

# HYBRID RANDOM FEATURES

Krzysztof Choromanski<sup>\*1,2</sup>, Haoxian Chen<sup>\*†2</sup>, Han Lin<sup>\*†2</sup>, Yuanzhe Ma<sup>\*†2</sup>, Arijitohanobish<sup>\*†‡</sup>,  
 Deepali Jain<sup>1</sup>, Michael S Ryoo<sup>1</sup>, Jake Varley<sup>1</sup>, Andy Zeng<sup>1</sup>, Valerii Likhoshesterov<sup>1,3</sup>,  
 Dmitry Kalashnikov<sup>1</sup>, Vikas Sindhwani<sup>1</sup>, Adrian Weller<sup>3,4</sup>

<sup>1</sup>Google Brain Robotics, <sup>2</sup>Columbia University, <sup>3</sup>University of Cambridge, <sup>4</sup>The Alan Turing Institute

## ABSTRACT

We propose a new class of random feature methods for linearizing softmax and Gaussian kernels called *hybrid random features* (HRFs) that automatically adapt the quality of kernel estimation to provide most accurate approximation in the defined regions of interest. Special instantiations of HRFs lead to well-known methods such as trigonometric (Rahimi & Recht, 2007) or (recently introduced in the context of linear-attention Transformers) positive random features (Choromanski et al., 2021b). By generalizing Bochner’s Theorem for softmax/Gaussian kernels and leveraging random features for compositional kernels, the HRF-mechanism provides strong theoretical guarantees - unbiased approximation and strictly smaller worst-case relative errors than its counterparts. We conduct exhaustive empirical evaluation of HRF ranging from pointwise kernel estimation experiments, through tests on data admitting clustering structure to benchmarking implicit-attention Transformers (also for downstream Robotics applications), demonstrating its quality in a wide spectrum of machine learning problems.

## 1 INTRODUCTION & RELATED WORK

Consider the *softmax* and *Gaussian kernel* functions  $K : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  defined as follows:

$$\text{SM}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \exp(\mathbf{x}^\top \mathbf{y}), \quad K_{\text{gauss}}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{2}\right). \quad (1)$$

These two are prominent examples of functions used in the so-called *kernel methods* (Gretton et al., 2005; Zhang et al., 2018) and beyond, i.e. in softmax-sampling (Blanc & Rendle, 2018). Random features (RFs, Rahimi & Recht, 2007; Liu et al., 2020; Peng et al., 2021) yield a powerful mechanism for linearizing and consequently scaling up kernel methods with dot-product kernel decompositions disentangling  $\mathbf{x}$  from  $\mathbf{y}$  in the formulae for kernel value  $K(\mathbf{x}, \mathbf{y}) \approx \phi(\mathbf{x})^\top \phi(\mathbf{y})$  via data-agnostic probabilistic (random feature) maps  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ . The tight relationship between softmax and Gaussian kernels given by the transformation  $K_{\text{gauss}}(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x}\|_2^2}{2})\text{SM}(\mathbf{x}, \mathbf{y})\exp(-\frac{\|\mathbf{y}\|_2^2}{2})$  provides a mapping of any random feature vector  $\phi_{\text{SM}}(\mathbf{u})$  for the softmax kernel to the corresponding one  $\phi_{\text{gauss}}(\mathbf{u}) = \exp(-\frac{\|\mathbf{u}\|_2^2}{2})\phi_{\text{SM}}(\mathbf{u})$  for the Gaussian kernel, thus we will focus on the former kernel. The classic random feature map mechanism  $\phi_m^{\text{trig}}$  for the softmax kernel, obtained from Bochner’s Theorem applied to the Gaussian kernel (Rahimi & Recht, 2007), is of the form:

$$\phi_m^{\text{trig}}(\mathbf{u}) = \frac{1}{\sqrt{m}} \exp\left(-\frac{\|\mathbf{u}\|_2^2}{2}\right) (\sin(\omega_1^\top \mathbf{u}), \dots, \sin(\omega_m^\top \mathbf{u}), \cos(\omega_1^\top \mathbf{u}), \dots, \cos(\omega_m^\top \mathbf{u}))^\top, \quad (2)$$

where  $m$  stands for the number of random features and  $\omega_1, \dots, \omega_m \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ .

The above (common in most downstream use-cases) method for linearizing softmax/Gaussian kernels was recently shown to fail in some of the new impactful applications of scalable kernel methods such as implicit-attention Transformer architectures called Performers (Choromanski et al., 2021b). We denote by MSE the mean squared error of the estimator (i.e. its variance since all estimators considered in this paper are unbiased). The above mechanism struggles to accurately approximate close-to-zero kernel values, as characterized by the particularly large MSE in that region. This is a crucial problem since most of the entries of the attention matrices in Transformers’ models are very

\*Equal Contribution, Correspondence to kchoro@google.com.

† Authorship in alphabetical order

‡ Independent researcher. Work done during postdoc at Yale University

small and the approximators need to be particularly accurate there. Otherwise the renormalizers computed to make attention matrices row-stochastic (standard attention normalization procedure in Transformers) would be estimated imprecisely, potentially even by negative values.

The solution to this problem was presented in FAVOR+ mechanism (Choromanski et al., 2021b), where a new positive random feature map for unbiased softmax kernel estimation was applied:

$$f_m^+(u) = \frac{1}{2m} \exp\left(\frac{ku^2}{2}\right) \exp(i \angle u); \dots; \exp(i \angle_m u); \exp(-i \angle u); \dots; \exp(-i \angle_m u) \quad (3)$$

Even though, very accurate in estimating small softmax kernel values (which turns out to be crucial in making RFs work for Transformers training), this mechanism is characterized by larger MSE for large kernel values. In several applications of softmax kernels (in particular Transformers, where attention matrices typically admit sparse combinatorial structure with relatively few but critical large entries and several close-to-zero ones or softmax sampling) the algorithm needs to process simultaneously very small and large kernel values. The following natural questions arise:

Is it possible to get the best from both the mechanisms to obtain RFs-based estimators particularly accurate for both very small and large softmax kernel values? Furthermore, can those estimators be designed to have low variance in more generally pre-defined regions?

We give affirmative answers to both of the above questions by constructing a new class of random feature maps techniques called hybrid random features or HRFs. Theoretical methods used by us to develop HRFs: (a) provide a unifying perspective, where trigonometric random features from Bochner's Theorem and novel mechanisms proposed in (Choromanski et al., 2021b) are just special corollaries of the more general result from complex analysis, (b) integrate in the original way several other powerful probabilistic techniques such as Goemans-Williamson method (Goemans & Williamson, 2004) and random features for the compositional kernels (Daniely et al., 2017).

We provide detailed theoretical analysis of HRFs, showing in particular that they provide strictly more accurate worst-case softmax kernel estimation than previous algorithms and lead to computational gains. We also conduct their thorough empirical evaluation on tasks ranging from pointwise kernel estimation to downstream problems involving training Transformer-models or even end-to-end robotic controller-stacks including attention-based architectures.

**Related Work:** The literature on different random feature map mechanisms for Gaussian (and thus also softmax) kernel estimation is voluminous. Most focus has been put on reducing the variance of trigonometric random features from (Rahimi & Recht, 2007) via various Quasi Monte Carlo (QMC) methods, where directions and/or lengths of Gaussian vectors used to produce features are correlated, often through geometric conditions such as orthogonality (Choromanski et al., 2017; Rowland et al., 2018; Yu et al., 2016; Choromanski et al., 2019; Choromanski & Sindhwani, 2016). Our HRFs do not compete with those techniques (and can be in fact easily combined with them) since rather than focusing on improving sampling mechanism for a given approximation algorithm, they provide a completely new algorithm. The new application of random features for softmax kernel in Transformers proposed in (Choromanski et al., 2020; 2021b) led to fruitful research on the extensions and limitations of these methods. Schlag et al. (2021) replaced random features by sparse deterministic constructions (no longer approximating softmax kernel). Luo et al. (2021) observed that combining  $L_2$ -normalization of queries and keys for variance reduction of softmax kernel estimation with FFT-based implementations of relative position encoding and FAVOR+ mechanism from Performers helps in training. Trigonometric random features were applied for softmax sampling in (Rawat et al., 2019). Several other techniques such as Nyström method (Yang et al., 2012; Williams & Seeger, 2000; Rudi et al., 2015) were proposed to construct data-dependent feature representations. Even though, as we show in Sec. 2.4, certain instantiations of the HRF mechanism clearly benefit from some data analysis, our central goal remains an unbiased estimation of the softmax/Gaussian kernels which is no longer the case for those other techniques.

## 2 HYBRID RANDOM FEATURES

### 2.1 PRELIMINARIES

Whenever we do not say explicitly otherwise, all presented lemmas and theorems are new. We start with the following basic definitions and results.

**Definition 2.1 (Kernel with a Random Feature Map Representation)** We say that a kernel function  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  admits a random feature (RF) map representation if it can be written as

$$K(x; y) = E_{\mathbf{u}} \left[ \prod_{i=1}^m \chi_i(x; \mathbf{u}) \chi_i(y; \mathbf{u}) \right]; \quad (4)$$

for some  $\chi_i : \mathbb{R}^d \rightarrow \mathbb{R}$ , and where  $\mathbf{u}$  is sampled from some probabilistic distribution  $\mathcal{P}(\mathbb{R}^d)$ . The corresponding random feature map, for a given  $N$ , is defined as

$$\phi_m(\mathbf{u}) = \left( \chi_1(\mathbf{u}; \mathbf{u}_1), \dots, \chi_m(\mathbf{u}; \mathbf{u}_m) \right)^T; \quad (5)$$

where  $\mathbf{u}_m = (\mathbf{u}_1, \dots, \mathbf{u}_m)^T$ ,  $\mathbf{u}_i$  stands for vertical concatenation, and  $\mathbf{u}_m$  is iid.

Random feature maps can be used to unbiasedly approximate corresponding kernels, as follows:

$$\hat{K}(x; y) = \phi_m(x)^T \phi_m(y); \quad (6)$$

Using the above notation, a trigonometric random feature from Equation 2 can be encoded as applying  $\chi_1(\mathbf{u}; \mathbf{u}) = \sin(\mathbf{u}^T \mathbf{x})$ ,  $\chi_2(\mathbf{u}; \mathbf{u}) = \cos(\mathbf{u}^T \mathbf{x})$ . Similarly, positive random features can be encoded as taking  $\chi_1(\mathbf{u}; \mathbf{u}) = \frac{1}{\sqrt{2}} \exp(\mathbf{u}^T \mathbf{x})$ ,  $\chi_2(\mathbf{u}; \mathbf{u}) = \frac{1}{\sqrt{2}} \exp(-\mathbf{u}^T \mathbf{x})$ .

The following result from (Choromanski et al., 2021b) shows that the mean squared error (MSE) of the trigonometric estimator is small for large softmax kernel values and large for small softmax kernel values, whereas an estimator applying positive random features behaves in the opposite way.

Denote by  $\hat{SM}_m^{\text{trig}}(x; y)$  an estimator of  $SM(x; y)$  for  $x, y \in \mathbb{R}^d$  using trigonometric RFs and  $\mathbf{u}_1, \dots, \mathbf{u}_m \stackrel{\text{iid}}{\sim} \mathcal{N}(0; I_d)$ . Denote by  $\hat{SM}_m^{++}(x; y)$  its analogue using positive RFs. We have:

Lemma 2.2 (positive versus trigonometric RFs) Take  $\hat{SM}_m^{\text{trig}}(x; y) = \frac{1}{2m} \sum_{i=1}^m \cos(\mathbf{u}_i^T (x - y))$ ,  $\hat{SM}_m^{++}(x; y) = \frac{1}{2m} \sum_{i=1}^m \exp(\mathbf{u}_i^T (x - y))$ ,  $f_1(u) = (2m)^{-1} \exp(u^2)$ ,  $f_2(u) = (2m)^{-1} \exp(-u^2)$ ,  $f_3(u) = (1 - \exp(-u^2))^2$ . The MSEs of these estimators are:

$$\text{MSE}(\hat{SM}_m^{\text{trig}}(x; y)) = f_1(k_1 k_2) f_3(k_1 k_2); \quad \text{MSE}(\hat{SM}_m^{++}(x; y)) = f_2(k_1 k_2) f_3(k_1 k_2); \quad (7)$$

## 2.2 THE ALGORITHM

We are ready to present the mechanism Hybrid Random Features (HRFs). Denote by  $\mathcal{E} = \{\hat{SM}_m^k(x; y)\}_{k=1}^{p+1}$  a list of estimators of  $SM(x; y)$  (the so-called base estimators) and by  $\mathcal{B} = \{b^k(x; y)\}_{k=1}^p$  a list of estimators of  $k(x; y)$  for some functions  $k : \mathbb{R}^d \rightarrow \mathbb{R}^d \rightarrow [0, 1]$ , constructed independently from  $\mathcal{E}$ . Take the following estimator  $\hat{SM}(x; y)$ :

$$\hat{SM}^{\mathcal{E}}(x; y) = \sum_{k=1}^p b^k(x; y) \hat{SM}_m^k(x; y) + \frac{1}{p+1} \sum_{k=1}^p b^k(x; y) \hat{SM}_m^{p+1}(x; y); \quad (8)$$

In the next section, we explain in detail how base estimators are chosen. We call  $\hat{SM}^{\mathcal{E}}(x; y)$  a hybrid random feature (HRF) estimator  $\hat{SM}(x; y)$  parameterized by  $\mathcal{E}$ . The role of the coefficients is to dynamically (based on the input  $(x, y)$ ) prioritize or deprioritize certain estimators to promote those which are characterized by lower variance for a given input. Note that if elements of  $\mathcal{E}$  are unbiased estimators of  $SM(x; y)$ , then trivially  $\hat{SM}^{\mathcal{E}}(x; y)$  is also an unbiased estimator of  $SM(x; y)$ . Assume that each  $\hat{SM}_m^k(x; y)$  is of the form  $\hat{SM}_m^k(x; y) = \left( \chi_{1;m}^k(x) \right)^T \chi_{2;m}^k(y)$  for  $\chi_{j;m}^k(\mathbf{u}) = \frac{1}{\sqrt{m}} \left( \chi_{j;m}^{1;k}(\mathbf{u}), \dots, \chi_{j;m}^{t_k;k}(\mathbf{u}) \right)^T$ ,  $t_k > 0$  and  $\chi_{j;m}^{1;k}, \dots, \chi_{j;m}^{t_k;k} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ , where  $j \in \{1, 2\}$ . Assume also that  $k(x; y)$  can be written as:

$$k(x; y) = a_k + E \left[ \sum_{i=1}^{X^k} f_{1;k}^i(x) f_{2;k}^i(y) \right]; \quad (9)$$

for some scalars  $a_k \in \mathbb{R}$ , distribution  $\mathcal{P}(\mathbb{R}^d)$  (where  $\mathcal{P}(\mathbb{R}^d)$  stands for the set of probabilistic distributions on  $\mathbb{R}^d$ ), mappings  $f_{i;k}^i : \mathbb{R}^d \rightarrow \mathbb{R}$  and that the corresponding estimator  $b^k(x; y) = b_n^k(x; y)$  of  $k(x; y)$  is of the form:

$$b_n^k(x; y) = a_k + \left( \chi_{1;n}^k(x) \right)^T \chi_{2;n}^k(y); \quad (10)$$

for  $f_{j;n}^k(u) = \frac{1}{p} \sum_{i=1}^p f_{j;n}^{i;k}(u)$  and  $f_{j;n}^{i;k}(u) = (f_{j;k}^i(u; \tau_1), \dots, f_{j;k}^i(u; \tau_n))^T$ , and where  $\tau_j \in \{1, 2, \dots, g\}$ . Linearization of the coefficients given by Equation 9 is crucial to obtain linearization of the hybrid estimators, and consequently random feature map decomposition.

