamma$ . The matrix coefficients for  $Y$  are defined as in (1.2), and the sum in (5.6) is required to include multiplicities.



Note that  $L^2(Y) = L^2(Y)_{\text{van}} \oplus \mathbb{C} \cdot \mathbf{1}_Y$ . We define  $\mathcal{A}(Y) \subset L^2(Y)_{\text{van}}^K$  to be the dense subspace of functions with finite Fourier support, that is,

$$\mathcal{A}(Y) = \text{span}\{m_{v,\eta} \mid \pi \in \widehat{G}_s \text{ non-trivial}, v \in \mathcal{H}_\pi^K, \eta \in (\mathcal{H}_\pi^{-\infty})^\Gamma\}.$$

Then  $\mathcal{A}(Y) \subset L^2(Y)_{\text{van}}^{K,\infty}$  is dense and since  $C^\infty(Y)$  and  $L^2(Y)^\infty$  are topologically isomorphic, it follows that  $\mathcal{A}(Y)$  is dense in  $C(Y)_{\text{van}}^K$  as required.

The following lemma together with Remark 5.4 immediately implies Theorem 5.1.

**Lemma 5.5.** *Assume that  $Y$  is compact and  $Z$  has (I), and define  $\mathcal{A}(Y)$  as above. Then there exists a decomposition of  $\mathcal{A}(Y)$  satisfying (5.4)-(5.5).*

*Proof.* Let  $J$  denote the set of all factorizations  $Z^* \rightarrow Z$ , including also  $Z^* = Z$  which we give the index  $j_0 \in J$ . For  $j \in J$  we define  $\mathcal{A}(Y)_j \subset \mathcal{A}(Y)$  accordingly to be spanned by the matrix coefficients for which  $H_\eta = H_j^*$ . Then (5.4) holds.

We consider first the subspace  $\mathcal{A}(Y)_{j_0}$ . The map  $\phi \mapsto \phi^H$  from (4.1) corresponds on the spectral side to a map  $(\mathcal{H}_\pi^{-\infty})^\Gamma \rightarrow (\mathcal{H}_\pi^{-\infty})^H$ , which can be constructed as follows.

As  $H/\Gamma_H$  is compact, we can define for each  $\pi \in \widehat{G}_s$

$$(5.7) \quad \Lambda_\pi : (\mathcal{H}_\pi^{-\infty})^\Gamma \rightarrow (\mathcal{H}_\pi^{-\infty})^H, \quad \Lambda_\pi(\eta) = \int_{H/\Gamma_H} \eta \circ \pi(h^{-1}) d(h\Gamma_H)$$

by  $\mathcal{H}_\pi^{-\infty}$ -valued integration: the defining integral is understood as integration over a compact fundamental domain  $F \subset H$  with respect to the Haar measure on  $H$ ; as the integrand is continuous and  $\mathcal{H}_\pi^{-\infty}$  is a complete locally convex space, the integral converges in  $\mathcal{H}_\pi^{-\infty}$ . It follows from (5.7) that  $(m_{v,\eta})^H = m_{v,\Lambda_\pi(\eta)}$  for all  $v \in \mathcal{H}_\pi^\infty$  and  $\eta \in (\mathcal{H}_\pi^{-\infty})^\Gamma$ .

Let  $\phi \in \mathcal{A}(Y)_{j_0}$ , then it follows from (5.6) that

$$(5.8) \quad \phi^H = \sum_{\pi \neq \mathbf{1}} m_{v_\pi, \Lambda_\pi(\phi^\wedge(\pi))}.$$

Note that  $H_\eta = H$  for each distribution vector  $\eta = \Lambda_\pi(\phi^\wedge(\pi))$  in this sum, by the definition of  $\mathcal{A}(Y)_{j_0}$ . As  $Z$  has property (I) the summand  $m_{v_\pi, \Lambda_\pi(\phi^\wedge(\pi))}$  is contained in  $L^p(G/H)$  for  $p > p_H(\pi)$ , and by [18], Lemma 7.2, this containment is then valid for all  $K$ -finite generalized matrix coefficients  $m_{v, \Lambda_\pi(\phi^\wedge(\pi))}$  of  $\pi$ . Thus  $m_{v_\pi, \Lambda_\pi(\phi^\wedge(\pi))}$  generates a Harish-Chandra module inside  $L^p(G/H)$ . As  $m_{v_\pi, \Lambda_\pi(\phi^\wedge(\pi))}$  is  $K$ -finite, we conclude that it is a smooth vector. Hence  $\phi^H \in L^p(G/H)^\infty$ , and in view of Theorem 4.1 we obtain (5.1).

The proof of (5.5) for  $\phi \in \mathcal{A}(Y)_j$  for general  $j \in J$  is obtained by the same reasoning, where one replaces  $H$  by  $H_j^*$  in (5.7) and (5.8).  $\square$

This concludes the proof of Theorem 5.1.

## 6. $L^p$ -BOUNDS FOR GENERALIZED MATRIX COEFFICIENTS

From here on we assume that  $Z = G/H$  is wavefront and real spherical. Recall that we assumed that  $G$  is semi-simple and that we wrote  $\mathfrak{g} = \mathfrak{g}_1 \oplus \dots \oplus \mathfrak{g}_m$  for the decomposition of  $\mathfrak{g}$  into simple factors. It is no big loss of generality to assume that  $G = G_1 \times \dots \times G_m$  splits accordingly. We will assume that from now on.

Further we request that the lattice  $\Gamma < G$  is irreducible, that is, the projection of  $\Gamma$  to any normal subgroup  $J \subsetneq G$  is dense in  $J$ .

Let  $\pi$  be an irreducible unitary representation of  $G$ . Then  $\pi = \pi_1 \otimes \dots \otimes \pi_m$  with  $\pi_j$  and irreducible unitary representation of  $G_j$ . We start with a simple observation.

**Lemma 6.1.** *Let  $(\pi, \mathcal{H})$  be an irreducible unitary representation of  $G$  and  $0 \neq \nu \in (\mathcal{H}^{-\infty})^\Gamma$ . Then, if one constituent  $\pi_j$  of  $\pi$  is trivial, then  $\pi$  is trivial.*

*Proof.* The element  $\nu$  gives rise to a  $G$ -equivariant injection

$$(6.1) \quad \mathcal{H}^\infty \hookrightarrow C^\infty(Y), \quad v \mapsto (g\Gamma \mapsto \nu(\pi(g^{-1})v)).$$

Say  $\pi_j$  is trivial and let  $J := \prod_{i \neq j}^m G_i$ . Let  $\Gamma_J$  be the projection of  $\Gamma$  to  $J$ . Then (6.1) gives rise to a  $J$ -equivariant injection  $\mathcal{H}^\infty \hookrightarrow C^\infty(J/\Gamma_J)$ . As  $\Gamma_J$  is dense in  $J$ , the assertion follows.  $\square$

As we discussed earlier in Remark 3.5 it is no big loss of generality to request that all factorizations are basic. We assume this from now on.

