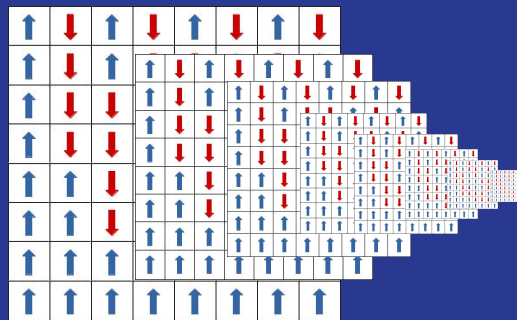
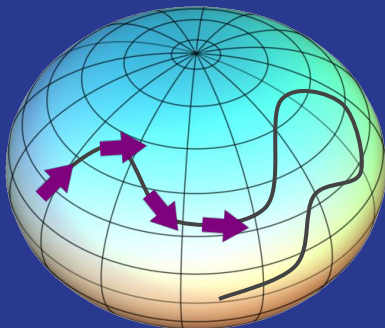


The  
Alan Turing  
Institute

# Unifying Orthogonal Monte Carlo Methods

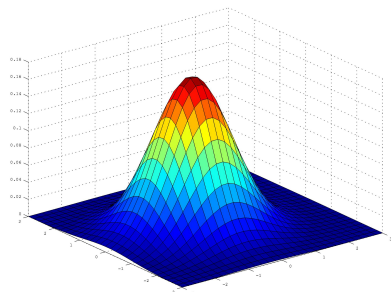
From **Kac's Random Walks** To **Hadamard Multi Rademachers**



Krzysztof Choromanski, Mark Rowland  
Wenyu Chen, Adrian Weller

# The Phenomenon of Orthogonal Monte Carlo Estimators

Estimation task:  $\mathbb{E}_{X \sim \mu} [f(X)]$



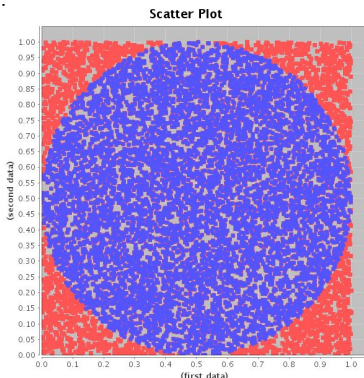
isotropic distribution  
(e.g. Gaussian)

Applications:

- dimensionality reduction (JLT-mechanisms)
- scaling kernel methods (random feature maps)
- hashing algorithms (e.g. LSH)
- (sliced) Wasserstein distances (WGANs, autoencoders...)
- reinforcement learning (ES algorithms)
- and many, many more...

Standard MC approach:

$$\frac{1}{N} \sum_{i=1}^N f(X_i), \text{ where } (X_i)_{i=1}^N \stackrel{\text{i.i.d.}}{\sim} \mu.$$



# The Phenomenon of Orthogonal Monte Carlo Estimators

**Estimation task:**  $\mathbb{E}_{X \sim \mu} [f(X)]$

Sampling from the Haar measure on the  $O(d)$  group

$$G_{\text{ort}} = \begin{pmatrix} g_{1,1}^{\text{ort}} & g_{1,2}^{\text{ort}} & \dots & g_{1,n}^{\text{ort}} \\ g_{2,1}^{\text{ort}} & g_{2,2}^{\text{ort}} & \dots & g_{2,n}^{\text{ort}} \\ \dots & \dots & \dots & \dots \\ g_{m,1}^{\text{ort}} & g_{m,n}^{\text{ort}} & \dots & g_{m,n}^{\text{ort}} \end{pmatrix}$$

isotropic distribution (e.g. Gaussian)



Expensive:  $O(n^3)$  time

**The Orthogonal Trick:** guarantees unbiasedness

$$\frac{1}{N} \sum_{i=1}^N f(X_i^{\text{ort}}), \text{ where } (X_i^{\text{ort}}) \sim \mu \text{ and } X_i^{\text{ort}} \perp X_j^{\text{ort}}.$$