We denote a hybrid estimator using random features for its base estimators and approximate coefficients as  $\text{SM}_{m;n}^{\text{hyb}}$ . Furthermore, for two vectors  $u, v$ , we denote their vectorized outer-product by  $uv^T$ . Finally, we denote by  $\oplus$  vectors-concatenation. Estimator  $\text{SM}_{m;n}^{\text{hyb}}(x; y)$  can be rewritten as a dot-product of two (hybrid) random feature vectors, as the next lemma shows.

Lemma 2.3. The HRF estimator  $\text{SM}_{m;n}^{\text{hyb}}(x; y)$  satisfies  $\text{SM}_{m;n}^{\text{hyb}}(x; y) = \langle \phi(x), \phi(y) \rangle$ , where  $\phi_j(z) = \frac{1}{m} \sum_{k=1}^p \sum_{i=1}^g f_{j;n}^{i;k}(z)$  for  $j \in \{1, 2, \dots, g\}$  is given as  $\phi_j(z) = \frac{1}{m} \sum_{k=1}^p \sum_{i=1}^g f_{j;n}^{i;k}(z)$  and:

$$\begin{aligned} \phi_j^1(z) &= \frac{1}{m} \sum_{k=1}^p \sum_{i=1}^g f_{j;n}^{i;k}(z) \\ \phi_j^2(z) &= \frac{1}{mn} \sum_{k=1}^p \sum_{i=1}^g \sum_{l=1}^g \sum_{f=1}^g \sum_{t=1}^g f_{j;n}^{i;k}(z) f_{j;n}^{l;t}(z) \\ \phi_j^3(z) &= \frac{1}{m} \sum_{k=1}^p \sum_{i=1}^g \sum_{l=1}^g \sum_{f=1}^g \sum_{t=1}^g \sum_{p=1}^g f_{j;n}^{i;k}(z) f_{j;n}^{l;p}(z) \\ \phi_j^4(z) &= \frac{1}{mn} \sum_{k=1}^p \sum_{i=1}^g \sum_{l=1}^g \sum_{f=1}^g \sum_{t=1}^g \sum_{p=1}^g f_{j;n}^{i;k}(z) f_{j;n}^{l;p}(z) \end{aligned} \quad (11)$$

Bipolar estimators: A prominent special case of the general hybrid estimator defined above is the one where  $E = (\text{SM}_m^{++}(x; y); \text{SM}_m^{\text{trig}}(x; y))$ . Thus consider the following estimator  $\text{SM}_{m;n}^{\text{hyb}}$ :

$$\text{SM}_{m;n}^{\text{hyb}}(x; y) = b_n(x; y) \text{SM}_m^{++}(x; y) + (1 - b_n(x; y)) \text{SM}_m^{\text{trig}}(x; y); \quad (12)$$

The question arises whether  $\text{SM}_{m;n}^{\text{hyb}}$  defined in such a way can outperform both  $\text{SM}_m^{++}$  and  $\text{SM}_m^{\text{trig}}$ . That of course depends also on the choice of  $R^d \rightarrow R$ . If we consider a (common) normalized setting where all input vectors have the same norm, we can rewrite  $\langle x; y \rangle = r \cos(\theta_{x,y})$ , where  $\theta_{x,y}$  is an angle between  $x$  and  $y$  and  $\|x\| = \|y\| = r$ . By our previous analysis we know that  $\text{SM}_m^{++}$  becomes perfect for  $\theta = \pi/2$  and  $\text{SM}_m^{\text{trig}}$  becomes perfect for  $\theta = 0$ . That suggests particularly simple linear dependence of  $b_n$  to guarantee vanishing variance for both the critical values:  $b_n = 0$  and  $b_n = 1$ . It remains to show that such a coefficient can be linearized. It turns out that this can be done with a particularly simple random feature map mechanism, as we will show later, leading to the so-called angular hybrid variant (see: Section 2.4).

### 2.3 CHOOSING BASE ESTIMATORS FOR HRFs: COMPLEX EXPONENTIAL ESTIMATORS

In this section, we explain how we choose base estimators in Equation 8. We denote by  $\tau \in \mathbb{C}$  a complex number such that  $\tau^2 = -1$ . The following lemma is a gateway to construct unbiased base estimators.

Lemma 2.4. Let  $\mathbf{z} \in \mathbb{C}^d$ . Then for every  $\mathbf{z} = (z_1, \dots, z_d)^T \in \mathbb{C}^d$  the following holds:

$$E[\exp(\tau \langle \mathbf{z}, \mathbf{z} \rangle)] = \exp\left(-\frac{\|\mathbf{z}\|^2}{2}\right); \quad (13)$$

For a complex vector  $\mathbf{z} = (z_1, \dots, z_d)^T \in \mathbb{C}^d$ , we denote  $\mathbf{z}^2 = (\tau z_1^2, \dots, \tau z_d^2)^T$ . Consider  $\mathbf{z} = A\mathbf{x} + (A^T)^{-1}\mathbf{y}$  for an invertible (in  $\mathbb{C}^{d \times d}$ ) matrix  $A \in \mathbb{C}^{d \times d}$ . Using Lemma 2.4, we get:

$$\exp\left(\frac{(A\mathbf{x})^2}{2}\right) \exp\left(\frac{((A^T)^{-1}\mathbf{y})^2}{2}\right) \text{SM}(\mathbf{x}; \mathbf{y}) = E[\exp(\tau \langle A\mathbf{x} + (A^T)^{-1}\mathbf{y}, A\mathbf{x} + (A^T)^{-1}\mathbf{y} \rangle)] \quad (14)$$

Thus for  $\mathbf{u} \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m \exp\left(\frac{(M\mathbf{u}_i)^2}{2}\right) (\exp(\tau \langle \mathbf{u}_1, M\mathbf{u}_i \rangle); \dots; \exp(\tau \langle \mathbf{u}_m, M\mathbf{u}_i \rangle))^T$  and  $\mathbf{u}_i \in \mathbb{N}(0; I_d)$ :

$$\text{SM}(\mathbf{x}; \mathbf{y}) = E[\langle \frac{1}{m} \sum_{i=1}^m \mathbf{u}_i, \frac{1}{m} \sum_{i=1}^m \mathbf{u}_i \rangle]; \quad (15)$$

We obtain the new random feature map mechanism providing unbiased estimation of the softmax kernel. In general it is asymmetric since it applies different parameterization to  $x$  and  $y$ , i.e. one using  $A$ , the other one  $(A^T)^{-1}$ . Asymmetric random features is a simple generalization of the model using the same for both  $x$  and  $y$ , presented earlier. The resulting estimator  $\mathcal{SM}_m^{\text{cexp}}(x; y)$ , that we call complex exponential (CE) will serve as base estimators, and are defined as:

$$\mathcal{SM}_m^{\text{cexp}}(x; y) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m \langle Ax, (A^T)^{-1}y \rangle \quad (16)$$

For  $A = I_d$ , CE-estimator becomes trigonometric and for  $A = I_d$ , it becomes  $\mathcal{SM}_m^{++}(x; y)$ . Note also that an estimator based on this mechanism has variance equal to zero if

$$x = (A^{-1})(A^T)^{-1}y \quad (17)$$

In fact we can relax that condition. It is easy to check that for the variance to be zero, we only need:

$$\text{Re}(Ax) + \text{Re}((A^T)^{-1}y) = 0 \text{ and } \text{Im}(Ax) + \text{Im}((A^T)^{-1}y) = 0 \quad (18)$$

#### 2.4 ADAPTING HRF ESTIMATORS: HOW TO INSTANTIATE HYBRID RANDOM FEATURES

We will use complex exponential estimators from Section 2.3 as base estimators in different instantiations of HRF estimators (that by definition admit random feature map decomposition). The key to the linearization of the HRF estimators is then Equation 9 which implies linearization of the shifted versions of  $\phi$ -coefficients (Equation 10) and consequently - desired random feature map decomposition via the mechanism of compositional kernels (details in the Appendix). The questions remains how to construct  $\phi$ -coefficients to reduce variance of the estimation in the desired regions of interest.

Accurate approximation of small and large softmax kernel values simultaneously: As we explained earlier, in several applications of softmax estimators a desired property is to provide accurate estimation in two antipodal regions - small and large softmax kernel values. This can be done in particular by choosing  $\phi = 1$ ,  $\mathcal{SM}^1 = \mathcal{SM}^{++}$ ,  $\mathcal{SM}^2 = \mathcal{SM}^{\text{trig}}$  and  $(x; y) = \frac{x \cdot y}{\|x\| \|y\|}$ , where  $\frac{x \cdot y}{\|x\| \|y\|}$  stands for an angle between  $x$  and  $y$ . The key observation is that  $(x; y)$  can be unbiasedly approximated via the mechanism of random features (following from Goemans-Williamson algorithm (Goemans & Williamson, 2004)) as follows (details in the Appendix, Sec. D):

$$b(x; y) = \frac{1}{2} + \frac{1}{2} \sum_{i=1}^n \langle x, z_i \rangle \langle y, z_i \rangle \quad (19)$$

where  $z_{1:n}(z) = z_{2:n}(z) = z_n(z) \stackrel{\text{def}}{=} \frac{1}{\sqrt{2n}} (\text{sgn}(z_1); \dots; \text{sgn}(z_n))^\top$  for  $z_1; \dots; z_n \in \mathbb{N}(0; 1)$  and  $\|z\|^2 = 1$ . We call the resulting estimator of the softmax kernel angular hybrid. As we show in Sec. 3, this estimator is indeed particularly accurate in approximating both small and large softmax kernel values. For instance, for the prominent case when input vectors are of fixed length, its variance is zero for both the smallest ( $\frac{x \cdot y}{\|x\| \|y\|} = 1$ ) and largest ( $\frac{x \cdot y}{\|x\| \|y\|} = 0$ ) softmax kernel values (see: Fig. 1).

Gaussian Lambda Coefficients: Another tempting option is to instantiate coefficients with approximate Gaussian kernel values (estimated either via positive or trigonometric random features). In that setting, functions  $b^k(x; y)$  are defined as  $b^k(x; y) = \exp(-\frac{\|kx + M_k y\|_2^2}{2c_k^2})$ ; for some  $M_1; \dots; M_p \in \mathbb{R}^{d \times d}$  and  $c_1; \dots; c_p \in \mathbb{R}$ . Since this time coefficients are the values of the Gaussian kernel between vectors  $\frac{x}{c_k}$  and  $\frac{M_k y}{c_k}$ , we can take  $c_k = 0$  and define  $k_1; \dots; k_2$  as:

$$k_1(x) = \exp(-\frac{\|kx\|_2^2}{2c_k^2}) \mathcal{SM}(\frac{x}{c_k}) \text{ and } k_2(y) = \exp(-\frac{\|M_k y\|_2^2}{2c_k^2}) \mathcal{SM}(\frac{M_k y}{c_k}); \quad (20)$$

where  $\mathcal{SM}$  is the random feature map corresponding to a particular estimator of the softmax kernel. Note that the resulting hybrid estimator, that we call Gaussian hybrid, has nice pseudo-combinatorial properties. The coefficients  $k \in \mathbb{N}(0; 1]$  refer to  $\frac{\|M_k y\|_2}{\|y\|_2}$  (the ring window can be accurately controlled via hyperparameters). Thus matrices  $M_k$  can be coupled with the base estimators particularly accurate in the regions, where  $\frac{\|M_k y\|_2}{\|y\|_2} = 0$ . The mechanism can be thus trivially applied to the setting of accurate approximation of both small and large softmax kernel values for inputs of fixed length, yet it turns out to be characterized by the larger variance than its angular hybrid counterpart and cannot provide variance vanishing property for both  $\frac{x \cdot y}{\|x\| \|y\|} = 0$  and  $\frac{x \cdot y}{\|x\| \|y\|} = 1$ .

Figure 1: MSEs for three softmax kernel estimators (from left to right  $\mathcal{SM}^{\text{trig}}$ ,  $\mathcal{SM}^{++}$  and angular hybrid) for  $m = 10; n = 1$  and input of lengths  $s = 5$ . MSEs are given as functions of: an angle  $\theta \in [0; \pi]$  between  $x$  and  $y$  and  $r$  (symmetrized along  $\theta$  for length axis). For each plot, we marked in grey its slice for a fixed  $r$  to illustrate that only for the angular hybrid estimator, the MSE goes to zero for both  $\theta = 0$  and  $\theta = \pi$ .

Adaptation to data admitting clustering structure: Assume now that the inputs come from a distribution  $P(X)$  admitting certain clustering structure with clusters  $X_1; \dots; X_u \subset \mathbb{R}^d$  of centers  $x_1; \dots; x_u \in \mathbb{R}^d$  and that the analogous fact is true for the corresponding clusters  $Y_1; \dots; Y_w \subset \mathbb{R}^d$  and centers  $y_1; \dots; y_w \in \mathbb{R}^d$ . For a fixed pair of clusters' centers  $(x_i; y_j)$ , one can associate a complex exponential estimator  $\mathcal{SM}_{(x_i; y_j)}^{\text{cexp}}$  with the corresponding matrix  $A \in \mathbb{C}^{d \times d}$  (see: Sec. 2.3) that minimizes the following loss, where  $\|z\|_2$  for  $z \in \mathbb{C}^d$  is defined as  $\|z\|_2 = \sqrt{\sum_{i=1}^d |z_i|^2}$ :

$$l(A) = \|Ax_i + (A^T)^{-1}y_j\|_2^2 \quad (21)$$

Note that, by Equation 17, if one can find such that  $l(A) = 0$ , then  $\mathcal{SM}_{(x_i; y_j)}^{\text{cexp}}$  makes no error in approximating softmax kernel value on the centers  $x_i; y_j$ . Furthermore, if sampled points are close to the centers and  $l(A)$  is small, the corresponding variance of the estimator will also be small. Thus one can think about  $\mathcal{SM}_{(x_i; y_j)}^{\text{cexp}}$  as an estimator adapted to a pair of clusters  $(x_i; y_j)$ . We can add additional constraints regarding  $A$ , i.e.  $A \in \mathbb{R}^{d \times d}$  or  $A$  being diagonal (or both) for closed-form formulae of the optimal  $A$ , yet not necessarily zeroing  $l(A)$  (details in the Appendix, Sec. H.3). Complex exponential estimators defined in this way can be combined together, becoming base estimators in the HRF estimator. The corresponding coefficients can be chosen to be Gaussian as described in the previous paragraph (with optimized matrices so that a coefficient refers for the corresponding pair of clusters' centers), but there is another much simpler choice of index different base estimators by pairs of clusters' centers  $(x_i; y_j)$ . One can simply take  $\mathbb{1}_{(x_i; y_j)}(x) = 1$  if  $x$  belongs to the cluster of center  $x_i$  and  $\mathbb{1}_{(x_i; y_j)}(x) = 0$  otherwise. Vectors  $\mathbb{1}_{(x_i; y_j)}(y)$  (that also become scalars in this scenario) are defined analogously (the symbol vanishes since is deterministic). Thus effectively the hybrid estimator activates precisely this random feature mechanism that is optimal for the particular pair of clusters. This scheme is feasible if cluster-membership calculation is computationally acceptable, otherwise aforementioned Gaussian coefficients can be applied.