Further we request from now that the cycle  $H/\Gamma_H \subset Y$  is compact. This technical condition ensures that the vector valued average map (5.7) converges.

**Lemma 6.2.** *Let  $(\pi, \mathcal{H})$  be a non-trivial irreducible unitary representation of  $G$ . Let  $\nu \in (\mathcal{H}_\pi^{-\infty})^\Gamma$  such that  $\eta := \Lambda_\pi(\nu) \in (\mathcal{H}_\pi^{-\infty})^H$  is non-zero. Then  $H_\eta/H$  is finite.*

*Proof.* As all factorizations are basic we find  $\mathfrak{h}_\eta = \mathfrak{h}_I$ , and as  $\pi$  is irreducible it infinitesimally embeds into  $C^\infty(G/H_\eta)$ . It follows that  $\pi_i$  is trivial for  $i \in I$ . Hence Lemma 6.1 implies  $I = \emptyset$  and  $\mathfrak{h}_\eta = \mathfrak{h}$ .  $\square$

In the sequel we use the Plancherel theorem (see [14])

$$L^2(G/\Gamma)^K \simeq \int_{\widehat{G}_s}^{\oplus} \mathcal{V}_{\pi,\Gamma} d\mu(\pi),$$

where  $\mathcal{V}_{\pi,\Gamma} \subset (\mathcal{H}_{\pi}^{-\infty})^{\Gamma}$  is a finite dimensional subspace and of constant dimension on each connected component in the continuous spectrum (parametrization by Eisenstein series), and where the Plancherel measure  $\mu$  has support

$$\widehat{G}_{\Gamma,s} := \text{supp}(\mu) \subset \widehat{G}_s.$$

Given an irreducible lattice  $\Gamma \subset G$  we define

$$p_H(\Gamma) := \sup\{p_H(\pi) : \pi \in \widehat{G}_{\Gamma,s}\}$$

and record the following.

**Lemma 6.3.**  $p_H(\Gamma) < \infty$ .

*Proof.* For a unitary representation  $(\pi, \mathcal{H})$  and vectors  $v, w \in \mathcal{H}$  we form the matrix coefficient  $\pi_{v,w}(g) := \langle \pi(g)v, w \rangle$ . We first claim that there exists a  $p < \infty$  (in general depending on  $\Gamma$ ) such that for all non-trivial  $\pi \in \widehat{G}_{\Gamma,s}$  one has  $\pi_{v,w} \in L^p(G)$  for all  $K$ -finite vectors  $v, w$ . In case  $G$  has property (T) this follows (independently of  $\Gamma$ ) from [7]. The remaining cases are  $\text{SO}_{\epsilon}(n, 1)$  and  $\text{SU}(n, 1)$  (up to covering) of real rank one, and they are treated in [6].

The claim can be interpreted geometrically via the leading exponent  $\Lambda_V \in \mathfrak{a}^*$  which is attached to the Harish-Chandra module of  $\mathcal{H}$  (see [18], Section 6). The lemma now follows from Prop. 4.2 and Thm. 6.3 in [18] (see the proof of Thm. 7.6 in [18] how these two facts combine to result in integrability).  $\square$

Let  $1 \leq p < \infty$ . Let us say that a subset  $\Lambda \subset \widehat{G}_s$  is  $L^p$ -bounded provided that  $m_{v,\eta} \in L^p(Z_{\eta})$  for all  $\pi \in \Lambda$  and  $v \in \mathcal{H}_{\pi}^{\infty}$ ,  $\eta \in (\mathcal{H}_{\pi}^{-\infty})^H$ . By definition we thus have that  $\widehat{G}_{\Gamma,s}$  is  $L^p$ -bounded for  $p > p_H(\Gamma)$ .

In this section we work under the following:

**Hypothesis A:** For every  $1 \leq p < \infty$  and every  $L^p$ -bounded subset  $\Lambda \subset \widehat{G}_s$  there exists a compact subset  $\Omega \subset G$  and constants  $c, C > 0$  such that the following assertions hold for all  $\pi \in \Lambda$ ,  $\eta \in (\mathcal{H}_{\pi}^{-\infty})^H$  and  $v \in \mathcal{H}_{\pi}^K$ :

$$(A1) \quad \|m_{v,\eta}\|_{L^p(Z_{\eta})} \leq C \|m_{v,\eta}\|_{\infty},$$

$$(A2) \quad \|m_{v,\eta}\|_{\infty} \leq c \|m_{v,\eta}\|_{\infty, \Omega_{\eta}}$$

where  $\Omega_{\eta} = \Omega H_{\eta} / H_{\eta}$ .

In the sequel we are only interested in the following choice of subset  $\Lambda \subset \widehat{G}_s$ , namely

$$(6.2) \quad \Lambda := \{\pi \in \widehat{G}_{\Gamma,s} \mid \Lambda_\pi(\nu) \neq 0 \text{ for some } \nu \in \mathcal{V}_{\pi,\Gamma}\}.$$

An immediate consequence of Hypothesis A is:

**Lemma 6.4.** *Assume that  $p > p_H(\Gamma)$ . Then there is a  $C > 0$  such that for all  $\pi \in \widehat{G}_{\Gamma,s}$ ,  $v \in \mathcal{H}_\pi^K$ ,  $\nu \in (\mathcal{H}_\pi^{-\infty})^\Gamma$  and  $\eta := \Lambda_\pi(\nu) \in (\mathcal{H}_\pi^{-\infty})^H$  one has*

$$\|\phi_\pi^H\|_{L^p(Z_\eta)} \leq C \|\phi_\pi\|_\infty$$

where  $\phi_\pi(g\Gamma) := \nu(\pi(g^{-1})v)$ .

*Proof.* Recall from (4.2), that integration is a bounded operator from  $L^\infty(Y) \rightarrow L^\infty(Z)$ . Hence the assertion follows from (A1).  $\square$

Recall the Cartan-Killing form  $\kappa$  on  $\mathfrak{g} = \mathfrak{k} + \mathfrak{s}$  and choose a basis  $X_1, \dots, X_l$  of  $\mathfrak{k}$  and  $X'_1, \dots, X'_s$  of  $\mathfrak{s}$  such that  $\kappa(X_i, X_j) = -\delta_{ij}$  and  $\kappa(X'_i, X'_j) = \delta_{ij}$ . With that data we form the standard Casimir element

$$\mathcal{C} := -\sum_{j=1}^l X_j^2 + \sum_{j=1}^s (X'_j)^2 \in \mathcal{U}(\mathfrak{g}).$$

Set  $\Delta_K := \sum_{j=1}^l X_j^2 \in \mathcal{U}(\mathfrak{k})$  and obtain the commonly used Laplace element

$$(6.3) \quad \Delta = \mathcal{C} + 2\Delta_K \in \mathcal{U}(\mathfrak{g})$$

which acts on  $Y = G/\Gamma$  from the left.