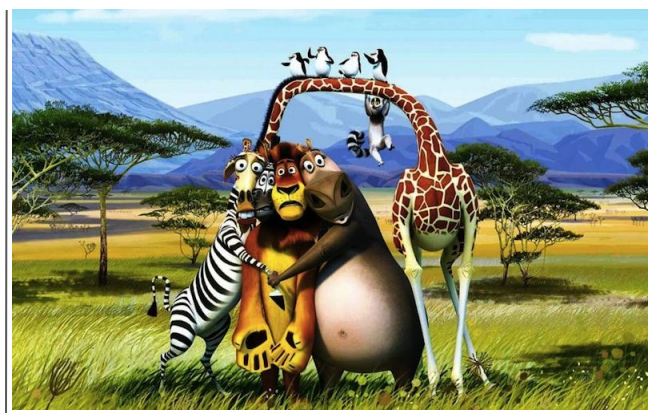


# of samples of the MC estimator  $\leq$  dim



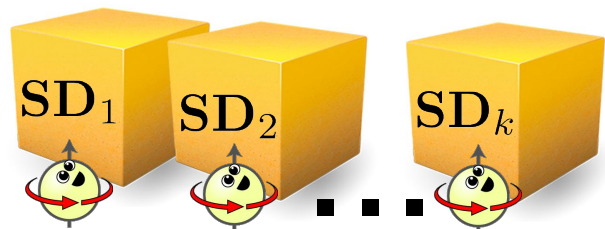
often implies better accuracy

# Towards Computational Efficiency: The Zoo of Approximate MCs

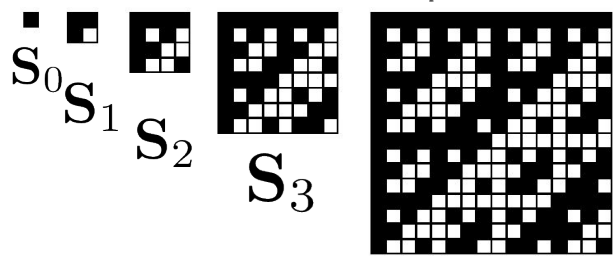


# Towards Computational Efficiency:

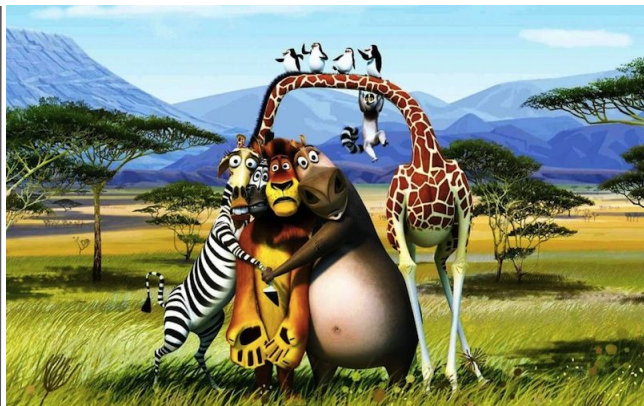
## The Zoo of Approximate MCs



$$M_{SR}^{(k)} = \prod_{i=1}^k SD_i^{(\mathcal{R})} \rightarrow |\lambda_i| = 1$$
$$\lambda_i \sim Unif\{-1, +1\}$$

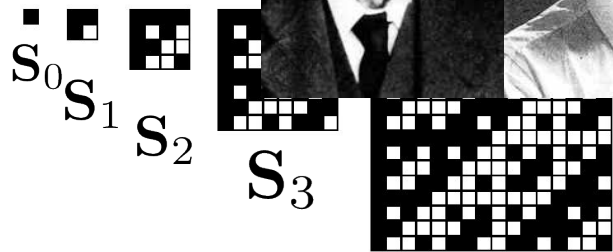
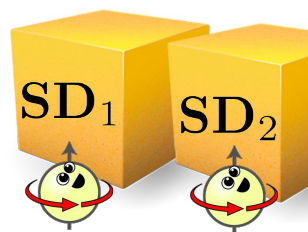