We conclude with an observation that in principle (to reduce the number of base estimators if needed) one can also construct HRF estimators that take base estimators defined in the above way only for the most massive pairs of clusters and add an additional default base estimators (for instance)

### 3 THEORETICAL GUARANTEES

We provide here theoretical guarantees regarding HRF estimators. We focus on the bipolar setting with two base estimators  $\mathcal{SM}^{\text{trig}}$  and  $\mathcal{SM}^{++}$ . The following is true:

Theorem 3.1 (MSE of the bipolar hybrid estimator) Take the bipolar hybrid estimator  $\mathcal{SM}_{m,n}^{\text{hyb}}(x; y)$ , where  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$  are chosen independently i.e. their random projections are chosen independently (note that we always assume  $b_n$  that) is constructed independently from  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$ ). Then the following holds:

$$\text{MSE}(\mathcal{SM}_{m,n}^{\text{hyb}}(x; y)) = E[b_n^2(x; y)]\text{MSE}(\mathcal{SM}_m^{++}(x; y)) + E[(1 - b_n(x; y))^2]\text{MSE}(\mathcal{SM}_m^{\text{trig}}(x; y)) \quad (22)$$

Furthermore, if  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$  apply the exact same sets of random projections, the mean squared error of the hybrid estimator is further reduced, namely we have:

$$\begin{aligned} \text{MSE}(\mathcal{SM}_{m,n}^{\text{hyb}}(x; y)) &= E[b_n^2(x; y)]\text{MSE}(\mathcal{SM}_m^{++}(x; y)) + E[(1 - b_n(x; y))^2]\text{MSE}(\mathcal{SM}_m^{\text{trig}}(x; y)) \\ &\quad - \frac{2}{m} \text{SM}^2(x; y) (1 - \cos(\kappa x_2^2 - \kappa y_2^2)) E[b_n(x; y)(1 - b_n(x; y))] \end{aligned} \quad (23)$$

The exact formula of  $\text{MSE}(\mathcal{SM}_m^{++}(x; y))$  and  $\text{MSE}(\mathcal{SM}_m^{\text{trig}}(x; y))$  is given in Lemma 2.2.

We can apply this result directly to the angular hybrid estimator, since:

Lemma 3.2. For the angular hybrid estimator the following holds:

$$E[b_n^2(x; y)] = \frac{x \cdot y}{n} - \frac{x \cdot y}{n} \frac{x \cdot y}{n} + \frac{1}{n} ; E[b_n(x; y)] = \frac{x \cdot y}{n} \quad (24)$$

We see that, as mentioned before, the variance of the angular hybrid estimator is zero for both  $x \cdot y = 0$  and  $x \cdot y = 1$  if inputs  $x; y$  have the same length. We need one more definition.

Definition 3.3. Assume that the inputs to the estimators are taken from some given bounded set  $C \subseteq \mathbb{R}^d$ . For a given estimator  $\mathcal{SM}$  on feature vectors  $x; y \in C$ , we define its max-relative-error with respect to  $C$  as  $\epsilon_C(\mathcal{SM}) = \max_{x; y \in C} \epsilon_{x; y}(\mathcal{SM})$ ; where  $\epsilon_{x; y}(\mathcal{SM}) = \frac{\text{MSE}(\mathcal{SM}(x; y))}{\text{SM}(x; y)}$ . Denote

by  $S(r)$  a sphere centered at  $\mathbf{0}$  and of radius  $r$ . Define  $\epsilon_{S(r)}(\mathcal{SM}) \stackrel{\text{def}}{=} \max_{x; y \in S(r)} \epsilon_{x; y}(\mathcal{SM})$  for  $x; y \in S(r)$  and such that  $\epsilon_{S(r)}(\mathcal{SM})$  (note that the mean squared errors of the considered estimators depend only on the angle  $x \cdot y$  for  $x; y$  chosen from a fixed sphere).

This definition captures the critical observation that in several applications of the softmax kernel estimation, e.g. efficient softmax sampling or linear-attention Transformers (Choromanski et al., 2021b), small relative errors are a much more meaningful measure of the quality of the method than small absolute errors. It also enables us to find hidden symmetries between different estimators:

Lemma 3.4. The following holds:

$$\epsilon_{S(r)}(\mathcal{SM}_m^{\text{trig}}) = \epsilon_{S(r)}(\mathcal{SM}_m^{++}) = \frac{1}{2m} \exp(2r^2 \sin^2(\frac{\pi}{2})) - \frac{1}{m} \exp(-4r^2 \sin^2(\frac{\pi}{2})) ; \quad (25)$$

and consequently  $\lim_{r \rightarrow 0} \epsilon_{S(r)}(\mathcal{SM}_m^{\text{trig}}) = \lim_{r \rightarrow 0} \epsilon_{S(r)}(\mathcal{SM}_m^{++}) = \frac{1}{2m} W(r) - \frac{1}{m} W(r) :$

$$\epsilon_{S(r)}(\mathcal{SM}_m^{\text{trig}}) = \epsilon_{S(r)}(\mathcal{SM}_m^{++}) = \lim_{r \rightarrow 0} \epsilon_{S(r)}(\mathcal{SM}_m^{\text{trig}}) = \lim_{r \rightarrow 0} \epsilon_{S(r)}(\mathcal{SM}_m^{++}) = \frac{1}{2m} W(r) - \frac{1}{m} W(r) \quad (26)$$

Our main result shows that HRF estimators can be applied to reduce the max-relative-error of the previously applied estimators of the softmax kernel. In the next theorem we show that the max-relative-error of the angular hybrid estimator scales as  $\frac{1}{m} \exp(2r^2)$  in the length of its inputs as opposed to  $\exp(2r^2)$  as it is the case for  $\mathcal{SM}_m^{\text{trig}}$  and  $\mathcal{SM}_m^{++}$  (see: Lemma 3.4). Furthermore, the max-relative-error scales as  $\frac{1}{m}$  and  $\frac{1}{mn}$  as  $r \rightarrow 0$  and  $r \rightarrow 1$  respectively, in particular goes to 0 in both critical cases. This is not true for  $\mathcal{SM}_m^{\text{trig}}$  nor for  $\mathcal{SM}_m^{++}$ .

Theorem 3.5. The max-relative-error of the angular hybrid estimator for the inputs on the sphere  $S(r)$  of radius  $r \leq 1$  satisfies for  $W(r) = \exp(2r^2) - \exp(-4r^2) :$

$$\epsilon_{S(r)}(\mathcal{SM}_{m,n}^{\text{anghyb}}) = \frac{1}{r} \frac{1}{2m} W(r) - \frac{1}{n} \frac{1}{n} + \frac{1}{n} \frac{1}{n} = \quad (27)$$

$$\text{Furthermore, } \lim_{r \rightarrow 0} \epsilon_{S(r)}(\mathcal{SM}_{m,n}^{\text{anghyb}}) = \lim_{r \rightarrow 0} \epsilon_{S(r)}(\mathcal{SM}_{m,n}^{\text{anghyb}}) = \frac{1}{2mn} W(r) :$$

Additional implications of Theorem 3.5, using the fact that  $(atmn)$ -dimensional HRFs can be constructed in time  $\mathcal{O}(nd + md + mn)$  (regular RFs need  $\mathcal{O}(mnd)$ ), are in the Appendix (Sec. D.4).

## 4 EXPERIMENTS

In this section we conduct exhaustive evaluation of the mechanism of HRFs on several tasks.

#### 4.1 POINTWISE SOFTMAX KERNEL ESTIMATION

We start with ablation studies regarding empirical relative errors of different softmax kernel estimators across different inputs' lengths and angles  $x, y$ . The results are presented in Fig. 2. Ablations over more lengths are presented in the Appendix (Sec. G). The angular hybrid estimator most accurately approximates softmax kernel and has smaller max-relative-error than other estimators.

Figure 2: Pointwise estimation of  $\mathcal{SM}(x; y)$  for the same-length 64-dim inputs ( $\ell = 1 : 0$  and  $r = 1 : 5$ ) and various angles  $x, y$ . Red-dotted lines are for marking zero-level. We used 10000 estimated softmax values in each subplot. The true value and 5th and 95th quantile estimated values are shown by the left y-axis, and the empirical relative errors are shown by the right y-axis. Trigonometric estimator and FAVOR+ applied 128 random features. To make fair comparison, for the hybrid variant the configuration leading to the similar number of FLOPS operations per random feature map creation was applied. Similar gains as for the angular are obtained by the Gaussian hybrid variant.

Comparison with QMC-methods: Even though, as we explained in Section 1, our algorithm does not compete with various methods for variance reduction based on different sampling techniques (since those methods, as orthogonal to ours, can be easily incorporated into our framework), we decided to compare HRFs also with them. We tested angular hybrid HRF variant as well as well-established QMC baselines: orthogonal random features (ORF)(Yu et al., 2016) (regular variant ORF-reg as well as the one applying Hadamard matrices ORF-Had) and QMCs based on random Halton sequences (Halton-R)(Avron et al., 2016). For HRFs, we tested two sub-variants: with orthogonal (HRF-ort) and regular iid random features (HRF-iid). We computed empirical MSEs by averaging over 100 randomly sampled pairs of vectors from two UCI datasets: wine and Boston. HRFs outperform all other methods, as we see in Table 1.

Table 1: Comparison of different estimators of the softmax kernel on the datapoints from two UCI datasets: wine and Boston in terms of the MSE (measured in  $10^{-3}$  units). The non-HRF estimators apply 512 random features and the HRF-ones are set up to match the non-HRFs in terms of the number of FLOPS (for a fair comparison). We also reported standard deviations.

Datasets	HRF-ort		HRF-iid		ORF-reg		ORF-Had		Halton-R	
Wine	0:70	0:08	0:85	0:05	1:00	0:04	1:10	0:06	4:02	0:1
Boston	0:72	0:06	0:79	0:08	1:05	0:03	1:14	0:02	2:53	0:08

#### 4.2 LANGUAGE MODELING

In this section, we apply HRFs to the language modeling task, training a 2-layer LSTM with hidden size  $h = 200$  and applying RFs for softmax sampling in training as described in (Rawat et al., 2019). Experimental results are obtained over 100 runs. Experimental details and additional validation results are in the Appendix (Sec. I, Table 4).

Figure 3: Statistical metrics measuring softmax matrix approximation quality on PennTree Bank. For standard estimators, the number of random features are 64, 128, 256, 512. To make fair comparison, for the hybrid variants, the configurations leading to the similar number of FLOPS operations per random feature map creation were applied. Negative fractions were not reported for FAVOR+ since by definition they are equal to zero.

Statistical Metrics. We computed d-Wasserstein distance and the Kolmogorov-Smirnov (KS) metric (Ramdas et al., 2017) between a true softmax distribution induced by a query on all the classes

and the approximate one given by different RF-mechanisms. We also compute the fraction of all probability estimates with undesired negative values. The results on the PennTree Bank Marcus et al. (1993) are presented in Fig. 3. In the Appendix (Sec. I) we present additional results for the WikiText2 dataset. We see that HRFs lead to most accurate softmax distribution estimation.

#### 4.3 TRAINING SPEECH MODELS WITH HRF-CONFORMERS PERFORMERS

We also tested HRFs on speech models with LibriSpeech ASR corpus (Panayotov et al., 2015). We put implicit Performer’s attention into 7-layer Conformer-Transducer encoder (Gulati et al., 2020) and compared softmax kernel estimators in terms of word error rate (WER) metric, commonly used to evaluate speech models. We tested HRFs using angular hybrid variant as well as those clustering-based. For the latter ones, the clusters were created according to k-means algorithm after first 1000 steps of training (and were frozen afterwards) and corresponding matrices are constructed to minimize the loss given in Equation 21. The results are presented in Table 2. HRFs produce smallest WER models and the clustering-based variants fully exercising the most general formulation on HRFs (with general complex estimators) turn out to be the best. Additional details regarding speech experiments as well as additional experiments with clustering-based HRFs are presented in the Appendix (Sec. J and Sec. H.1 respectively).

Table 2: Comparison of WERs of Conformer-Transducer applying different RF-mechanisms for the implicit attention. For methods other than clustering-based HRFs (HRF-C), numbers next to method names denote the values of  $m$  or  $(m; n)$ . Method HRF-A stands for the angular hybrid variant. Numbers next to HRF-C correspond to the number of clusters constructed in the query and key space respectively. HRF-C uses random features. We also report standard deviations averaged over different training runs.

	HRF-C(3,3)	HRF-C(2, 3)	HRF-C(3,2)	HRF-C(2,2)	HRF-A(16,8)
WER	1:72 0:02%	1:75 0:03%	1:83 0:03%	1:85 0:04%	2:03 0:08%
	HRF-A(8, 8)	FAVOR+ 432	FAVOR+ 256	Trig 432	Trig 256
WER	2:05 0:05%	2:65 0:06%	2:77 0:04%	3:12 0:05%	3:3 0:06%

#### 4.4 DOWNSTREAM ROBOTICS EXPERIMENTS

In Robotics, we leverage the accuracy of HRFs to obtain inference improvements, critical for on-robot deployment. To abstract from a particular hardware characteristic, we measure it in the number of FLOPS. We conduct two experiments targeting: (a) quadruped locomotion and (b) robotic-arm manipulation. In the former, HRFs are applied as a replacement of the default mechanism using positive random features in the class of implicit-attention architectures for vision processing called Implicit Attention Policies (or IAPs) (Choromanski et al., 2021a). In the latter, HRFs become a part of the regular Vision-Performer stack used to process high-resolution (500 x 500 x 3) input. The results are presented in Fig. 4, where HRFs need fewer FLOPS to obtain same quality policies as baselines. All details regarding both experimental setups are given in the Appendix (Sec. K). Videos of the HRF-trained robotic policies are included in the supplementary material.

Figure 4: Left: Step-stone locomotion task. Comparison of the training curves for: the IAP using angular hybrid estimator of  $m = n = 8$  and IAP applying regular FAVOR+ mechanism from (Choromanski et al., 2021b) with  $m = 256$ . Both policies are of similar quality, yet HRF-method requires fewer FLOPS to run its trained policy. The visualization of the HRF policy in action and its attention (with iterated out pixels marked in red) is in the bottom right corner. Right: Similar setting (and conclusions) but for the robotic-arm manipulation task. The additional regular RF-conformations did not train by producing large loss due to large variance of the underlying softmax kernel estimators.

## 5 CONCLUSION

We presented a new class of random feature techniques, called random features for softmax/Gaussian kernel estimation that more accurately approximate these kernels than previous algorithms. We also demonstrated their robustness in a wide range of applications from softmax sampling to Transformers/attention training, also for downstream Robotics tasks.

Reproducibility Statement: Section 4 and the Appendix include all experimental details (e.g. hyperparameter setup) needed to reproduce all the results presented in the paper. The part of the code that we could make publicly available can be found in the following git repository: [https://github.com/HL-hanlin/HRF\\_ICLR2022](https://github.com/HL-hanlin/HRF_ICLR2022).

Ethics Statement: HRFs can be used in principle to train massive Transformer models with large number of parameters and proportional compute resources. Thus they should be used responsibly, in particular given CO<sub>2</sub> emission challenges that scientific community tries to address now.

## ACKNOWLEDGEMENTS

AW acknowledges support from a Turing AI Fellowship under grant EP/V025379/1, The Alan Turing Institute, and the Leverhulme Trust via CFI. AS acknowledges AWS compute resources from Onur Kara.