Let  $d \in \mathbb{N}$ . For  $1 \leq p \leq \infty$ , it follows from [2], Section 3, that Sobolev norms on  $L^p(Y)^\infty \subset C^\infty(Y)$  can be defined by

$$\|f\|_{p,2d}^2 = \sum_{j=0}^d \|\Delta^j f\|_p^2.$$

Basic spectral theory allows one to define  $\|\cdot\|_{p,d}$  more generally for any  $d \geq 0$ .

Let us define

$$s := \dim \mathfrak{s} = \dim G/K = \dim \Gamma \backslash G/K$$

and

$$r := \dim \mathfrak{a} = \text{rank}_{\mathbb{R}}(G/K),$$

where  $\mathfrak{a} \subset \mathfrak{s}$  is maximal abelian.

We denote by  $C_b(Y)$  the space of continuous bounded functions on  $Y$  and by  $C_b(Y)_o$  the subspace with vanishing integral.

**Proposition 6.5.** *Assume that*

- (1)  $Z$  is a wavefront real spherical space with all factorizations basic,
- (2)  $G = G_1 \times \dots \times G_m$  with all  $\mathfrak{g}_i$  simple and non-compact.
- (3)  $\Gamma < G$  is irreducible and  $Y_H$  is compact,
- (4) Hypothesis A is valid.

Then the map

$$\text{Av}_H : C_b^\infty(Y)_o^K \rightarrow L^p(Z)^K; \text{Av}_H(\phi) = \phi^H$$

is continuous. More precisely, for all

- (1)  $k > s + 1$  if  $Y$  is compact.
- (2)  $k > \frac{r+1}{2}s + 1$  if  $Y$  is non-compact and  $\Gamma$  is arithmetic

there exists a constant  $C = C(p, k) > 0$  such that

$$\|\phi^H\|_{L^p(Z)} \leq C \|\phi\|_{\infty, k} \quad (\phi \in C_b^\infty(Y)_o^K)$$

*Proof.* For all  $\pi \in \widehat{G}$  the operator  $d\pi(\mathcal{C})$  acts as a scalar  $\lambda_\pi$  and we set

$$|\pi| := |\lambda_\pi| \geq 0.$$

Let  $\phi \in C_b^\infty(Y)_o^K$  and write  $\phi = \phi_d + \phi_c$  for its decomposition in discrete and continuous Plancherel parts. We assume first that  $\phi = \phi_d$ .

In case  $Y$  is compact we have Weyl's law: There is a constant  $c_Y > 0$  such that

$$\sum_{|\pi| \leq R} m(\pi) \sim c_Y R^{s/2} \quad (R \rightarrow \infty).$$

Here  $m(\pi) = \dim \mathcal{V}_{\pi, \Gamma}$ . We conclude that

$$(6.4) \quad \sum_{\pi} m(\pi)(1 + |\pi|)^{-k} < \infty$$

for all  $k > s/2 + 1$ . In case  $Y$  is non-compact, we let  $\widehat{G}_{\mu, d}$  be the discrete support of the Plancherel measure. Then assuming  $\Gamma$  is arithmetic, the upper bound in [15] reads:

$$\sum_{\substack{\pi \in \widehat{G}_{\mu, d} \\ |\pi| \leq R}} m(\pi) \leq c_Y R^{rs/2} \quad (R > 0).$$

For  $k > rs/2 + 1$  we obtain (6.4) as before.

As  $\phi$  is in the discrete spectrum we decompose it as  $\phi = \sum_{\pi} \phi_{\pi}$  and obtain with Hypothesis (A1)

$$\|\phi^H\|_p \leq \sum_{\pi} \|\phi_{\pi}^H\|_p \leq C \sum_{\pi} \|\phi_{\pi}\|_{\infty}.$$

The last sum we estimate as follows:

$$\begin{aligned} \sum_{\pi} \|\phi_{\pi}\|_{\infty} &= \sum_{\pi} (1 + |\pi|)^{-k/2} (1 + |\pi|)^{k/2} \|\phi_{\pi}\|_{\infty} \\ &\leq C \sum_{\pi} (1 + |\pi|)^{-k/2} \|\phi_{\pi}\|_{\infty, k} \end{aligned}$$

with  $C > 0$  a constant depending only on  $k$  (we allow universal positive constants to change from line to line). Applying the Cauchy-Schwartz inequality combined with (6.4) we obtain

$$\|\phi^H\|_p \leq C \left( \sum_{\pi} \|\phi_{\pi}\|_{\infty, k}^2 \right)^{\frac{1}{2}}$$

with  $C > 0$ . With Hypothesis (A2) we get the further improvement:

$$\|\phi^H\|_p \leq C \left( \sum_{\pi} \|\phi_{\pi}\|_{\Omega, \infty, k}^2 \right)^{\frac{1}{2}}$$

To finish the proof we apply the Sobolev lemma on  $K \setminus G$ . Here Sobolev norms are defined by the central operator  $\mathcal{C}$ , whose action agrees with the left action of  $\Delta$ . It follows that  $\|f\|_{\infty, \Omega} \leq C \|f\|_{2, k_1, \Omega}$  with  $k_1 > \frac{s}{2}$  for  $K$ -invariant functions  $f$  on  $G$ . This gives

$$\|\phi^H\|_p \leq C \left( \sum_{\pi} \|\phi_{\pi}\|_{\Omega, 2, k+k_1}^2 \right)^{\frac{1}{2}} = C \|\phi\|_{\Omega, 2, k+k_1} \leq C \|\phi\|_{\infty, k+k_1}$$

which proves the proposition for the discrete spectrum.

If  $\phi = \phi_c$  belongs to the continuous spectrum, where multiplicities are bounded (see [14]), the proof is simpler. Let  $\mu_c$  be the restriction of the Plancherel measure to the continuous spectrum. As this is just Euclidean measure on  $r$ -dimensional space we have

$$(6.5) \quad \int_{\widehat{G}_s} (1 + |\pi|)^{-k} d\mu_c(\pi) < \infty$$

if  $k > r/2$ . We assume for simplicity in what follows that  $m(\pi) = 1$  for almost all  $\pi \in \text{supp } \mu_c$ . As  $\sup_{\pi \in \text{supp } \mu_c} m(\pi) < \infty$  the proof is easily adapted to the general case.

Let

$$\phi = \int_{\widehat{G}_s} \phi_{\pi} d\mu_c(\pi).$$

As  $\|\phi^H\|_{\infty} \leq \|\phi\|_{\infty}$  we conclude with Lemma 6.4, (6.5) and Fubini's theorem that

$$\phi^H = \int_{\widehat{G}_s} \phi_{\pi}^H d\mu_c(\pi)$$

and, by the similar chain of inequalities as in the discrete case

$$\|\phi^H\|_p \leq C\|\phi\|_{\infty, k+k_1}$$

with  $k > \frac{r}{2}$  and  $k_1 > \frac{s}{2}$ . This concludes the proof.  $\square$

## 7. ERROR TERM ESTIMATES

Recall  $\mathbf{1}_R$ , the characteristic function of  $B_R$ . The first error term for the lattice counting problem can be expressed by

$$\text{err}(R, \Gamma) := \sup_{\substack{\phi \in C_b(Y) \\ \|\phi\|_{\infty} \leq 1}} \left| \left\langle \frac{\mathbf{1}_R^{\Gamma}}{|B_R|} - \mathbf{1}_Y, \phi \right\rangle \right| \quad (R > 0),$$

and our goal is to give an upper bound for  $\text{err}(R, \Gamma)$  as a function of  $R$ .