## REFERENCES

- Unitree Robotics. URL <http://www.unitree.cc/>.
- Haim Avron, Vikas Sindhwani, Jiyan Yang, and Michael W. Mahoney. Quasi-monte carlo feature maps for shift-invariant kernels. *J. Mach. Learn. Res.* 17:120:1–120:38, 2016. URL <http://jmlr.org/papers/v17/14-538.html>.
- Guy Blanc and Steffen Rendle. Adaptive sampled softmax with kernel based sampling. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 589–598. PMLR, 2018. URL <http://proceedings.mlr.press/v80/blanc18a.html>.
- Y. Cho and L. K. Saul. Analysis and extension of arc-cosine kernels for large margin classification. Technical Report CS2012-0972, Department of Computer Science and Engineering, University of California, San Diego., 2012.
- Krzysztof Choromanski and Vikas Sindhwani. Recycling randomness with structure for sublinear time kernel expansions. In Maria-Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pp. 2502–2510. JMLR.org, 2016. URL <http://proceedings.mlr.press/v48/choromanski16.html>.
- Krzysztof Choromanski, Mark Rowland, Wenyu Chen, and Adrian Weller. Unifying orthogonal Monte Carlo methods. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1203–1212. PMLR, 2019. URL <http://proceedings.mlr.press/v97/choromanski19a.html>.
- Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, David Belanger, Lucy Colwell, and Adrian Weller. Masked language modeling for proteins via linearly scalable long-context transformers. preprint arXiv:2006.03555, 2020.
- Krzysztof Choromanski, Deepali Jain, Jack Parker-Holder, Xingyou Song, Valerii Likhoshesterov, Anirban Santara, Aldo Pacchiano, Yunhao Tang, and Adrian Weller. Unlocking pixels for reinforcement learning via implicit attention. *CoRR*, abs/2102.04353, 2021a. URL <https://arxiv.org/abs/2102.04353>.
- Krzysztof Marcin Choromanski, Mark Rowland, and Adrian Weller. The unreasonable effectiveness of structured random orthogonal embeddings. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, 4-9 December 2017, Long Beach, CA, USA, pp. 219–228, 2017.

- Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Taras Sadowski, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin Belanger, Lucy J. Colwell, and Adrian Weller. Rethinking attention with performers. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021b. URL <https://openreview.net/forum?id=Ua6zuk0WRH>.
- Amit Daniely, Roy Frostig, Vineet Gupta, and Yoram Singer. Random features for compositional kernels. CoRR, abs/1703.07872, 2017. URL <http://arxiv.org/abs/1703.07872>.
- Michel X. Goemans and David P. Williamson. Approximation algorithms for MAX-3-CUT and other problems via complex semidefinite programming. *Comput. Syst. Sci.* 68(2):442–470, 2004. doi: 10.1016/j.jcss.2003.07.012. URL <https://doi.org/10.1016/j.jcss.2003.07.012>.
- Arthur Gretton, Ralf Herbrich, Alexander J. Smola, Olivier Bousquet, and Bernhard Schölkopf. Kernel methods for measuring independence. *Mach. Learn. Res.* 6:2075–2129, 2005. URL <http://jmlr.org/papers/v6/gretton05a.html>.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. Conformer: Convolution-augmented transformer for speech recognition. In Helen Meng, Bo Xu, and Thomas Fang Zheng (eds.), *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020*, pp. 5036–5040. ISCA, 2020. doi: 10.21437/Interspeech.2020-3015. URL <https://doi.org/10.21437/Interspeech.2020-3015>.
- Hakan Inan, Khashayar Khosravi, and Richard Socher. Tying word vectors and word classifiers: A loss framework for language modeling. *Proceedings of the 5th International Conference on Learning Representations*, 2017.
- Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan A. K. Suykens. Random features for kernel approximation: A survey in algorithms, theory, and beyond. CoRR, abs/2004.11154, 2020. URL <https://arxiv.org/abs/2004.11154>.
- Shengjie Luo, Shanda Li, Tianle Cai, Di He, Dinglan Peng, Shuxin Zheng, Guolin Ke, Liwei Wang, and Tie-Yan Liu. Stable, fast and accurate: Kernelized attention with relative positional encoding. *Advances in Neural Information Processing Systems*, 2021.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics* 19(2):313–330, 1993. URL <https://aclanthology.org/J93-2004>.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer Sentinel Mixture Models. *5th International Conference on Learning Representations*, abs/1609.07843, 2017.
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An ASR corpus based on public domain audio books. *2015 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2015, South Brisbane, Queensland, Australia, April 19-24, 2015*, pp. 5206–5210. IEEE, 2015. doi: 10.1109/ICASSP.2015.7178964. URL <https://doi.org/10.1109/ICASSP.2015.7178964>.
- Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A. Smith, and Lingpeng Kong. Random feature attention. *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=QtTKTdVrFBB>.
- Orr Press and Lior Wolf. Using the Output Embedding to Improve Language Models, 2017.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In John C. Platt, Daphne Koller, Yoram Singer, and Sam T. Roweis (eds.), *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*

- pp. 1177–1184. Curran Associates, Inc., 2007. URL <http://papers.nips.cc/paper/3182-random-features-for-large-scale-kernel-machines>
- Aaditya Ramdas, Nicolás García Trillos, and Marco Cuturi. On Wasserstein two-sample testing and related families of nonparametric tests. *Entropy*, 19(2):47, 2017. doi: 10.3390/e19020047. URL <https://doi.org/10.3390/e19020047>
- Ankit Singh Rawat, Jiecao Chen, Felix X. Yu, Ananda Theertha Suresh, and Sanjiv Kumar. Sampled softmax with random Fourier features. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Aubert-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 13834–13844, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/e43739bba7cdb577e9e3e4e42447f5a5-Abstract.html>
- Mark Rowland, Krzysztof Choromanski, François Chalus, Aldo Pacchiano, Afonso S. Barros, Richard E. Turner, and Adrian Weller. Geometrically coupled Monte Carlo sampling. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Niccolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 195–205, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/b3e3e393c77e35a4a3f3cbd1e429b5dc-Abstract.html>
- Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Less is more on noncomputational regularization. In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp. 1657–1665, 2015. URL <https://proceedings.neurips.cc/paper/2015/hash/03e0704b5690a2dee1861dc3ad3316c9-Abstract.html>
- Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9355–9366. PMLR, 2021. URL <http://proceedings.mlr.press/v139/schlag21a.html>
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018. URL <http://arxiv.org/abs/1807.03748>
- Christopher K. I. Williams and Matthias W. Seeger. Using the Nyström method to speed up kernel machines. In Todd K. Leen, Thomas G. Dietterich, and Volker Tresp (eds.), *Advances in Neural Information Processing Systems 13, Papers from Neural Information Processing Systems (NIPS) 2000, Denver, CO, USA*, pp. 682–688. MIT Press, 2000. URL <https://proceedings.neurips.cc/paper/2000/hash/19de10adbaa1b2ee13f77f679fa1483a-Abstract.html>
- Tianbao Yang, Yu-Feng Li, Mehrdad Mahdavi, Rong Jin, and Zhi-Hua Zhou. Nyström method vs random Fourier features: A theoretical and empirical comparison. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C. Burges, Léon Bottou, and Kilian Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States*, pp. 485–493, 2012. URL <https://proceedings.neurips.cc/paper/2012/hash/621bf66ddb7c962aa0d22ac97d69b793-Abstract.html>
- Felix X. Yu, Ananda Theertha Suresh, Krzysztof Marcin Choromanski, Daniel N. Holtmann-Rice, and Sanjiv Kumar. Orthogonal random features. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 1975–1983, 2016. URL <http://papers.nips.cc/paper/6246-orthogonal-random-features>

Qinyi Zhang, Sarah Filippi, Arthur Gretton, and Dino Sejdinovic. Large-scale kernel methods for independence testing. *Stat. Comput.* 28(1):113–130, 2018. doi: 10.1007/s11222-016-9721-7. URL <https://doi.org/10.1007/s11222-016-9721-7>.

## APPENDIX: HYBRID RANDOM FEATURES

## A PROOF OF LEMMA 2.4

Proof. From the independence of  $z_i$ , we get:

$$E[\exp(\sum_{i=1}^d z_i)] = \prod_{i=1}^d E[\exp(z_i)] \quad (28)$$

Thus it suffices to show that for  $z \sim \mathcal{N}(0, 1)$  and any  $z \in \mathbb{C}$  the following holds:

$$E[\exp(gz)] = \exp\left(\frac{z^2}{2}\right) \quad (29)$$

Define  $f(z) = E[\exp(gz)]$ . Note that  $f(z) = \exp\left(\frac{z^2}{2}\right)$  for  $z = ix$ , where  $i^2 = -1$  and  $x \in \mathbb{R}$  which follows from the formula of the characteristic function of the Gaussian distribution. Similarly,  $f(z) = \exp\left(\frac{z^2}{2}\right)$  for  $z \in \mathbb{R}$  which follows from the formula of the moment generating function of the Gaussian distribution. We now use the fact from complex analysis that if two analytic functions  $f, g: \mathbb{C} \rightarrow \mathbb{C}$  are identical on uncountably many points then they are equal. To complete the proof (we leave to the reader checking that  $f(z)$  and  $g(z) \stackrel{\text{def}}{=} \exp\left(\frac{z^2}{2}\right)$  are analytic).  $\square$

## B PROOF OF LEMMA 2.3

Proof. The following is true:

$$\begin{aligned} \mathbb{E} \mathcal{M}_{m;n}^{\text{hyb}}(x; y) &= \prod_{k=1}^p \frac{a_k}{m} \left( \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \right) \left( \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \right) + \\ &\quad \frac{1}{mn} \prod_{k=1}^p \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \left( \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \right) + \\ &\quad \frac{1}{m} \prod_{k=1}^p a_k \left( \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \right) \left( \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \right) + \\ &\quad \frac{1}{mn} \prod_{k=1}^p \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \left( \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \right) + \\ &\quad \frac{1}{mn} \prod_{k=1}^p \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \left( \mathbb{E} \left[ \prod_{i=1}^{t_{k;1}} z_{i;1} \right] \cdots \mathbb{E} \left[ \prod_{i=1}^{t_{k;t_k}} z_{i;t_k} \right] \right) + \end{aligned} \quad (30)$$

The statement of the lemma follows directly from that. Note that the key trick is an observation that a product of two functions of  $x, y$  that can be linearized on expectation (i.e. rewritten with  $x$  and  $y$  disentangled) is itself a function that can be linearized on expectation. In the setting where the former are approximated by random features, the latter is approximated by their cartesian product (the random feature mechanism for the compositional kernels from (Daniely et al., 2017))  $\square$

## C BIPOLAR HRF ESTIMATORS: PROOF OF THEOREM 3.1

Proof. The following holds:

$$\begin{aligned} \text{Var}(\mathbb{E} \mathcal{M}_{m;n}^{\text{hyb}}(x; y)) &= \text{Var} \left( \mathbb{E} \left[ \mathcal{M}_m^{++}(x; y) \right] \right) + \text{Var} \left( \mathbb{E} \left[ \mathcal{M}_m^{\text{trig}}(x; y) \right] \right) + \\ &\quad 2 \text{Cov} \left( \mathbb{E} \left[ \mathcal{M}_m^{++}(x; y) \right]; \mathbb{E} \left[ \mathcal{M}_m^{\text{trig}}(x; y) \right] \right) \end{aligned} \quad (31)$$

We will focus now on the covariance term. To simplify the notation, we will drop index  $b_n$ .

Assume first easier to analyze case where  $\mathcal{M}_m^{\text{trig}}(x; y)$  and  $\mathcal{M}_m^{++}(x; y)$  are chosen independently. Then the following is true:

$$\begin{aligned}
& \text{Cov}(b(\cdot) \mathcal{SM}_m^{++}(x; y); (1 - b(\cdot)) \mathcal{SM}_m^{\text{trig}}(x; y)) = \\
& E[b(\cdot)(1 - b(\cdot)) \mathcal{SM}_m^{++}(x; y) \mathcal{SM}_m^{\text{trig}}(x; y)] - E[b(\cdot) \mathcal{SM}_m^{++}(x; y)] E[(1 - b(\cdot)) \mathcal{SM}_m^{\text{trig}}(x; y)] \\
& = E[b(\cdot)(1 - b(\cdot))] E[\mathcal{SM}_m^{++}(x; y)] E[\mathcal{SM}_m^{\text{trig}}(x; y)] \\
& = (SM(x; y))^2 \text{Var}(b(\cdot)) \tag{32}
\end{aligned}$$

Now assume that  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$  use the exact same random projections. Then, using similar analysis as before, we get:

$$\begin{aligned}
& \text{Cov}(b(\cdot) \mathcal{SM}_m^{++}(x; y); (1 - b(\cdot)) \mathcal{SM}_m^{\text{trig}}(x; y)) = \\
& E[\mathcal{SM}_m^{++}(x; y) \mathcal{SM}_m^{\text{trig}}(x; y)] E[b(\cdot)(1 - b(\cdot))] - (SM(x; y))^2 E[b(\cdot)] E[(1 - b(\cdot))] \tag{33}
\end{aligned}$$

This time however it is no longer the case that  $E[\mathcal{SM}_m^{++}(x; y) \mathcal{SM}_m^{\text{trig}}(x; y)] E[b(\cdot)(1 - b(\cdot))] = E[\mathcal{SM}_m^{++}(x; y)] E[\mathcal{SM}_m^{\text{trig}}(x; y)] = (SM(x; y))^2$  since  $\mathcal{SM}_m^{++}(x; y)$  and  $\mathcal{SM}_m^{\text{trig}}(x; y)$  are no longer independent. In order to compute  $E[\mathcal{SM}_m^{++}(x; y) \mathcal{SM}_m^{\text{trig}}(x; y)] E[b(\cdot)(1 - b(\cdot))]$ , we will first introduce useful denotation.

Denote by  $\mathbf{!}_1, \dots, \mathbf{!}_m \stackrel{\text{iid}}{\sim} \mathcal{N}(0; I_d)$  the random projections sampled to construct  $\mathcal{SM}_m^{++}(x; y)$  and  $\mathcal{SM}_m^{\text{trig}}(x; y)$ . Denote  $Y_i = \cosh(\mathbf{!}_i^\top(x + y)) \stackrel{\text{def}}{=} \frac{\exp(\mathbf{!}_i^\top(x + y)) + \exp(-\mathbf{!}_i^\top(x + y))}{2}$ . We have:

$$\mathcal{SM}_m^{++}(x; y) = \exp\left(\frac{kxk^2 + kyk^2}{2} \frac{Y_1 + \dots + Y_m}{m}\right) \tag{34}$$

If we denote  $Z_i = \cos(\mathbf{!}_i^\top(x - y))$ , then we can write  $\mathcal{SM}_m^{\text{trig}}(x; y)$  as:

$$\mathcal{SM}_m^{\text{trig}}(x; y) = \exp\left(\frac{kxk^2 + kyk^2}{2} \frac{Z_1 + \dots + Z_m}{m}\right) \tag{35}$$

We can then rewrite  $E[\mathcal{SM}_m^{++}(x; y) \mathcal{SM}_m^{\text{trig}}(x; y)]$  as:

$$E[\mathcal{SM}_m^{++}(x; y) \mathcal{SM}_m^{\text{trig}}(x; y)] = \frac{1}{m^2} \prod_{i \in J} E[Y_i Z_i] + \frac{X^n}{i=1} E[Y_i Z_i]^5 = \tag{36}$$

$$\frac{1}{m^2} \prod_{i=1}^m (SM(x; y))^2 + m E[\cosh(\mathbf{!}_i^\top(x + y)) \cos(\mathbf{!}_i^\top(x - y))] ;$$

where  $\mathbf{!}_i \sim \mathcal{N}(0; I_d)$ . The equality follows from the unbiasedness of  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$  and the fact that different  $\mathbf{!}_i$  are chosen independently. Thus it remains to compute  $E[\cosh(\mathbf{!}_i^\top(x + y)) \cos(\mathbf{!}_i^\top(x - y))]$ . Note first that  $E[\cosh(\mathbf{!}_i^\top(x + y)) \cos(\mathbf{!}_i^\top(x - y))] = E[\exp(\mathbf{!}_i^\top(x + y)) \cos(\mathbf{!}_i^\top(x - y))]$  since  $\mathbf{!}_i \sim \mathcal{N}(0; I_d)$  and  $\cosh$  is an even function. Denote  $z = x + y + i(x - y)$ . We have:

$$\begin{aligned}
& E[\exp(\mathbf{!}_i^\top(x + y)) \cos(\mathbf{!}_i^\top(x - y))] = \text{Re} E[\exp(\mathbf{!}_i^\top z)] = \text{Re} \left[ \prod_{i=1}^d \exp\left(\frac{Z_i^2}{2}\right) \right] = \\
& \text{Re} \exp\left(\frac{\prod_{j=1}^d (x_j + y_j)^2 + 2i(x_j^2 - y_j^2) \prod_{j=1}^{\neq j} (x_j - y_j)^2}{2}\right) = (SM(x; y))^2 \cos(kxk_2^2 - kyk_2^2) \tag{37}
\end{aligned}$$