According to the decomposition  $C_b(Y) = C_b(Y)_o \oplus \mathbb{C}\mathbf{1}_Y$  we decompose functions as  $\phi = \phi_o + \phi_1$  and obtain

$$\text{err}(R, \Gamma) = \sup_{\substack{\phi \in C_b(Y) \\ \|\phi\|_{\infty} \leq 1}} \frac{|\langle \mathbf{1}_R^{\Gamma}, \phi_o \rangle|}{|B_R|} = \sup_{\substack{\phi \in C_b(Y) \\ \|\phi\|_{\infty} \leq 1}} \frac{|\langle \mathbf{1}_R, \phi_o^H \rangle|}{|B_R|}.$$

Further, from  $\|\phi_o\|_{\infty} \leq 2\|\phi\|_{\infty}$  we obtain that  $\text{err}(R, \Gamma) \leq 2\text{err}_1(R, \Gamma)$  with

$$\text{err}_1(R, \Gamma) := \sup_{\substack{\phi \in C_b(Y)_o \\ \|\phi\|_{\infty} \leq 1}} \frac{|\langle \mathbf{1}_R^{\Gamma}, \phi \rangle|}{|B_R|} = \sup_{\substack{\phi \in C_b(Y)_o \\ \|\phi\|_{\infty} \leq 1}} \frac{|\langle \mathbf{1}_R, \phi^H \rangle|}{|B_R|}.$$

**7.1. Smooth versus non-smooth counting.** Like in the classical Gauss circle problem one obtains much better estimates for the remainder term if one uses a smooth cutoff. Let  $\alpha \in C_c^{\infty}(G)$  be a non-negative test function with normalized integral. Set  $\mathbf{1}_{R,\alpha} := \alpha * \mathbf{1}_R$  and define

$$\text{err}_{\alpha}(R, \Gamma) := \sup_{\substack{\phi \in C_b(Y)_o^K \\ \|\phi\|_{\infty} \leq 1}} \frac{|\langle \mathbf{1}_{R,\alpha}^{\Gamma}, \phi \rangle|}{|B_R|} = \sup_{\substack{\phi \in C_b(Y)_o^K \\ \|\phi\|_{\infty} \leq 1}} \frac{|\langle \mathbf{1}_{R,\alpha}, \phi^H \rangle|}{|B_R|}.$$

**Lemma 7.1.** *Let  $k > s+1$  if  $Y$  is compact and  $k > \frac{r+1}{2}s+1$  otherwise. Let  $p > p_H(\Gamma)$  and  $q$  be such that  $\frac{1}{p} + \frac{1}{q} = 1$ . Then there exists  $C > 0$  such that*

$$(7.1) \quad \text{err}_{\alpha}(R, \Gamma) \leq C\|\alpha\|_{1,k}|B_R|^{-\frac{1}{p}}$$

for all  $R \geq 1$  and all  $\alpha \in C_c^{\infty}(G)$ .

*Proof.* First note that

$$\langle \mathbf{1}_{R,\alpha}, \phi^H \rangle = \langle \mathbf{1}_{R,\alpha}, (-\mathbf{1} + \Delta)^{k/2} (-\mathbf{1} + \Delta)^{-k/2} \phi^H \rangle.$$

With  $\psi = (-\mathbf{1} + \Delta)^{-k/2}\phi$  we have  $\|\psi\|_{\infty,k} \leq C\|\phi\|_{\infty}$  for some  $C > 0$ . We thus obtain

$$\begin{aligned} \text{err}_{\alpha}(R, \Gamma) &\leq C \sup_{\substack{\psi \in C_b(Y)_0^K \\ \|\psi\|_{\infty,k} \leq 1}} \frac{|\langle \mathbf{1}_{R,\alpha}, (-\mathbf{1} + \Delta)^{k/2}\psi^H \rangle|}{|B_R|} \\ &\leq \frac{C}{|B_R|} \sup_{\substack{\psi \in C_b(Y)_0^K \\ \|\psi\|_{\infty,k} \leq 1}} |\langle \mathbf{1}_{R,\alpha}, (-\mathbf{1} + \Delta)^{k/2}\psi^H \rangle| \end{aligned}$$

Moving  $(-\mathbf{1} + \Delta)^{k/2}$  to the other side we get with Hölder's inequality and Proposition 6.5 that

$$\text{err}_{\alpha}(R, \Gamma) \leq \frac{C}{|B_R|} \|(-\mathbf{1} + \Delta)^{k/2}\alpha * \mathbf{1}_R\|_q.$$

Finally,

$$\|(-\mathbf{1} + \Delta)^{k/2}\alpha * \mathbf{1}_R\|_q \leq C\|\alpha\|_{1,k}\|\mathbf{1}_R\|_q$$

and with  $\|\mathbf{1}_R\|_q = |B_R|^{\frac{1}{q}}$ , the lemma follows.  $\square$

**Remark 7.2.** In the literature results are sometimes stated not with respect to  $\text{err}(R, \Gamma)$  but the pointwise error term  $\text{err}_{pt}(R, \Gamma) = |\mathbf{1}_R^{\Gamma}(\mathbf{1}) - |B_R||$ . Likewise we define  $\text{err}_{pt,\alpha}(R, \Gamma)$ . Let  $B_Y$  be a compact neighborhood of  $\mathbf{1}\Gamma \in Y$  and note that

$$\text{err}_{pt,\alpha}(R, \Gamma) \leq |B_R| \sup_{\substack{\phi \in L^1(B_Y) \\ \|\phi\|_1 \leq 1}} |\langle \frac{\mathbf{1}_{R,\alpha}^{\Gamma}}{|B_R|} - \mathbf{1}_Y, \phi \rangle| \quad (R > 0).$$

The Sobolev estimate  $\|\phi\|_{\infty} \leq C\|\phi\|_{1,k}$ , for  $K$ -invariant functions  $\phi$  on  $B_Y$  and with  $k = \dim Y/K$  the Sobolev shift, then relates these error terms:

$$\text{err}_{pt,\alpha}(R, \Gamma) \leq |B_R| \sup_{\substack{\phi \in C_b^{\infty}(Y) \\ \|\phi\|_{\infty,-k} \leq 1}} |\langle \frac{\mathbf{1}_R^{\Gamma}}{|B_R|} - \mathbf{1}_Y, \phi \rangle|.$$

We then obtain

$$\text{err}_{pt,\alpha}(R, \Gamma) \leq C|B_R|^{1-\frac{1}{p}} \quad (R > 0)$$

in view of (7.1).

We return to the error bound in Lemma 7.1 and would like to compare  $\text{err}_1(R, \Gamma)$  with  $\text{err}_{\alpha}(R, \Gamma)$ . For that we note (by the triangle inequality) that

$$|\text{err}_1(R, \Gamma) - \text{err}_{\alpha}(R, \Gamma)| \leq \sup_{\substack{\phi \in C_b(Y)_0^K \\ \|\phi\|_{\infty} \leq 1}} \frac{|\langle \mathbf{1}_{R,\alpha}^{\Gamma} - \mathbf{1}_R^{\Gamma}, \phi \rangle|}{|B_R|}.$$



Suppose that  $\text{supp } \alpha \subset B_\epsilon^G$  for some  $\epsilon > 0$ . Then Lemma 2.2 implies that  $\mathbf{1}_{R,\alpha}$  is supported in  $B_{R+\epsilon}$ , and hence

$$\begin{aligned} |\langle \mathbf{1}_{R,\alpha}^\Gamma - \mathbf{1}_R^\Gamma, \phi \rangle| &\leq \| \mathbf{1}_{R,\alpha}^\Gamma - \mathbf{1}_R^\Gamma \|_1 \\ &\leq \| \mathbf{1}_{R,\alpha} - \mathbf{1}_R \|_1 \\ &\leq |B_{R+\epsilon}|^{\frac{1}{2}} \| \mathbf{1}_{R,\alpha} - \mathbf{1}_R \|_2 \\ &\leq |B_{R+\epsilon}|^{\frac{1}{2}} |B_{R+\epsilon} \setminus B_R|^{\frac{1}{2}}. \end{aligned}$$

With Lemma 2.1 we get

$$|B_{R+\epsilon} \setminus B_R| \leq C\epsilon |B_R| \quad (R \geq 1, \epsilon < 1).$$

Thus we obtain that

$$|\text{err}_1(R, \Gamma) - \text{err}_\alpha(R, \Gamma)| \leq C\epsilon^{\frac{1}{2}}.$$

Combining this with the estimate in Lemma 7.1 we arrive at the existence of  $C > 0$  such that

$$\text{err}_1(R, \Gamma) \leq C(\epsilon^{-k} |B_R|^{-\frac{1}{p}} + \epsilon^{\frac{1}{2}})$$

for all  $R \geq 1$  and all  $0 < \epsilon < 1$ . The minimum of the function  $\epsilon \mapsto \epsilon^{-k} C + \epsilon^{1/2}$  is attained at  $\epsilon = (2kc)^{\frac{2}{2k+1}}$  and thus we get:

**Theorem 7.3.** *Under the assumptions of Proposition 6.5 the first error term  $\text{err}(R, \Gamma)$  for the lattice counting problem on  $Z = G/H$  can be estimated as follows: for all  $p > p_H(\Gamma)$  and  $k > s + 1$  for  $Y$  compact, resp.  $k > \frac{r+1}{2}s + 1$  otherwise, there exists a constant  $C = C(p, k) > 0$  such that*

$$\text{err}(R, \Gamma) \leq C |B_R|^{-\frac{1}{(2k+1)p}}$$

for all  $R \geq 1$ .

**Remark 7.4.** The point where we lose essential information is in the estimate (6.4) where we used Weyl's law. In the moment pointwise multiplicity bounds are available the estimate would improve. To compare the results with Selberg on the hyperbolic disc, let us assume that  $p_H(\Gamma) = 2$ . Then with  $r = 1$  and  $s = 2$  our bound is  $\text{err}(R, \Gamma) \leq C_\epsilon |B_R|^{-\frac{1}{14} + \epsilon}$  while Selberg showed  $\text{err}(R, \Gamma) \leq C_\epsilon |B_R|^{-\frac{1}{3} + \epsilon}$ .

## 8. TRIPLE SPACES

In this section we verify our Hypothesis A for triple space  $Z = G/H$  where  $G = G' \times G' \times G'$ ,  $H = \text{diag}(G')$  and  $G' = \text{SO}_e(1, n)$  for some  $n \geq 2$ . Observe that  $\text{SO}_e(1, 2) \cong \text{PSI}(2, \mathbb{R})$ . We take  $K' := \text{SO}(n, \mathbb{R}) < G'$  as a maximal compact subgroup and set  $K := K' \times K' \times K'$ . Further we set  $\mathfrak{s} := \mathfrak{s}' \times \mathfrak{s}' \times \mathfrak{s}'$ . A maximal abelian subspace  $\mathfrak{a} \subset \mathfrak{s}$  is then of the form

$$\mathfrak{a} = \mathfrak{a}'_1 \times \mathfrak{a}'_2 \times \mathfrak{a}'_3$$

with  $\mathfrak{a}'_i \subset \mathfrak{s}'$  one dimensional subspaces. We recall the following result from [8].

**Proposition 8.1.** *For the triple space the following assertion hold true:*

- (1)  $G = KAH$  if and only if  $\dim(\mathfrak{a}'_1 + \mathfrak{a}'_2 + \mathfrak{a}'_3) = 2$ .
- (2) Suppose that all  $\mathfrak{a}'_i$  are pairwise distinct. Then one has  $PH$  is open for all minimal parabolics  $P$  with Langlands-decomposition  $P = M_P A_P N_P$  and  $A_P = A$ .

We say that the choice of  $A$  is *generic* if all  $\mathfrak{a}'_i$  are distinct and  $\dim(\mathfrak{a}'_1 + \mathfrak{a}'_2 + \mathfrak{a}'_3) = 2$ .

The invariant measure  $dz$  on  $Z$  can then be estimated as

$$\int_Z f(z) dz \leq \int_K \int_A f(ka \cdot z_0) J(a) da dk \quad (f \in C_c(Z), f \geq 0)$$

with  $J(a)$ , the Jacobian a non-negative real valued function on  $A$  with

$$(8.1) \quad J(a) \leq \sup_{w \in \mathcal{W}} a^{2w\rho}$$

by Lemma 3.2.

**8.1. Proof of the Hypothesis A.** We first note that for all  $\pi \in \widehat{G}_s$  the space of  $H$ -invariants

$$(\mathcal{H}_\pi^{-\infty})^H = \mathbb{C}I.$$

is one-dimensional, see [5], Thm. 3.1.

Write  $\pi = \pi_1 \otimes \pi_2 \otimes \pi_3$  with each factor a  $K'$ -spherical unitary irreducible representation of  $G'$ . If we assume that  $\pi \neq \mathbf{1}$  has non-trivial  $H$ -fixed distribution vectors, then at least two of the factors  $\pi_i$  are non-trivial.

Let  $v_i$  be normalized  $K'$ -fixed vectors of  $\pi_i$  and set  $v = v_1 \otimes v_2 \otimes v_3$ . Since  $Z$  is a multiplicity one space, the functional  $I \in (\mathcal{H}_\pi^{-\infty})^H$  is unique up to scalars. Our concern is to obtain uniform  $L^p$ -bounds for the generalized matrix coefficients  $f_\pi := m_{v,I}$ :

$$f_\pi(g_1, g_2, g_3) := I(\pi_1(g_1)^{-1}v_1 \otimes \pi_2(g_2)^{-1}v_2 \otimes \pi_3(g_3)^{-1}v_3),$$

when  $\pi$  belongs to the set  $\Lambda$  of (6.2).