Thus we conclude that:

$$E[\mathcal{SM}_m^{++}(x; y) \mathcal{SM}_m^{\text{trig}}(x; y)] = (1 - \frac{1}{m})(SM(x; y))^2 + \frac{1}{m}(SM(x; y))^2 \cos(kxk_2^2 - kyk_2^2) \tag{38}$$

Therefore we get the formulae for the covariance term in both: the setting where random projections of  $\mathbf{SM}_m^{\text{trig}}(x; y)$  and  $\mathbf{SM}_m^{++}(x; y)$  are shared (variant II) and when they are not (variant I). The following is true for  $Z = \frac{1}{m}(1 - \cos(\|x\|_2 \|y\|_2))$ :  $E[b(\cdot)](1 - b(\cdot))$ :

$$\text{Cov}(b(\cdot)\mathbf{SM}_m^{++}(x; y); (1 - b(\cdot))\mathbf{SM}_m^{\text{trig}}(x; y)) = \begin{cases} (\text{SM}(x; y))^2 \text{Var}(b(\cdot)) & \text{for variant I} \\ (\text{SM}(x; y))^2 (\text{Var}(b(\cdot)) + Z) & \text{for variant II} \end{cases} \quad (39)$$

We also have the following:

$$\begin{aligned} \text{Var}(b(\cdot)\mathbf{SM}_m^{++}(x; y)) &= E[(b(\cdot))^2 (\mathbf{SM}_m^{++}(x; y))^2] - (E[b(\cdot)\mathbf{SM}_m^{++}(x; y)])^2 = \\ &= E[(b(\cdot))^2] \text{MSE}(\mathbf{SM}_m^{++}(x; y)) + (\text{SM}(x; y))^2 - (E[b(\cdot)])^2 (\text{SM}(x; y))^2 = \\ &= (\text{SM}(x; y))^2 \text{Var}(b(\cdot)) + E[(b(\cdot))^2] \text{MSE}(\mathbf{SM}_m^{++}(x; y)) \end{aligned} \quad (40)$$

and furthermore (by the analogous analysis):

$$\text{Var}((1 - b(\cdot))\mathbf{SM}_m^{\text{trig}}(x; y)) = (\text{SM}(x; y))^2 \text{Var}(b(\cdot)) + E[(1 - b(\cdot))^2] \text{MSE}(\mathbf{SM}_m^{\text{trig}}(x; y)) \quad (41)$$

By putting the derived formula for the above variance terms as well as covariance terms back in the Equation 31, we complete the proof of the theorem (note that the mean squared error of the hybrid estimator is its variance since it is unbiased).

□

## D ANGULAR HYBRID ESTIMATORS

### D.1 PROOF OF LEMMA 3.2

Proof. The formula for the expectation is directly implied by the formula (1) from Sec. 2.1 of Cho & Saul (2012) for the zeroth-order arc-cosine kernel  $k(x; y) \stackrel{\text{def}}{=} 1 - \frac{x \cdot y}{\|x\|_2 \|y\|_2}$ . It also follows directly from Goemans-Williamson algorithm (Goemans & Williamson, 2004). We will now derive the formula for  $c = E[b^2(x; y)]$ . Since  $b(\cdot)$  depends only on the angle  $\angle(x, y)$ , we will refer to it as:  $b(\cdot)$  and to its estimator as  $\hat{b}(\cdot)$ . Denote:

$$b_n^{\text{ang}}(z) = \frac{1}{n} (\text{sgn}(\langle z, e_1 \rangle); \dots; \text{sgn}(\langle z, e_n \rangle)) \quad (42)$$

Denote:  $X_i = (b_n^{\text{ang}}(x))[i] - b_n^{\text{ang}}(y)[i]$ . We have:

$$b(\cdot) = \frac{1}{2} \left( 1 + \prod_{i=1}^n X_i \right) \quad (43)$$

Note first that by the construction of  $b(\cdot)$ , we have:  $E[b(\cdot)] = 0$  and thus:  $E[\prod_{i=1}^n X_i] = 1 - \frac{2}{n}$ . Therefore we conclude that:

$$\begin{aligned} c &= \frac{1}{4} E \left[ \left( 1 + \prod_{i=1}^n X_i \right)^2 \right] = \frac{1}{4} \left( 1 + 2 \sum_{i=1}^n E[X_i^2] + \sum_{i \neq j} E[X_i^2] E[X_j^2] \right) \\ &= \frac{1}{4} \left( 1 + 2 \left( 1 - \frac{2}{n} \right) + n \frac{1}{n^2} + n(n-1) \frac{1}{n^2} \left( 1 - \frac{2}{n} \right)^2 \right) = \frac{1}{4} \left( 4 - \frac{4}{n} + \frac{1}{n} \left( 1 - \frac{2}{n} \right)^2 \right) \\ &= \frac{1}{4} \left( 4 - \frac{4}{n} + \frac{1}{n} \left( 1 - \frac{4}{n} + \frac{4}{n^2} \right) \right) = \frac{1}{4} \left( 4 - \frac{4}{n} + \frac{1}{n} - \frac{4}{n^2} + \frac{4}{n^3} \right) \end{aligned} \quad (44)$$

□

## D.2 EXPLICIT FORMULA FOR THE MSE OF THE ANGULAR HYBRID ESTIMATOR

Theorem 3.1 and Lemma 3.2 immediately imply the following result.

**Theorem D.1 (MSE of the angular hybrid estimator)** Take the angular hybrid estimator  $\mathbf{SM}_{m;n}^{\text{anghyb}}(x; y)$ , where  $\mathbf{SM}_m^{\text{trig}}(x; y)$  and  $\mathbf{SM}_m^{++}(x; y)$  are chosen independently i.e. their random projections are chosen independently (note that we always assume  $\mathbf{h}(z)$  is chosen independently from  $\mathbf{SM}_m^{\text{trig}}(x; y)$  and  $\mathbf{SM}_m^{++}(x; y)$ ). Then the following holds:

$$\begin{aligned} \text{MSE}(\mathbf{SM}_{m;n}^{\text{anghyb}}(x; y)) &= \frac{1}{n} + \frac{1}{n} \frac{1}{2m} \exp(kz k^2) \text{SM}^2(x; y) (1 - \exp(-kz k^2))^2 + \\ &= \frac{1}{n} + \frac{1}{n} \frac{1}{2m} \exp(kz k^2) \text{SM}^2(x; y) (1 - \exp(-kz k^2))^2 \end{aligned} \quad (45)$$

for  $z = x \cdot y$  and  $k = \|x + y\|$ . Furthermore, if  $\mathbf{SM}_m^{\text{trig}}(x; y)$  and  $\mathbf{SM}_m^{++}(x; y)$  apply the exact same sets of random projections, the mean squared error of the hybrid estimator is further reduced, namely we have:

$$\begin{aligned} \text{MSE}(\mathbf{SM}_{m;n}^{\text{anghyb}}(x; y)) &= \frac{1}{n} + \frac{1}{n} \frac{1}{2m} \exp(kz k^2) \text{SM}^2(x; y) (1 - \exp(-kz k^2))^2 + \\ &= \frac{1}{n} + \frac{1}{n} \frac{1}{2m} \exp(kz k^2) \text{SM}^2(x; y) (1 - \exp(-kz k^2))^2 \\ &\quad - \frac{2}{m} \text{SM}^2(x; y) (1 - \cos(kx k_2 \cdot ky k_2)) - \frac{1}{n} - \frac{1}{n} \end{aligned} \quad (46)$$

Thus if  $kx k_2 = ky k_2 = r$  then regardless of whether the same sets of random projections are applied or not, we get:

$$\begin{aligned} \text{MSE}(\mathbf{SM}_{m;n}^{\text{anghyb}}(x; y)) &= \frac{1}{n} + \frac{1}{n} \frac{1}{2m} \exp(8r^2 \cos^2(\frac{\alpha}{2})) - 2r^2 \\ &= \frac{1}{n} + \frac{1}{n} \frac{1}{2m} \exp(2r^2) (1 - \exp(-4r^2 \sin^2(\frac{\alpha}{2})))^2 \end{aligned} \quad (47)$$

## D.3 PROOF OF THEOREM 3.5

**Proof.** Note first that from the derived above formula of the MSE of the lambda-angular bipolar hybrid estimator and the definition of the max-relative-error, we obtain:

$$s_{(r)}(\mathbf{SM}_{m;n}^{\text{anghyb}}) = \frac{\exp(r^2)}{2m} \frac{1}{\max_{\alpha \in [0; \pi]} h_r(\alpha)} \quad (48)$$

where:

$$h_r(\alpha) = a_r(\alpha) + a_r(\alpha) \quad (49)$$

and  $a_r(\alpha)$  is defined as:

$$a_r(\alpha) = -\left(-\frac{1}{n} + \frac{1}{n}\right) \exp(2r^2 \cos(\alpha)) - \frac{1}{n} \exp(-4r^2 \cos^2(\frac{\alpha}{2}))^2 \quad (50)$$

Therefore we have:

$$s_{(r)}(\mathbf{SM}_{m;n}^{\text{anghyb}}) = \frac{\exp(r^2)}{2m} \frac{1}{\max_{\alpha \in [0; \pi]} a_r(\alpha)} \quad (51)$$

Notice that:

$$a_r(\alpha) = b_r(\alpha) (1 - \exp(-4r^2))^2 \quad (52)$$

where:

$$b_r(\theta) = b_r^1(\theta) + b_r^2(\theta) \quad (53)$$

and

$$b_r^1(\theta) = \left(1 - \frac{1}{n}\right) \frac{2}{r} \exp(2r^2 \cos(\theta)); \quad (54)$$

$$b_r^2(\theta) = \frac{1}{n} - \exp(2r^2 \cos(\theta)); \quad (55)$$

Therefore:

$$\max_{\theta \in [0; \frac{\pi}{2}]} b_r(\theta) = \max_{\theta \in [0; \frac{\pi}{2}]} b_r^1(\theta) + \max_{\theta \in [0; \frac{\pi}{2}]} b_r^2(\theta) \quad (56)$$

Denote:  $b^1 = \max_{\theta \in [0; \frac{\pi}{2}]} b_r^1(\theta)$  and  $b^2 = \max_{\theta \in [0; \frac{\pi}{2}]} b_r^2(\theta)$ . Note that:

$$\frac{db_r^1(\theta)}{d\theta} = \exp(2r^2 \cos(\theta)) \left(1 - \frac{1}{n}\right) \frac{2}{r} (1 - r^2 \sin(\theta)) \quad (57)$$

and

$$\frac{db_r^2(\theta)}{d\theta} = \exp(2r^2 \cos(\theta)) \frac{1}{n} (1 - 2r^2 \sin(\theta)) \quad (58)$$

Thus, from the properties of function:  $\sin(\theta)$  and the fact that  $r \leq 1$ , we conclude that both derivatives are first non-negative, then non-positive and then non-negative and that the unique local maximum on the interval  $[0; \frac{\pi}{2}]$  is achieved for  $\theta = \frac{\pi}{2}$ . Note also that  $b_r^1(\frac{\pi}{2}) = b_r^1(0)$  and  $b_r^1(0) = b_r^2(0) = 0$ ,  $b_r^1(\theta) = \left(1 - \frac{1}{n}\right) \frac{2}{r} \exp(-2r^2)$ ,  $b_r^2(\theta) = \frac{1}{n} \exp(-2r^2)$ . We conclude that global maximum for  $b_r^i$  on the interval  $[0; \frac{\pi}{2}]$  for  $i = 1, 2$  is achieved either in its unique local maximum on that interval or for  $\theta = \frac{\pi}{2}$ . Let us consider first  $b_r^1$ . In its local maximum on  $[0; \frac{\pi}{2}]$  we have:

$$\sin(\theta) = \frac{1}{r^2} \quad (59)$$

Since  $\sin(\theta) \leq \frac{1}{2}$  on  $[0; \frac{\pi}{2}]$ , we get:

$$\left(\frac{1}{r^2}\right)^2 \leq \frac{1}{2r^2}; \quad (60)$$

i.e.:

$$r \leq \frac{1}{\sqrt{2}} \quad (61)$$

Therefore:

$$b_r^1(\theta) = \left(1 - \frac{1}{n}\right) \frac{1}{2r} \exp(2r^2) = b_r^1(\theta) \quad (62)$$

We thus conclude that:

$$\max_{\theta \in [0; \frac{\pi}{2}]} b_r^1(\theta) = \left(1 - \frac{1}{n}\right) \frac{1}{2r} \exp(2r^2) \quad (63)$$

By the completely analogous analysis applied to  $b_r^2$  we obtain:

$$\max_{\theta \in [0; \frac{\pi}{2}]} b_r^2(\theta) = \frac{1}{2n} \frac{1}{r} \exp(2r^2) \quad (64)$$

Now, using Equation 51, Equation 52, and Equation 56, we obtain:

$$S_{(r)}(\mathcal{SM}_{m;n}^{\text{anghyb}}) = \frac{\exp(2r^2)}{2mr} \left(1 - \exp(-4r^2)\right) \frac{1}{n} \frac{1}{n} + \frac{1}{n} \frac{1}{n} \quad (65)$$

and that completes the first part of the proof (proof of Inequality 27). The equations on the limits are directly implied by the fact that:

$$S_{(r)}(\mathcal{SM}_{m;n}^{\text{anghyb}}) = \frac{\exp(2r^2)}{2m} h_r(\theta); \quad (66)$$

□

#### D.4 DISCUSSION OF THEOREM 3.5

Theorem 3.5 leads to yet another important conclusion not discussed in the main body of the paper. Note that asymptotically as  $d \rightarrow 0$  or  $d \rightarrow \infty$ , the relative error decreases (as a function of  $m$  and  $n$ ) as  $\frac{1}{mn} = \left(\frac{1}{d_{\text{est}}}\right)$ , where  $d_{\text{est}}$  stands for the number of random features of the resulting estimator. Note also that for the regular estimator the rate of decrease is also  $\frac{1}{d_{\text{est}}}$  (here we treat all other arguments of the formula for the relative error as constants; we already know that the dependence on them of the relative error for the HRF estimators is superior to the one for regular estimators). The key difference from the computational point of view is that for the HRF estimator, the random feature map can be constructed in  $\mathcal{O}(md + md + mn)$ , whereas for the regular estimator it requires  $\mathcal{O}(mnd)$  (see: Sec. L). Thus for the regular estimator random feature map computation is much more expensive. An important application, where random feature map computations takes substantial time of the overall compute time are implicit-attention Transformers, such as Performers (Choromanski et al., 2021b). This is also the setting, where an overwhelming fraction of the approximate softmax kernel values will be extreme (very small or large) since an overwhelming fraction of the attention matrix entries after some training will have extreme values (very small or large). In the setting, where the lengths of queries and keys do not vary much (for instance a prominent case where they all have the same fixed length) that corresponds to angle values  $\theta \rightarrow 0$  or  $\theta \rightarrow \pi$ . This is exactly the scenario from the second part of Theorem 3.5.

#### E PROOF OF LEMMA 3.4

Proof. Equation 25 follows directly from Lemma 2.2. Notice that the relative error  $\text{err}(\mathcal{SM}_m^{\text{trig}})$  is an increasing function of  $\sin^2(\frac{\theta}{2})$  and thus is largest for  $\theta = \pi$ . Plugging in this value into the formula of the relative error gives us the expression from the statement of the lemma. Completely analogous analysis holds for  $\mathcal{SM}_m^{++}$ .  $\square$

#### F GAUSSIAN HYBRID ESTIMATORS

In this section we will provide more results regarding Gaussian hybrid estimators. We will focus on the bipolar scenario and within this scenario on the so-called *canonical* variant described below.