We decompose  $\Lambda = \Lambda_0 \cup \Lambda_1 \cup \{\mathbf{1}\}$  with  $\Lambda_0 \subset \Lambda$  the set of  $\pi \in \Lambda$  with all  $\pi_i$  non-trivial, and  $\Lambda_1$  the set of  $\pi$ 's with exactly one  $\pi_i$  to be trivial.

Consider first the case where  $\pi \in \Lambda_1$ , i.e. one  $\pi_i$  is trivial, say  $\pi_3$ . Then  $\pi_2 = \pi_1^*$ . We identify  $Z \simeq G' \times G'$  via  $(g, h) \mapsto (\mathbf{1}, g, h)H$  and obtain

$$f_\pi(g, h) = \langle \pi_1(g)v_1, v_1 \rangle,$$

a spherical function. Note that  $Z_\eta \simeq G'$  and Hypothesis A follows from standard properties about  $K'$ -spherical functions on  $G'$ . To be more specific let  $G' = N'A'K'$  be an Iwaswa-decomposition with middle-projection  $\mathbf{a} : G' \rightarrow A'$ , then

$$f_\pi(g, h) = \varphi_{\lambda_1}(g) := \int_{K'} \mathbf{a}(k'g)^{\lambda_1 - \rho'} dk'.$$

We use Harish-Chandra's estimates  $|\varphi_\nu(a)| \leq a^\nu \varphi_0(a)$  and  $\varphi_0(a) \leq Ca^{-\rho}(1 + |\log a|)^d$  for  $a \in A'$  in positive chamber. The condition of  $\pi \in \Lambda_1$  implies that  $\rho - \operatorname{Re} \lambda_1 > 0$  is bounded away from zero and Hypothesis A follows in this case.

Suppose now that  $\pi \in \Lambda_0$ , i.e. all  $\pi_i$  are non-trivial.

For a simplified exposition we assume that  $n = 2$ , i.e.  $G' = \operatorname{PSl}(2, \mathbb{R})$ , and comment at the end for the general case. Then  $\pi_i = \pi_{\lambda_i}$  are principal series for some  $\lambda_i \in i\mathbb{R}^+ \cup [0, 1)$  with  $\mathcal{H}_\pi^\infty = C^\infty(\mathbb{S}^1)$  in the compact realization. Set  $\lambda = (\lambda_1, \lambda_2, \lambda_3)$  and set  $\pi = \pi_\lambda$ .

In order to analyze  $f_\pi$  we use  $G = KAH$  and thus assume that  $g = a = (a_1, a_2, a_3) \in A$ . We work in the compact model of  $\mathcal{H}_{\pi_i} = L^2(\mathbb{S}^1)$  and use the explicit model for  $I$  in [3]: for  $h_1, h_2, h_3$  smooth functions on the circle one has

$$I(h_1 \otimes h_2 \otimes h_3) = \frac{1}{(2\pi)^3} \int_0^{2\pi} \int_0^{2\pi} \int_0^{2\pi} h_1(\theta_1) h_2(\theta_2) h_3(\theta_3) \cdot \mathcal{K}(\theta_1, \theta_2, \theta_3) d\theta_1 d\theta_2 d\theta_3,$$

where

$$\mathcal{K}(\theta_1, \theta_2, \theta_3) = |\sin(\theta_2 - \theta_3)|^{(\alpha-1)/2} |\sin(\theta_1 - \theta_3)|^{(\beta-1)/2} |\sin(\theta_1 - \theta_2)|^{(\gamma-1)/2}.$$

In this formula one has  $\alpha = \lambda_1 - \lambda_2 - \lambda_3$ ,  $\beta = -\lambda_1 + \lambda_2 - \lambda_3$  and  $\gamma = -\lambda_1 - \lambda_2 + \lambda_3$  where  $\lambda_i \in i\mathbb{R} \cup (-1, 1)$  are the standard representation parameters of  $\pi_i$ . According to [5], Cor. 2.1, the kernel  $\mathcal{K}$  is absolutely integrable.

Set

$$A' := \left\{ a_t := \begin{pmatrix} t & 0 \\ 0 & \frac{1}{t} \end{pmatrix} \mid t > 0 \right\} < G'$$

Then  $A'_i = k_{\phi_i} A' k_{\phi_i}^{-1}$  with  $\phi_i \in [0, 2\pi]$  and

$$k_\phi = \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix}.$$

Set  $a_{t,i} = k_{\phi_i} a_t k_{\phi_i}^{-1}$ .

Returning to our analysis of  $f_\pi$  we now take  $h_i(t_i, \theta_i) = [\pi_1(a_{t_i, i})v_i](\theta_i)$  and remark that

$$h_i(t_i, \theta_i) = \frac{1}{(t_i^2 + \sin^2(\theta_i - \phi_i)(\frac{1}{t_i^2} - t_i^2))^{\frac{1}{2}(1+\lambda_i)}}.$$

Let us set  $|\pi| := \pi_{\operatorname{Re} \lambda_1} \otimes \pi_{\operatorname{Re} \lambda_2} \otimes \pi_{\operatorname{Re} \lambda_3}$ . Say  $|\pi| \geq |\pi'|$  if  $|\operatorname{Re} \alpha| \leq |\operatorname{Re} \alpha'|$ ,  $|\operatorname{Re} \beta| \leq |\operatorname{Re} \beta'|$  and  $|\operatorname{Re} \gamma| \leq |\operatorname{Re} \gamma'|$ . Our formulas then show for  $|\pi| \leq |\pi'|$  that

$$(8.2) \quad |f_\pi(a)| \leq f_{|\pi|}(a) \leq f_{|\pi'|}(a) \quad (a \in A).$$

Let  $c_i := 1 - |\operatorname{Re} \lambda_i|$  for  $i = 1, 2, 3$ . The fundamental estimate in [21], Thm. 3.2, then yields a constant  $d$ , independent of  $\pi$ , and a constant  $C = C(\pi) > 0$  such that for  $a = (a_{t_1, 1}, a_{t_2, 2}, a_{t_3, 3})$  one has

$$(8.3) \quad |f_\pi(a)| \leq C \frac{(1 + |\log t_1| + |\log t_2| + |\log t_3|)^d}{[\cosh \log t_1]^{c_1} \cdot [\cosh \log t_2]^{c_2} \cdot [\cosh \log t_3]^{c_3}}.$$

In view of (8.2) the constant  $C(\pi)$  depends only on the distance of  $\operatorname{Re} \lambda_i$  to the trivial representation. Looking at the integral representation of  $f_\pi$  with the kernel  $\mathcal{K}$  we deduce a lower bound without the logarithmic factor, i.e. the bound is essentially sharp. Hence (8.1) together with the fact that all  $f_\pi$  for  $\pi \in \Lambda_0$  are in  $L^p(Z)$  for some  $p < \infty$  implies that

$$(8.4) \quad \inf_{\pi \in \Lambda_0} c_i(\pi) > 0.$$

Further (8.2) and (8.3) together give

$$(8.5) \quad \sup_{\pi \in \Lambda_0} C(\pi) < \infty.$$

In particular we get both

$$(8.6) \quad \sup_{\pi \in \Lambda_0} \|f_\pi\|_p < \infty.$$

and

$$(8.7) \quad \sup_{\pi \in \Lambda_0} \|f_\pi\|_\infty < \infty.$$

On the other hand for  $g = \mathbf{1} = (\mathbf{1}, \mathbf{1}, \mathbf{1})$ , the value  $f_\pi(\mathbf{1})$  is obtained by applying  $I$  to the constant function  $\mathbf{1} = \mathbf{1} \otimes \mathbf{1} \otimes \mathbf{1}$ . This value has been computed explicitly by Bernstein and Reznikov in [3] as

$$\frac{\Gamma((\alpha + 1)/4)\Gamma((\beta + 1)/4)\Gamma((\gamma + 1)/4)\Gamma((\delta + 1)/4)}{\Gamma((1 - \lambda_1)/2)\Gamma((1 - \lambda_2)/2)\Gamma((1 - \lambda_3)/2)}$$

where  $\alpha, \beta, \gamma$  are as before and  $\delta = -\lambda_1 - \lambda_2 - \lambda_3$ . Stirling approximation,

$$|\Gamma(\sigma + it)| = \text{const.} e^{-\frac{\pi}{2}|t|} |t|^{\sigma - \frac{1}{2}} (1 + O(|t|^{-1}))$$

as  $|t| \rightarrow \infty$  and  $\sigma$  is bounded, yields a lower bound for  $f_\pi(\mathbf{1})$ :

$$(8.8) \quad \inf_{\pi \in \Lambda_0} |f_\pi(\mathbf{1})| > 0.$$

As  $\|f_\pi\|_\infty \geq |f_\pi(\mathbf{1})|$  the assertion (A1) of Hypothesis A is readily obtained from (8.6) and (8.8). Likewise (A2) with  $\Omega = \{\mathbf{1}\}$  follows from (8.7) and (8.8).

In general for  $G' = \text{SO}_e(1, n)$  one needs to compute the Bernstein-Reznikov integral. This was accomplished in [9].

**Theorem 8.2.** *Let  $Z = G' \times G' \times G' / \text{diag}(G')$  for  $G' = \text{SO}_e(1, n)$  and assume that  $H/\Gamma_H$  is compact. Then the first error term  $\text{err}(R, \Gamma)$  for the lattice counting problem on  $Z = G/H$  can be estimated as follows: for all  $p > p_H(\Gamma)$  there exists a  $C = C(p) > 0$  such that*

$$\text{err}(R, \Gamma) \leq C |B_R|^{-\frac{1}{(6n+3)p}}$$

for all  $R \geq 1$ .

**8.2. Cubic lattices.** Here we let  $G_0 = \text{SO}_e(1, 2)$  with the quadratic  $Q$  form defining  $G_0$  having integer coefficients and anisotropic over  $\mathbb{Q}$ , for example

$$Q(x_0, x_1, x_2) = 2x_0^2 - 3x_1^2 - x_2^2.$$

Then, according to Borel,  $\Gamma_0 = G_0(\mathbb{Z})$  is a uniform lattice in  $G_0$ .

Next let  $k$  be a cubic Galois extension of  $\mathbb{Q}$ . Note that  $k$  is totally real. An example of  $k$  is the splitting field of the polynomial  $f(x) = x^3 + x^2 - 2x - 1$ . Let  $\sigma$  be a generator of the Galois group of  $k|\mathbb{Q}$ . Let  $\mathcal{O}_k$  be the ring of algebraic integers of  $k$ . We define  $\Gamma < G = G_0^3$  to be the image of  $G_0(\mathcal{O}_k)$  under the embedding

$$G_0(\mathcal{O}_k) \ni \gamma \mapsto (\gamma, \gamma^\sigma, \gamma^{\sigma^2}) \in G.$$

Then  $\Gamma < G$  is a uniform irreducible lattice with trace  $H \cap \Gamma \simeq \Gamma_0$  a uniform lattice in  $H \simeq G_0$ .

## 9. OUTLOOK

We discuss some topics of harmonic analysis on reductive homogeneous spaces which are currently open and would have immediate applications to lattice counting.

**9.1. A conjecture which implies Hypothesis A.** Hypothesis A falls in the context of a more general conjecture about the growth behavior of families of Harish-Chandra modules.

We let  $Z = G/H$  be a real spherical space. Denote by  $A_Z^- \subset A_Z$  the compression cone of  $Z$  (see Section 3) and recall that wavefront means that  $A^- A_H / A_H = A_Z^-$  which, however, we do not assume for the moment.

We use  $V$  to denote Harish-Chandra modules for the pair  $(\mathfrak{g}, K)$  and  $V^\infty$  for their unique moderate growth smooth Fréchet globalizations. These  $V^\infty$  are global objects in the sense that they are  $G$ -modules whereas  $V$  is defined in algebraic terms. We write  $V^{-\infty}$  for the strong dual of  $V^\infty$ . We say that  $V$  is  $H$ -distinguished provided that the space of  $H$ -invariants  $(V^{-\infty})^H$  is non-trivial.

It is no big loss of generality to assume that  $A_Z^-$  is a sharp cone, as the edge of this cone is in the normalizer of  $H$  and in particular acts on the finite dimensional space of  $H$ -invariants.

As  $A_Z^-$  is pointed it is a fundamental domain for the little Weyl group and as such a simplicial cone (see [16]). If  $\mathfrak{a}_Z^- = \log A_Z^-$ , then we write  $\omega_1, \dots, \omega_r$  for a set of generators (spherical co-roots) of  $\mathfrak{a}_Z^-$ .