If  $x$  and  $y$  have fixed  $L_2$ -norm ( $\|x\|_2 = \|y\|_2 = r$ ), we can propose a normalized bipolar Gaussian hybrid estimator such that  $\phi(x; y) = 0$  if  $x; y = 0$  and  $\phi(x; y) = 1$  if  $x; y = \pi$ . Furthermore the variance of that estimator, that we call  $\text{SM}_{m;n}^{\text{gausshyb}}$ , will be zeroed out for  $x; y = 0$  or  $x; y = \pi$ , but not for both.

The coefficient-function  $\phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as:

$$\phi(x; y) = \frac{1}{2} \exp\left(-\frac{r^2}{2} \langle x, y \rangle\right) \quad (67)$$

where  $\langle \cdot, \cdot \rangle$  is given as:

$$\langle x, y \rangle = 1 - \exp(-2r^2) \quad (68)$$

The hyperparameter controls the smoothness of the Gaussian kernel.

The exponential in  $\phi$  can be estimated either by trigonometric or positive random features. It is easy to see that for the fixed input lengths, in the former case the variance of the estimator is zero for  $x; y = 0$  and in the latter it is zero for  $x; y = \pi$ , but the variance never zeroes out for both  $x; y = 0$  and  $x; y = \pi$ . In principle, one can derive entire hierarchy of HRF estimators, where the coefficients in an HRF estimator from one level of hierarchy are estimated with the use of HRF estimators from the higher level. In that context, angular hybrid estimator can be applied to provide approximation of the coefficients of the normalized bipolar Gaussian hybrid estimator leading to variance zeroing out for both  $x; y = 0$  and  $x; y = \pi$ . However such nested constructions are not the topic of this paper. From now on we will assume that trigonometric random features are applied to estimate  $\phi$ -coefficient, Therefore we have:

$$b(x; y) = \frac{1}{n} \sum_{i=1}^n \cos(\theta_i(x - y)) \quad (69)$$

and that leads to the following choice of  $\theta_i$ :  $\theta_i \in [0, 2\pi]$  from Equation 10:

$$\begin{cases} \theta_i = \frac{1}{n} \\ \theta_{1:n}(z) = \theta_{2:n}(z) = \frac{1}{n} (\sin(\theta_1 z); \cos(\theta_1 z); \dots; \sin(\theta_n z); \cos(\theta_n z)) \end{cases} \quad (70)$$

where  $\theta_1, \dots, \theta_n \sim \mathcal{N}(0, I_d)$  and  $i^2 = 1$ .

Using Theorem 3.1, we conclude that:

**Theorem F.1** (MSE of the normalized bipolar Gaussian estimator). Take the normalized bipolar Gaussian estimator  $\mathcal{SM}_{m;n}^{\text{gaussHyb}}(x; y)$ , where  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$  are chosen independently i.e. their random projections are chosen independently (note that we always assume that  $b(x; y)$  is chosen independently from  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$ ). Denote:  $\theta = x - y$ . Then the following holds:

$$\text{MSE}(\mathcal{SM}_{m;n}^{\text{gaussHyb}}(x; y)) = E[b^2(x; y)] \text{MSE}(\mathcal{SM}_m^{++}(x; y)) + E[(1 - b(x; y))^2] \text{MSE}(\mathcal{SM}_m^{\text{trig}}(x; y)) \quad (71)$$

where

$$E[(b(x; y))^2] = \frac{1}{2} (1 - \exp(-\frac{2}{2} k^2 k^2))^2 + \frac{1}{2 \cdot 2^n} (1 - \exp(-2k^2 k^2))^2 \quad (72)$$

$$E[(1 - b(x; y))^2] = \frac{1}{2} ((1 - \exp(-\frac{2}{2} k^2 k^2))^2 + \frac{1}{2 \cdot 2^n} (1 - \exp(-2k^2 k^2))^2) \quad (73)$$

Furthermore, if  $\mathcal{SM}_m^{\text{trig}}(x; y)$  and  $\mathcal{SM}_m^{++}(x; y)$  apply the exact same sets of random projections, the mean squared error of the hybrid estimator remains the same.

**Proof.** We can calculate  $E[(b(x; y))^2]$  and  $E[(1 - b(x; y))^2]$  as follows:

$$\begin{aligned} E[(1 - b(x; y))^2] &= \frac{1}{2} E[(1 + \frac{1}{n} \sum_{i=1}^n \cos(\theta_i(x - y)))^2] \\ &= \frac{(1 - \frac{1}{2})^2}{2} + \frac{2(1 - \frac{1}{2})}{2} E[\cos(\theta_1(x - y))] + \frac{1}{2} E[(\frac{1}{n} \sum_{i=1}^n \cos(\theta_i(x - y)))^2] \\ &= \frac{(1 - \frac{1}{2})^2}{2} + \frac{2(1 - \frac{1}{2})}{2} \exp(-\frac{2}{2} k^2 k^2) + \frac{1}{2} \exp(-2k^2 k^2) + \frac{1}{2 \cdot 2^n} (1 - \exp(-2k^2 k^2))^2 \\ &= \frac{1}{2} ((1 - \exp(-\frac{2}{2} k^2 k^2))^2 + \frac{1}{2 \cdot 2^n} (1 - \exp(-2k^2 k^2))^2 \end{aligned} \quad (74)$$

$$\begin{aligned} E[(b(x; y))^2] &= 1 - 2E[1 - b(x; y)] + E[(1 - b(x; y))^2] \\ &= 1 + \frac{2}{2} \exp(-\frac{2}{2} k^2 k^2) + E[(1 - b(x; y))^2] \\ &= \frac{1}{2} (1 - \exp(-\frac{2}{2} k^2 k^2))^2 + \frac{1}{2 \cdot 2^n} (1 - \exp(-2k^2 k^2))^2 \end{aligned} \quad (75)$$

By plugging in these two formulae in the expression for MSE, we obtain the first half of the theorem.

Note that if  $\mathcal{SM}_m^{\text{trig}}(\cdot; r)$  and  $\mathcal{SM}_m^{++}(\cdot; r)$  apply the exact same sets of random projections, then  $(1 - \cos(\|x\|_2 \|y\|_2))$  in  $\text{MSE}(\mathcal{SM}_{m;n}^{\text{gausshyb}}(x; y))$  becomes zero in Theorem F.1, therefore the MSE remains the same. We thus obtain the second part of the theorem and that completes the proof.  $\square$

Now let us assume the setting of the same-lengths inputs. We denote the length by  $r$  and by an angle between inputs by  $\theta$ . We can rewrite Eq. 67 and provide equivalent definition by replacing  $\|x\|_2 \|y\|_2$  with their angle  $\theta$  and norm  $r$ :

$$b(\cdot; r) = \frac{1 - \exp(-2r^2 \sin(\frac{\theta}{2})^2)}{2} \quad (76)$$

We obtain the following version of the above theorem:

**Theorem F.2** (MSE of the normalized bipolar Gaussian hybrid estimator for same-length inputs)

Assume that  $\|x\|_2 = \|y\|_2 = r$  and denote:  $\theta = \angle_{x,y}$ . Take the normalized bipolar Gaussian

hybrid estimator  $\mathcal{SM}_{m;n}^{\text{gausshyb}}(\cdot; r)$ , where  $\mathcal{SM}_m^{\text{trig}}(\cdot; r)$  and  $\mathcal{SM}_m^{++}(\cdot; r)$  are chosen independently

i.e. their random projections are chosen independently (note that we always assume that  $\mathcal{SM}_m^{\text{trig}}$

is chosen independently from  $\mathcal{SM}_m^{\text{trig}}(\cdot; r)$  and  $\mathcal{SM}_m^{++}(\cdot; r)$ ). Then the following holds:

$$\text{MSE}(\mathcal{SM}_{m;n}^{\text{gausshyb}}(\cdot; r)) = E[b^2(\cdot; r)] \text{MSE}(\mathcal{SM}_m^{++}(\cdot; r)) + E[(1 - b(\cdot; r))^2] \text{MSE}(\mathcal{SM}_m^{\text{trig}}(\cdot; r)) \quad (77)$$

where

$$E[b(\cdot; r)] = \frac{1}{2} (1 - \exp(-2r^2 \sin(\frac{\theta}{2})^2)) + \frac{1}{2} \frac{1}{2^n} (1 - \exp(-4r^2 \sin(\frac{\theta}{2})^2)) \quad (78)$$

$$E[(1 - b(\cdot; r))^2] = \frac{1}{2} ((1 - \exp(-2r^2 \sin(\frac{\theta}{2})^2))^2) + \frac{1}{2} \frac{1}{2^n} (1 - \exp(-4r^2 \sin(\frac{\theta}{2})^2))^2 \quad (79)$$

Furthermore, if  $\mathcal{SM}_m^{\text{trig}}(\cdot; r)$  and  $\mathcal{SM}_m^{++}(\cdot; r)$  apply the exact same sets of random projections, the mean squared error of the hybrid estimator will remain the same.

**Proof.** For  $\|x\|_2 = \|y\|_2 = r$ , the following holds for  $\theta = \angle_{x,y}$ :

$$\|x\|_2 \|y\|_2 \cos(\theta) = 4r^2 \sin(\frac{\theta}{2})^2 \quad (80)$$

And this theorem holds trivially by replacing  $\|x\|_2 \|y\|_2 \cos(\theta)$  with the above formula in Theorem F.1.  $\square$

## G POINTWISE SOFTMAX KERNEL ESTIMATION EXPERIMENTS

In Fig. 5 we present complete ablation studies regarding empirical relative errors of different RF-based estimators of the softmax kernel over different lengths of inputs and angles  $\theta_{x,y}$ .

## H COMPLEX EXPONENTIAL ESTIMATORS FOR CLUSTERED DATA

### H.1 SYNTHETIC EXPERIMENTS ON DATA ADMITTING CLUSTERING STRUCTURE

Here we test HRF estimators customized to data with clustering structure, as described in Sec. 2.4. Inputs  $x$  are taken from two 50-dimensional 1000-element Gaussian clusters and the inputs  $y$  from two other 50-dimensional 1000-element Gaussian clusters. We use empirical MSE metric and constrain matrix  $A$  to be real diagonal for a compact closed-form formulae (see: Sec. H.3). We created four different clusters' configurations corresponding to different values of  $\sigma(A)$  (see: Eq. 21) for all pairs of clusters and with different concentrations around centers controlled by the standard deviation (small/large values of  $\sigma$  and  $\mu$ ). As shown in Fig. 6, for all variants HRF estimators adapted to clustered data consistently provide notable accuracy gains, even for larger values of  $\sigma$  and less concentrated clusters. Additional experimental details are given in Sec. H.3.

Figure 5: Pointwise estimation of  $\sigma_{M(x; y)}$  for the same-length 64-dim inputs ( $t = 1:0$ ,  $r = 1:25$  and  $r = 1:5$ ) and various angles  $\theta_{x, y}$ . We used  $s = 10000$  estimated softmax values in each subplot. The true softmax value and the 5<sup>th</sup> and 95<sup>th</sup> quantile estimated values are shown by the left y-axis, and the empirical relative errors are shown by the right y-axis. Trigonometric estimator and FAVOR+ applied 128 random features. To make fair comparison, for the hybrid variant the configuration leading to the similar number of FLOPS operations per random feature map creation was applied.

Figure 6: Comparing different estimators for data with clustering structure. The empirical MSE is obtained by averaging over 20 trials. For the HRFs, we apply deterministic coefficients. The fact that the empirical error for FAVOR+ is not perfectly monotonic in  $m$  was first observed in (Luo et al., 2021) (see: Fig. 1: (b)).

## H.2 INSTANTIATING COMPLEX EXPONENTIAL ESTIMATORS FOR CLUSTERED DATA

Assume that the inputs (queries) can be modeled by  $n_q$  clusters and that the inputs (keys) can be modeled by  $n_k$  clusters (e.g. via  $k$ -means clustering algorithm). Denote the center of each cluster as  $r_i \in \mathbb{R}^d$ , ( $i = 1; \dots; n_k$ ) and  $r_j \in \mathbb{R}^d$ , ( $j = 1; \dots; n_q$ ). Then there exist  $n_q n_k$  pairs of  $(r_i; r_j)$ , ( $i = 1; \dots; n_q; j = 1; \dots; n_k$ ), so we can construct  $n_q n_k$  softmax kernel estimators to estimate cross-group softmax kernel values.

Consider  $z = Ar_i + (A^>)^{-1}r_j$  for an invertible (in  $\mathbb{C}^d$ ) matrix  $A \in \mathbb{C}^{d \times d}$ . An estimator based on this mechanism has variance equal to

$$z = Ar_i + (A^>)^{-1}r_j = 0 \quad (81)$$

From now on we constrain  $A$  to be diagonal, so  $A = A^>$ . We can rewrite the above equation as:

$$z = Ar_i + (A)^{-1}r_j = 0 \quad (82)$$

Since position  $(k; k)$  in matrix  $A$  is a complex number of the form  $\alpha_k + \beta_k i$ , we need to satisfy the following equation for each  $k = 1; \dots; d$ :

$$(\alpha_k + \beta_k i)r_{i;k} + \frac{1}{\alpha_k + \beta_k i}r_{j;k} = 0 \quad (83)$$

$$(\alpha_k + \beta_k i)r_{i;k} + \frac{\alpha_k}{\alpha_k^2 + \beta_k^2}r_{j;k} + \frac{\beta_k}{\alpha_k^2 + \beta_k^2}i r_{j;k} = 0 \quad (84)$$

where  $r_{i;k}$  is the  $k$ -th entry of vector  $r_i$ .

We can simplify this equation and separate into real and imaginary part:

$$\begin{aligned} \text{Re} : & \left( \frac{2}{k} + \frac{2}{k} \right) k r_{i;k} + k r_{j;k} = 0 \\ \text{Im} : & \left( \frac{2}{k} + \frac{2}{k} \right) k r_{i;k} - k r_{j;k} = 0 \end{aligned} \quad (85)$$

Our goal now is to choose values for  $k$  and  $k$  given  $r_i$  and  $r_j$ .

If  $r_{i;k} r_{j;k} > 0$ , we can set:

$$k = \frac{p}{r_{j;k} = r_{i;k}} \quad (86)$$

And if  $r_{i;k} r_{j;k} < 0$ , we can set:

$$k = \frac{p}{r_{j;k} = r_{i;k}} \quad (87)$$

When  $r_{i;k} = r_{j;k} = 0$ , we can take  $k = k = 1$ . If  $r_{i;k} = 0; r_{j;k} \neq 0$  or the opposite, then we cannot satisfy the above equation perfectly. We can take to some large positive value and set  $k = 0$  when  $r_{i;k} = 0; r_{j;k} \neq 0$ , and set  $k = k$  to some small positive value close to zero when  $r_{i;k} \neq 0; r_{j;k} = 0$ .

When restricting  $A = \text{diag}(a_1; \dots; a_d)$  to  $\mathbb{R}^d \times \mathbb{R}^d$ , it is easily seen that to minimize  $\|A r_i + (A^{-1})^T r_j\|^2$ , we can take

$$a_k = \frac{q}{r_{j;k} = r_{i;k}} \quad (88)$$

if  $r_{i;k} \neq 0$ . And if  $r_{i;k} = 0; r_{j;k} \neq 0$ ,  $a_k$  can be set to a large positive number. We can set  $a_k = 1$  if  $r_{i;k} = 0; r_{j;k} = 0$ .

In the analysis below we denote matrix  $A$  calculated from Eq. 85 and Eq. 88 (depending on whether we are using real or complex matrices) given  $r_i$  and  $r_j$  as  $A^{ij}$ .