Set  $\bar{Q} := \theta(Q)$  where  $\theta$  is the Cartan involution determined by the choice of  $K$ . Note that  $V/\bar{\mathfrak{q}}V$  is a finite dimensional  $\bar{Q}$  module, in particular a finite dimensional  $A_Z$ -module. Let  $\Lambda_1, \dots, \Lambda_N \in \mathfrak{a}_Z^*$  be the  $\mathfrak{a}_{Z, \mathbb{C}}$ -weight spectrum. Then we define the  $H$ -spherical exponent  $\Lambda_V \in \mathfrak{a}_Z^*$  of  $V$  by

$$\Lambda_V(\omega_i) := \max_{1 \leq j \leq N} \operatorname{Re} \Lambda_j(\omega_i).$$

Further attached to  $V$  is a “logarithmic” exponent  $d \in \mathbb{N}$ . Having this data we recall the main bound from [21]

$$|m_{v, \eta}(a \cdot z_0)| \lesssim a^{\Lambda_V} (1 + \|\log a\|)^{d_V} \quad (a \in A_Z^-).$$

**Conjecture 9.1.** *Fix a  $K$ -type  $\tau$ , a constant  $C > 0$ , and a compact subset  $\Omega \subset G$ . Then there exists a compact set  $\Omega_A \subset A_Z^-$  such that for all Harish-Chandra modules  $V$  with  $\|\Lambda_V\| \leq C$ , and all  $v \in V[\tau]$  one has*

$$\begin{aligned} \max_{\substack{a \in A_Z^- \\ g \in \Omega}} |m_{v, \eta}(ga \cdot z_0)| a^{-\Lambda_V} (1 + \|\log a\|)^{-d} = \\ \max_{\substack{a \in \Omega_A \\ g \in \Omega}} |m_{v, \eta}(ga \cdot z_0)| a^{-\Lambda_V} (1 + \|\log a\|)^{-d}. \end{aligned}$$

It is easily seen that, if true this will imply Hypothesis A.

**Remark 9.2.** It might well be that a slightly stronger conjecture is true. For that we recall that a Harish-Chandra module  $V$  has a unique

minimal globalization, the analytic model  $V^\omega$ . The space  $V^\omega$  is an increasing union of subspaces  $V_\epsilon$  for  $\epsilon \rightarrow 0$ . The parameter  $\epsilon$  parametrizes left  $G$ -invariant neighborhoods  $\Xi_\epsilon \subset G_{\mathbb{C}}$  of  $\mathbf{1}$  which decrease with  $\epsilon \rightarrow 0$ . Further  $V_\epsilon$  consists of those vectors  $v \in V^\omega$  for which the orbit map  $G \rightarrow V^\omega$ ,  $g \mapsto g \cdot v$  extends to a holomorphic map on  $\Xi_\epsilon$ . For fixed  $\epsilon, C > 0$  the strengthened conjecture would be that there exists a compact subset  $\Omega_A$  such that for all Harish-Chandra modules  $V$  with  $\|\Lambda_V\| \leq C$  and all  $v \in V_\epsilon$  one has

$$\begin{aligned} & \max_{a \in A_{\mathbb{Z}}^-} |m_{v,\eta}(a \cdot z_0)| a^{-\Lambda_V} (1 + \|\log a\|)^{-d} = \\ & \max_{a \in \Omega_A} |m_{v,\eta}(a \cdot z_0)| a^{-\Lambda_V} (1 + \|\log a\|)^{-d}. \end{aligned}$$

Note that the compact set  $\Omega$  is no longer needed, as  $\Omega \cdot V_\epsilon \subset V_\epsilon$ .

**9.2. Spectral geometry of  $Z_\eta$ .** In the general context of a reductive real spherical space it may be possible to establish both main term counting and the error term bound, with the arguments presented here for wavefront spaces, provided the following two key questions allow affirmative answers.

In what follows  $Z = G/H$  is a real reductive spherical space and  $V$  denotes an irreducible Harish-Chandra module and  $\eta \in (V^{-\infty})^H$ .

**Question A:** *Is  $H_\eta$  reductive?*

**Question B:** *If for  $v \in V$  the generalized matrix coefficient  $m_{v,\eta}$  is bounded, then there exists a  $1 \leq p < \infty$  such that  $m_{v,\eta} \in L^p(Z_\eta)$ .*

## REFERENCES

- [1] Y. Benoist and H. Oh, *Effective equidistribution of  $S$ -integral points on symmetric varieties*, Ann. Inst. Fourier (Grenoble) **62** (2012), no. **5**, 1889–1942.
- [2] J. Bernstein and B. Krötz, *Smooth Fréchet globalizations of Harish-Chandra modules*, Israel J. Math. **199** (2014), 35 – 111.
- [3] J. Bernstein and A. Reznikov, *Estimates of automorphic functions*, Moscow Math. J. **4** (1) (2004), 19–37
- [4] P. Bravi and G. Pezzini, *The spherical systems of the wonderful reductive subgroups*, arxiv: 1109.6777
- [5] J.-L. Clerc and B. Ørsted, *Conformally invariant trilinear forms on the sphere*. Ann. Inst. Fourier **61** (2011), 18071838.
- [6] L. Clozel, *Démonstration de la conjecture  $\tau$* , Invent. Math. **151** (2003), no. **2**, 297–328.
- [7] M. Cowling, *Sur les coefficients des représentations unitaires des groupes de Lie simples*, Lecture Notes in Mathematics **739** (1979), 132–178.
- [8] T. Danielsen, B. Krötz and H. Schlichtkrull, *Decomposition theorems for triple spaces*, Geom. Dedicata (to appear).

- [9] A. Deitmar, *Invariant triple products*, Int. J. Math. Math. Sci. 2006, Art. ID 48274, 22 pp.
- [10] W. Duke, Z. Rudnick and P. Sarnak, *Density of integer points on affine homogeneous varieties*, Duke Math. J. **71** (1993), no. 1, 143–179.
- [11] A. Eskin and C. McMullen, *Mixing, counting, and equidistribution in Lie groups*, Duke Math. J. 71 (1993), no. 1, 181V209.
- [12] A. Eskin, S. Mozes and N. Shah, *Unipotent flows and counting lattice points on homogeneous varieties*, Ann. of Math. (2) **143** (1996), 253–299.
- [13] A. Gorodnik and A. Nevo, *Counting lattice points*, J. Reine Angew. Math. **663** (2012), 127–176.
- [14] Harish-Chandra, *Automorphic Forms on Semisimple Lie Groups*, Springer LNM 62.
- [15] L. Ji, *The Weyl upper bound on the discrete spectrum of locally symmetric spaces*, J. Diff. Geom. **51** (1) (1999), 97–147.
- [16] F. Knop and B. Krötz, *A  $k$ -rational local structure theorem*, in preparation
- [17] F. Knop, B. Krötz, E. Sayag and H. Schlichtkrull, *Simple compactifications and polar decomposition for real spherical spaces*, arXiv:1402.3467
- [18] ———, *Volume growth, temperedness and integrability of matrix coefficients on a real spherical space*, arXiv: 1407.8006
- [19] F. Knop, B. Krötz and H. Schlichtkrull, *The local structure theorem for real spherical varieties*, arXiv:1310.6390.
- [20] B. Krötz, E. Sayag and H. Schlichtkrull, *Vanishing at infinity on homogeneous spaces of reductive type*, arXiv: 1211.2781.
- [21] ———, *Decay of matrix coefficients on reductive homogeneous space of spherical type*, Math. Z. DOI 10.1007/s00209-014-1313-7 (to appear).
- [22] B. Krötz and H. Schlichtkrull, *Multiplicity bounds and the subrepresentation theorem for real spherical spaces*, Trans. Amer. Math. Soc. (to appear).
- [23] Y. Sakellaridis and A. Venkatesh, *Periods and harmonic analysis on spherical varieties*, arXiv:1203.0039