From the analysis in the main body of the paper we know that for arbitrary vectors  $x, y \in \mathbb{R}^d$ :

$$\exp\left(\frac{k A^{ij} x k^2}{2}\right) \exp\left(\frac{k (A^{ij})^{-1} y k^2}{2}\right) \text{SM}(x; y) = E[\exp(\sum_{i,j} (A^{ij} x + (A^{ij})^{-1} y))] \quad (89)$$

Therefore, the softmax kernel can be estimated as:

$$\text{SM}(x; y) \approx \frac{m_{A^{ij}}(x) \cdot m_{(A^{ij})^{-1}}(y)}{m_{(A^{ij})^{-1}}(y)} \quad (90)$$

with  $m_{A^{ij}}(x)$  and  $m_{(A^{ij})^{-1}}(y)$  given by:

$$m_{A^{ij}}(x) \stackrel{\text{def}}{=} \frac{1}{m} \exp\left(\frac{k A^{ij} x k^2}{2}\right) (\exp(\sum_{i,j} A^{ij} x); \dots; \exp(\sum_{i,j} A^{ij} x)) \quad (91)$$

$$m_{(A^{ij})^{-1}}(y) \stackrel{\text{def}}{=} \frac{1}{m} \exp\left(\frac{k (A^{ij})^{-1} y k^2}{2}\right) (\exp(\sum_{i,j} (A^{ij})^{-1} y); \dots; \exp(\sum_{i,j} (A^{ij})^{-1} y)) \quad (92)$$

for  $i, j, \dots, m \in \{1, \dots, d\}$ .

Such a softmax kernel estimator has variance equal to zero if  $r_i$  and  $r_j$  are close if we are using the complex matrix and has small variance if we are using the real matrix. For other pairs of vectors  $(x; y)$  close to  $r_i$  and  $r_j$ , the variance is close to zero, or relatively small.

We can also take the Gaussian coefficient estimator that gets its maximum value when  $r_i; y = r_j$  as  $m_{ij}(x; y)$ . Such an estimator could be constructed with in a similar way.

It can be easily linearized since:

$$\begin{aligned} s_{ij}(x; y) &= \exp\left(\frac{kA^{ij}x + (A^{ij})^{-1}yk^2}{2^2}\right) = \\ & \exp\left(\frac{kA^{ij}xk^2}{2^2}\right)\exp\left(\frac{k(A^{ij})^{-1}yk^2}{2^2}\right)\text{SM}\left(\frac{x}{k}; \frac{y}{k}\right); \end{aligned} \quad (93)$$

where  $\text{SM}(x; y)$  could be estimated with similar random feature maps as given by Eq. 90. Therefore, the Gaussian coefficients can be estimated as:

$$s_{ij}(x; y) = \exp\left(\frac{kA^{ij}x + (A^{ij})^{-1}yk^2}{2^2}\right) \hat{m}_{A^{ij}}(x) \hat{m}_{(A^{ij})^{-1}}(y) \quad (94)$$

with  $\hat{m}_{A^{ij}}(x)$  and  $\hat{m}_{(A^{ij})^{-1}}(y)$  given by:

$$\hat{m}_{A^{ij}}(x) = \frac{1}{m} \exp\left(\frac{kA^{ij}xk^2}{2}\right) (\exp(-\lambda_1 A^{ij}x); \dots; \exp(-\lambda_m A^{ij}x)) \quad (95)$$

$$\hat{m}_{(A^{ij})^{-1}}(y) = \frac{1}{m} \exp\left(\frac{k(A^{ij})^{-1}yk^2}{2}\right) (\exp(-\lambda_1 (A^{ij})^{-1}y); \dots; \exp(-\lambda_m (A^{ij})^{-1}y)) \quad (96)$$

for  $\lambda_1; \dots; \lambda_m \sim \mathcal{N}(0; I_d)$  chosen independently and that were applied in base estimators.

We summarize this section with the following two constructions of the hybrid estimators for the clustered data that are implied by the above analysis.

**Definition H.1 (Hybrid Gaussian-Mixtures estimators for clustered data)** Assume that inputs (queries) and (keys) can be modeled by  $n_q$  and  $n_k$  clusters respectively with centers  $r_i \in \mathbb{R}^d$ , ( $i = 1; \dots; n_q$ ,  $j = 1; \dots; n_k$ ). We denote  $A^{ij}$  as the complex matrix satisfying Eq. (85) with center  $r_i; r_j$ . Furthermore, we denote  $\mathbf{E} = (\text{SM}^{ij}(x; y))_{i=1; j=1}^{n_q; n_k}$  as a list of estimators of the softmax kernel  $\text{SM}^{ij}(x; y)$  and  $\mathbf{b} = (b^{ij}(x; y))_{i=1; j=1}^{n_q; n_k}$  as a list of estimators of  $s_{ij}(x; y)$  constructed independently from  $\mathbf{E}$ . We also use one additional base estimator ( $\text{SM}^{\text{trig}}(x; y)$  or  $\text{SM}^{++}(x; y)$ ) and denote the estimator of its softmax kernel  $\text{SM}^0(x; y)$  and -coefficient as  $b^0(x; y)$ . Then our hybrid estimator takes the following form:

$$\text{SM}^{\mathbf{E}}(x; y) = \sum_{i=1}^{n_q} \sum_{j=1}^{n_k} b^{ij}(x; y) \text{SM}^{ij}(x; y) + b^0(x; y) \text{SM}^0(x; y) \quad (97)$$

with constraint:

$$\sum_{i=1}^{n_q} \sum_{j=1}^{n_k} b^{ij}(x; y) + b^0(x; y) = 1 \quad (98)$$

and where the estimators of coefficients and base estimators are given as:

$$b^{ij}(x; y) = \hat{m}_{A^{ij}}(x) \hat{m}_{(A^{ij})^{-1}}(y) \quad (99)$$

$$\text{SM}^{ij}(x; y) = \hat{m}_{A^{ij}}(x) \hat{m}_{(A^{ij})^{-1}}(y) \quad (100)$$

for  $\hat{m}_{A^{ij}}(x)$ ;  $\hat{m}_{(A^{ij})^{-1}}(y)$ ;  $\hat{m}_{A^{ij}}(x)$ ;  $\hat{m}_{(A^{ij})^{-1}}(y)$  given by Eq. (91, 92, 95, 96).

Definition H.2 (Hybrid Zero-One-Mixtures estimators for clustered data) Assume that inputs (queries) and (keys) can be modeled by  $n_q$  and  $n_k$  clusters respectively with centers  $\mu_i \in \mathbb{R}^d$ , ( $i = 1; \dots; n_q$ ,  $j = 1; \dots; n_k$ ). We denote  $A^{ij}$  as the complex matrix satisfying Eq. (85) with centers  $\mu_i; \mu_j$ . Furthermore, we denote  $\mathbf{E} = (\mathbf{SM}^{ij}(x; y))_{i=1; j=1}^{n_q; n_k}$  as a list of estimators of the softmax kernel  $\mathbf{SM}^{ij}(x; y)$  and  $\mathbf{b} = (b^{ij}(x; y))_{i=1; j=1}^{n_q; n_k}$  as a list of estimators of  $ij(x; y)$  constructed independently from  $\mathbf{E}$ .

Then our hybrid estimator takes the following form:

$$\mathbf{SM}^{\mathbf{E}}(x; y) = \sum_{i=1}^{n_q} \sum_{j=1}^{n_k} b^{ij}(x; y) \mathbf{SM}^{ij}(x; y) \quad (101)$$

with constraint:

$$\sum_{i=1}^{n_q} \sum_{j=1}^{n_k} b^{ij}(x; y) = 1 \quad (102)$$

and where the estimators of coefficients and base estimators are given as:

$$b^{ij}(x; y) = \mu_i(x) \mu_j(y) \quad (103)$$

$$\mathbf{SM}^{ij}(x; y) = \frac{m_{A^{ij}}(x)}{m_{(A^{ij})^{-1}}(y)} \quad (104)$$

for  $\frac{m_{A^{ij}}(x)}{m_{(A^{ij})^{-1}}(y)}$  given by Eq. (91, 92), with  $\mu_i(x)$  being a scalar indicating whether  $x$  belongs to the  $i$ -th cluster of the  $n_q$  clusters and similarly with  $\mu_j(y)$  being a scalar indicating whether  $y$  belongs to the  $j$ -th of the  $n_k$  clusters.

In our experiments we used the formulation from the second definition with  $\mathbf{E}$  constrained to be real diagonal. In other words, we are using Eq. 88 to construct  $\mathbf{E}$ .

### H.3 ADDITIONAL EXPERIMENTAL DETAILS

To create the synthetic data, we first generate two random vectors  $X_0$  and  $Y_0$  in  $\mathbb{R}^{50}$ . Then we generate four random orthogonal matrices  $S_i \in \mathbb{R}^{50 \times 50}$  where  $i = 1; 2$  and  $j = 1; 2$  and use  $X_i = O_i^{-1} X_0$  and  $Y_j = O_j^{-1} Y_0$  for  $i = 1; 2$  as the mean vectors for our Gaussian clusters. For each pair of clusters, the minimal values of Eq. 21 may be non-zero  $\text{sign}(x_{i;k}) = \text{sign}(y_{j;k})$  for some  $k$  due to the usage of real diagonal matrices. To control the values, we manually adjust the sign of each element of  $X_i$  and  $Y_j$ . The data for input  $x$  is the combination of two clusters, each consisting 1000 data points following Gaussian distribution with mean vectors  $X_1, X_2$  and the common covariance matrix  $\Sigma$ . And similar constructions are done for input  $y$ . We create four synthetic data sets, with different values of  $\sigma$  and magnitudes of the standard deviation. After this step, we normalize the data points by controlling their  $\ell_2$  norms to make the true softmax value be within a reasonable range. After the data is created, we first use K-means clustering algorithm `kmeans` to cluster  $x$  and  $y$ , and we obtain centers  $(\mu_i; \mu_j)_{i,j \in \{1,2\}}$ . Then we use the set of real diagonal matrices to minimize Eq. 21 for all pairs of centers of  $x$  and  $y$  to generate four complex exponential base estimators. To choose coefficients, we use the more efficient approach described above, i.e. using indicator vectors, and these are already computed in the clustering process. For each data set, we compare the performance of these estimators by calculating the mean square error of them, and we do this multiple times with different numbers of random features used. All the results are obtained by averaging over 20 repetitions. In addition to the average performance which can be found in Fig. 6, we also record the maximal and minimal empirical MSE over 20 repetitions for all estimators in the Table 3, where  $s_1 = [0.08; 0.05; 0.04; 0.03]$ ;  $s_2 = [0.11; 0.09; 0.08; 0.07]$ ;  $s_3 = [0.05; 0.03; 0.03; 0.02]$ ;  $s_4 = [0.07; 0.06; 0.05; 0.04]$  representing the lists of values for all pairs of clusters for different synthetic data sets.

Data Set	Number of Random Features	Hybrid-Cluster	FAVOR+	Trigonometric
$s = s_1, \quad = 1$	20	[0.0016, 0.0183]	[0.0071, 0.1455]	[0.0039, 0.0808]
	50	[0.0011, 0.0081]	[0.0026, 0.2541]	[0.0031, 0.0510]
	80	[0.0005, 0.0052]	[0.0032, 0.0282]	[0.0010, 0.0281]
	100	[0.0006, 0.0036]	[0.0022, 0.0528]	[0.0013, 0.0170]
	120	[0.0004, 0.0025]	[0.0028, 0.0534]	[0.0012, 0.0135]
	150	[0.0003, 0.0020]	[0.0014, 0.0582]	[0.0007, 0.0134]
	200	[0.0002, 0.0010]	[0.0013, 0.0205]	[0.0004, 0.0087]
$s = s_2, \quad = 1$	20	[0.0032, 0.0265]	[0.0095, 0.2071]	[0.0036, 0.0989]
	50	[0.0012, 0.0079]	[0.0036, 0.2123]	[0.0019, 0.0565]
	80	[0.0007, 0.0059]	[0.0037, 0.0570]	[0.0008, 0.0208]
	100	[0.0007, 0.0040]	[0.0016, 0.0280]	[0.0010, 0.0164]
	120	[0.0005, 0.0025]	[0.0013, 0.0262]	[0.0007, 0.0246]
	150	[0.0004, 0.0024]	[0.0010, 0.0407]	[0.0004, 0.0190]
	200	[0.0003, 0.0018]	[0.0018, 0.0236]	[0.0004, 0.0134]
$s = s_3, \quad = \frac{p}{10}$	20	[0.0083, 0.0300]	[0.0362, 0.1218]	[0.0174, 0.0490]
	50	[0.0046, 0.0183]	[0.0125, 0.1262]	[0.0081, 0.0249]
	80	[0.0029, 0.0105]	[0.0108, 0.0281]	[0.0044, 0.0143]
	100	[0.0024, 0.0056]	[0.0077, 0.0421]	[0.0041, 0.0095]
	120	[0.0020, 0.0046]	[0.0077, 0.0402]	[0.0034, 0.0077]
	150	[0.0015, 0.0041]	[0.0049, 0.0395]	[0.0025, 0.0073]
	200	[0.0011, 0.0031]	[0.0037, 0.0150]	[0.0018, 0.0050]
$s = s_4, \quad = \frac{p}{10}$	20	[0.0097, 0.0294]	[0.0332, 0.1553]	[0.0179, 0.0485]
	50	[0.0049, 0.0172]	[0.0176, 0.1345]	[0.0058, 0.0260]
	80	[0.0031, 0.0113]	[0.0113, 0.0415]	[0.0038, 0.0110]
	100	[0.0028, 0.0061]	[0.0080, 0.0292]	[0.0031, 0.0100]
	120	[0.0022, 0.0045]	[0.0065, 0.0185]	[0.0028, 0.0128]
	150	[0.0016, 0.0044]	[0.0045, 0.0235]	[0.0021, 0.0094]
	200	[0.0012, 0.0031]	[0.0046, 0.0166]	[0.0015, 0.0067]

Table 3: Minimal and maximal empirical MSE with 20 repetitions

## I LANGUAGE MODELING TRAINING DETAILS AND ADDITIONAL RESULTS

For the Language Modeling tasks, we trained a 2-layer LSTM of hidden sizes 200 and 650 on the PennTree Bank (Marcus et al., 1993) and the WikiText2 dataset (Merity et al., 2017) respectively. We tied the weights of the word embedding layer and the decoder layer (Inan et al., 2017; Press & Wolf, 2017). Thus we could treat the language modeling problem as minimizing the cross entropy loss between the dot product of the model output (queries) and the class embedding (keys) obtained from the embedding layer with the target word.

We now present detailed statistics of HRFs on the Penn Tree Bank dataset. The mean and the standard deviation of the statistical metric results in Fig. 3, Sec. 4.2 is shown in Table 4. The distribution of the lengths of the keys and the queries are shown in Figure 7. For our experiments with the Angular Hybrid and the Gaussian Hybrid, the number of random features are 8, 16, 32, 64 for the base estimators, and 8 for the coefficient estimators.

Table 4: Results are computed over 10 runs on Penntree Bank Dataset. **Boldface** denotes the best one, and underline denotes the second best. Negative fractions for Favor+ is not reported as it produces positive random features.

1d-Wasserstein				
RF types	$d = 64$	$d = 128$	$d = 256$	$d = 512$
FAVOR+	$3171.35 \pm 69.04$	$3142.65 \pm 85.55$	$3050.58 \pm 108.25$	$2985.18 \pm 73.15$
Trig.	$1577.44 \pm 13.59$	$1582.50 \pm 13.42$	$1574.22 \pm 7.01$	$1576.50 \pm 5.83$
Gaussian	<u><math>1520.32 \pm 11.65</math></u>	<u><math>1521.10 \pm 8.79</math></u>	<u><math>1516.50 \pm 6.66</math></u>	<u><math>1515.41 \pm 5.96</math></u>
Angular	<b><math>1395.03 \pm 79.91</math></b>	<b><math>1445.24 \pm 58.27</math></b>	<b><math>1394.97 \pm 62.73</math></b>	<b><math>1425.88 \pm 65.82</math></b>
KS statistics				
FAVOR+	$0.573 \pm 0.0073$	$0.569 \pm 0.0098$	$0.558 \pm 0.013$	$0.551 \pm 0.0092$
Trig.	$0.424 \pm 0.0040$	$0.427 \pm 0.0041$	$0.425 \pm 0.0024$	$0.424 \pm 0.0023$
Gaussian	$0.411 \pm 0.0029$	<b><math>0.410 \pm 0.0019</math></b>	<b><math>0.410 \pm 0.0013</math></b>	<b><math>0.410 \pm 0.0014</math></b>
Angular	<b><math>0.404 \pm 0.017</math></b>	<u><math>0.419 \pm 0.010</math></u>	$0.418 \pm 0.0089$	$0.414 \pm 0.0093$
Negative fraction				
Trig.	$0.257 \pm 0.0057$	$0.258 \pm 0.0039$	$0.258 \pm 0.0021$	$0.257 \pm 0.0017$
Gaussian	$0.155 \pm 0.0039$	$0.159 \pm 0.0019$	$0.157 \pm 0.0011$	<b><math>0.151 \pm 0.0009</math></b>
Angular	<b><math>0.152 \pm 0.0041</math></b>	<b><math>0.158 \pm 0.0027</math></b>	<b><math>0.153 \pm 0.0018</math></b>	<u><math>0.157 \pm 0.0013</math></u>

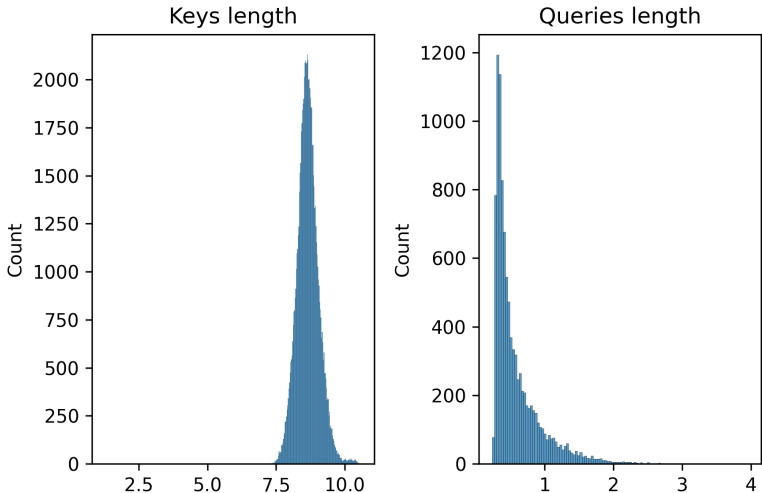


Figure 7: Distribution of the lengths of the keys and queries in the PennTree Bank dataset

Since the HRF mechanisms can be sensitive to the lengths of the keys and the queries, for the language modeling task on the WikiText2 dataset, we added a regularizer to constrain the lengths of the keys and the queries to be close to 1. However, constraining the lengths close to 1 hurts the model performance and so we chose to use a temperature scaling as in (Rawat et al., 2019). Thus before passing the computed dot product to the softmax layer, we scaled the dot product by  $\tau > 1$ . Our final loss function for the language modeling task was:

$$\mathcal{L}(\theta) = L_{CE} + \lambda_1 \mathbb{E}(\|q\|_2 - 1)^2 + \lambda_2 \mathbb{E}(\|k\|_2 - 1)^2 \quad (105)$$

where  $\|\cdot\|_2$  is the row-wise  $L_2$  norm of the appropriate matrices, and  $L_{CE}$  is the cross-entropy loss. The distribution of the lengths of the keys and the queries are shown in Figure 9.

Finally we present some additional results on the WikiText2 dataset. For our experiments on the WikiText2 dataset, we chose  $\lambda_1 = \lambda_2 = 2$  and  $\tau = 6$ . Our model trained with these hyperparameters achieve a perplexity score of 105.35 on the test set after 50 epochs.

We compute the 1-dimensional Wasserstein distance and the Kolmogrov-Smirnov (KS) metric (Ramdas et al., 2017) between the approximated softmax distribution and the true softmax distribution. The mean and the standard deviation of the statistical metric results in Fig. 8 is shown in Table 5. For our experiments with the Angular Hybrid and the Gaussian Hybrid, the number of random features are 8, 16, 32, 64 for the base estimators, and 8 for the coefficient estimators.

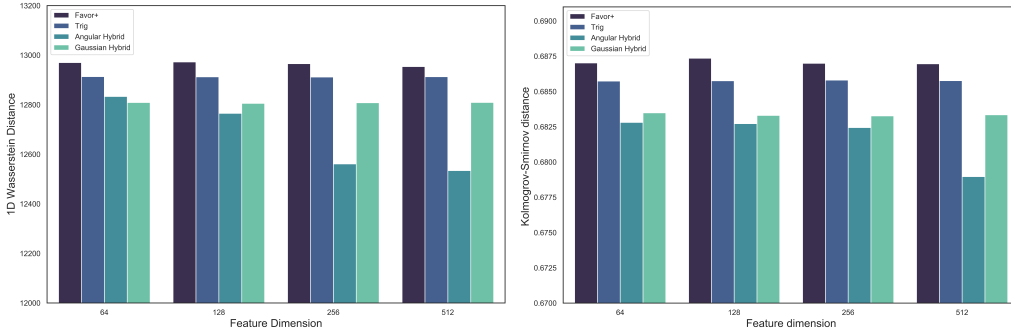


Figure 8: Statistical metrics measuring softmax matrix approximation quality on the WikiText2 dataset. For standard estimators, the number of random features are 64, 128, 256, 512. To make fair comparison, for the hybrid variants, the configurations leading to the similar number of FLOPS operations per random feature map creation were applied. Results reported over 10 runs.

Table 5: Results are computed over 10 runs on the Wikitext2 Dataset. **Boldface** denotes the best one, and underline denotes the second best.

1d-Wasserstein				
RF types	$d = 64$	$d = 128$	$d = 256$	$d = 512$
FAVOR+	$12970.38 \pm 119.49$	$12971.67 \pm 108.79$	$12965.82 \pm 89.62$	$12954.17 \pm 83.95$
Trig.	$12913.66 \pm 59.22$	$12912.65 \pm 57.42$	$12910.91 \pm 37.23$	$12913.25 \pm 35.12$
Gaussian	<b><math>12809.03 \pm 97.2</math></b>	$12805.71 \pm 78.91$	$12808.02 \pm 84.57$	$12809.40 \pm 65.88$
Angular	<u><math>12833.56 \pm 41.87</math></u>	<b><math>12765.44 \pm 38.91</math></b>	<b><math>12561.26 \pm 26.48</math></b>	<b><math>12534.48 \pm 25.88</math></b>
KS statistics				
FAVOR+	$0.687 \pm 0.0043$	$0.687 \pm 0.0051$	$0.687 \pm 0.0019$	$0.686 \pm 0.0067$
Trig.	$0.685 \pm 0.0071$	$0.685 \pm 0.0053$	$0.685 \pm 0.0029$	$0.685 \pm 0.0014$
Gaussian	<u><math>0.683 \pm 0.0092</math></u>	<u><math>0.683 \pm 0.0073</math></u>	<u><math>0.683 \pm 0.0049</math></u>	<u><math>0.682 \pm 0.0089</math></u>
Angular	<b><math>0.682 \pm 0.0079</math></b>	<b><math>0.682 \pm 0.0066</math></b>	<b><math>0.682 \pm 0.0037</math></b>	<b><math>0.678 \pm 0.0064</math></b>

The above tables (Table 4 and Table 5) show that the our hybrid variants consistently outperform Favor+ and the trigonometric random features.

## I.1 LANGUAGE MODELING USING HRF

In this subsection, we will describe how one can use HRF to train a LSTM on the language modeling task on the PennTree Bank dataset, similar to the experiments carried out in (Rawat et al., 2019).

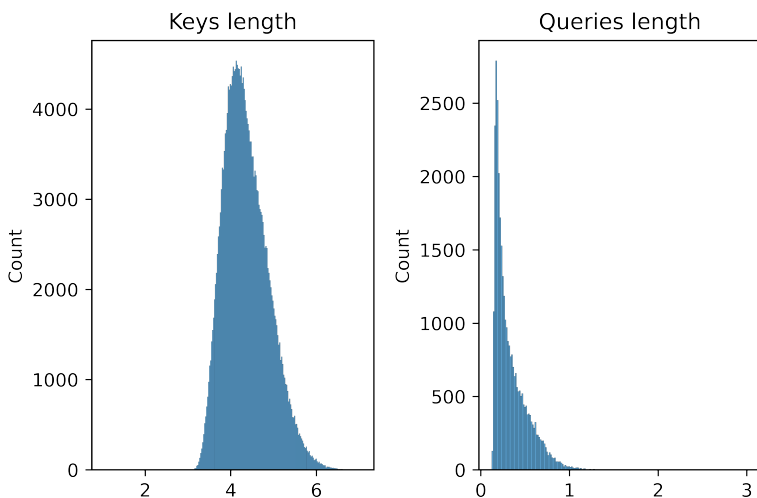


Figure 9: Distribution of the lengths of the keys and queries in the WikiText2 dataset

We trained a 2-layer LSTM model with hidden and output sizes of 200, and used the output as input embedding for sampled softmax. We sampled 40 negative (other than true) classes out of 10000 classes to approximate the expected loss in each training epoch. As observed in Fig. 2, the relative error of all three estimators grow exponentially fast with embedding norm  $r$ . Therefore, if we keep un-normalized embeddings during sampling, then even though we could get an unbiased estimation of the loss, the variance could be high. Such bias-variance trade off is also mentioned in paper (Rawat et al., 2019). To solve this issue, we used normalized input and class embeddings during sampling to generate biased sampling distribution with lower variance, while keeping un-normalized embeddings to calculate the loss. We trained our model for 80 epochs, with batch size equal to 20, dropout ratio in LSTM equal to 0.5. Implementation details could be seen in our Github repository.

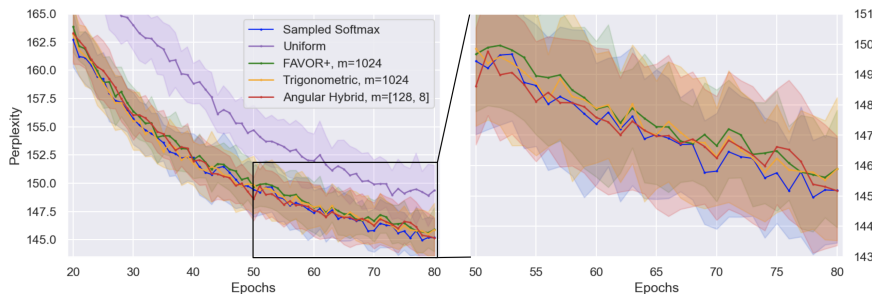


Figure 10: Language Modeling with softmax sampling: training results on PennTree Bank dataset (50-80 epoch window zoomed in on the right subfigure). The solid line for each method is estimated over 25 independent runs. The shaded areas represent perplexity within 1 standard deviation of the average. FAVOR+/trigonometric mechanisms used 1024 random features. To make fair comparison, for HRFs, the configurations leading to the similar number of FLOPS operations per random feature map creation were applied. 128 and 8 random features is used to estimate the softmax estimators and  $\beta$ -coefficient in HRFs respectively. Even though statistical metrics for HRFs are better in Fig. 3, the difference between different random feature estimators is not very significant.

We compared our hybrid estimator with Trigonometric, FAVOR+, the uniform method which sampled 40 negative classes with equal probability, as well as the sampled softmax method which uses the unbiased true probability calculated from full softmax as the sampling probability. Trigonometric and FAVOR+ mechanisms used 1024 random features. To make comparison fair, hybrid variants were using  $(m, n)$ -configurations characterized by the similar number of FLOPS needed to construct their corresponding random features as regular RF-estimators. We reported our comparison

results averaged over 25 independent runs in Fig. 10. The *perplexity* score in Fig. 10 is defined as  $2^{\text{cross entropy loss}}$ . We could conclude that even though the statistical metrics for HRFs are better in Fig. 3, the difference between HRFs and other random feature estimators in this specific softmax sampling downstream task is not very significant.

## J SPEECH EXPERIMENTS: ADDITIONAL DETAILS

Tested Conformer-Performer models consisted of  $l = 17$  conformer layers. Each attention layer used  $H = 8$  heads. The embedding dimensionality was  $p = 512$  and since dimensions were split equally among different heads, query/key dimensionality was set up to  $d_{QK} = 64$ . Each input sequence was of length  $L \approx 500$ . Padding mechanism was applied for all tested variants.

## K DOWNSTREAM ROBOTICS EXPERIMENTS

In both Robotics experiments we found query/key  $L_2$ -normalization technique particularly effective (queries and keys of  $L_2$ -norm equal to one) thus we applied it in all the experiments.

### K.1 VISUAL LOCOMOTION WITH QUADRUPEL ROBOTS

We evaluate HRF-based attention RL policies in a robotic task of learning locomotion on uneven terrains from vision input. We use the quadruped from Unitree called Laikago (lai). It has 12 actuated joints, 3 per leg. The task is to walk forward on a randomized uneven terrain that requires careful foot placement planning based on visual feedback. The ground is made of a series of step-stones with gaps in between. The step stones widths are fixed at 50 cm, the lengths are between [50, 80] cm in length, and the gap size between adjacent stones are between [10, 20] cm. It perceives the ground through 2 depth cameras attached to its body, one on the front and other on the belly facing downwards. We use Implicit Attention Policy (IAP) architecture (masking variant) described in Choromanski et al. (2021a) which uses Performer-based attention mechanism to process  $32 \times 24$  depth images from the 2 cameras.

We apply a hierarchical setup to solve this task. The high level uses IAP-rank with masking to process camera images and output the desired foot placement position. The low level employs a position-based swing leg controller, and a model predictive control (MPC) based stance leg controller, to achieve the foot placement decided by high level. The policies are trained with evolutionary strategies (ES). In Fig. 4 (left) we compare FAVOR+ with angular hybrid approximation in IAP. HRF with  $8 \times 8$  random projections ( $m = n = 8$ ) performs as well as the softmax kernel with 256 random projections. A series of image frames along the episode of a learned locomotion HRF-based IAP policy is shown bottom right of Fig. 4. On the top-left corner of the images, the input camera images are attached. The red part of the camera image is the area masked out by self-attention. The policy learns to pay attention to the gaps in the scene in order to avoid stepping on them.

The training curves in Fig. 4 (left) are obtained by averaging over 5 top runs. The shaded regions shows the standard deviation over multiple runs.

### K.2 ROBOTIC MANIPULATION: BI-MANUAL SWEEPING

Here we consider the bi-manual sweep robotic-arm manipulation. The two-robotic-arm system needs to solve the task of placing red balls into green bowls (all distributions are equally rewarded, so in particular robot can just use one bowl).

In this bi-manual sweeping task, the reward per step is between 0 and 1, the fraction of the blocks within the green bowls. The loss for this task is the negative log likelihood between the normalized softmax probability distribution of sampled actions (one positive example, and all others uniform negative counter-examples), and the ground truth one-hot probability distribution for a given observation from the training dataset. For more details see (van den Oord et al., 2018). This is a 12-DoF cartesian control problem, the state observation consists of the current cartesian pose of the arms and a  $500 \times 500 \times 3$  RGB image. An example of the camera image is given in Fig. 11.

The training curves Fig. 4 (right) are obtained by averaging over 3 learning rates:  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$ .