$$S_4 D = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$



# Towards Computational Efficiency:

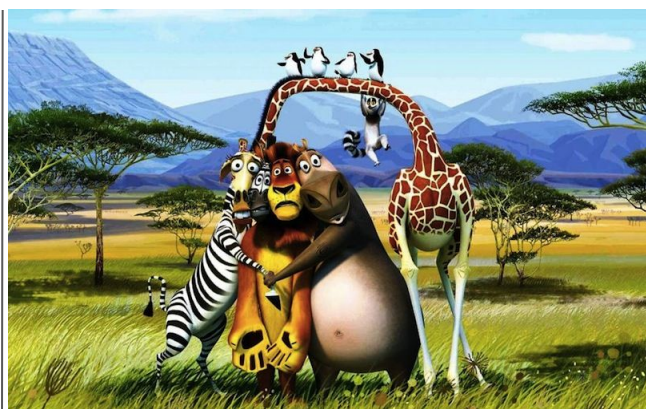
## The Zoo of Approximate MCs



$$\prod_{i=1}^k \text{SD}_i^{(\mathcal{R})} \rightarrow |\lambda_i| = 1$$

$\sim \text{Unif}\{-1, +1\}$

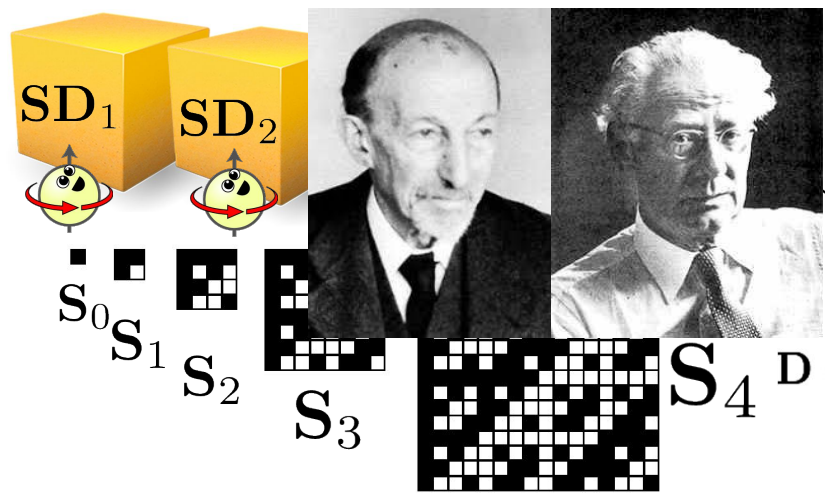
$$\mathbf{S}_4 \mathbf{D} = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$





# Towards Computational Efficiency:

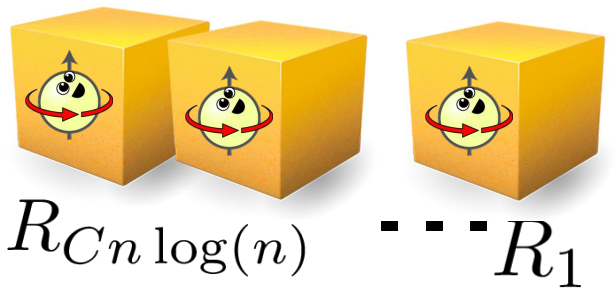
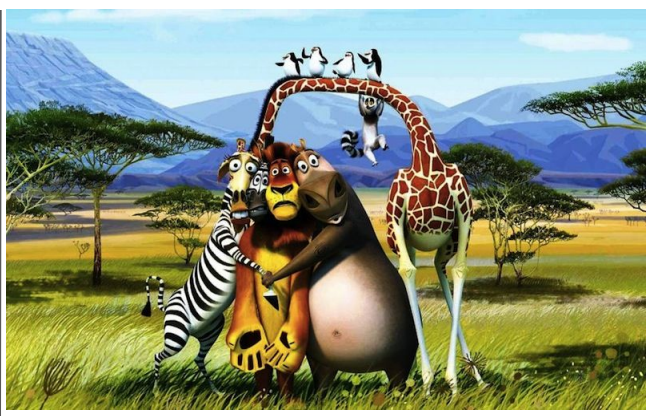
## The Zoo of Approximate MCs



$$\prod_{i=1}^k \text{SD}_i^{(\mathcal{R})} \rightarrow |\lambda_i| = 1$$

$$\sim \text{Unif}\{-1, +1\}$$

$$\mathbf{S}_4 \mathbf{D} = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$



$$\mathbf{R}[i, j] = \begin{cases} 1, & \text{if } i = j \text{ and } i \notin \{I, J\} \\ 0, & \text{if } i \neq j \text{ and } \{i, j\} \neq \{I, J\} \\ \cos \Theta, & \text{if } i = j \text{ and } i \in \{I, J\} \\ \sin \Theta, & \text{if } i = J, j = I \\ -\sin \Theta, & \text{if } i = I, j = J \end{cases}$$

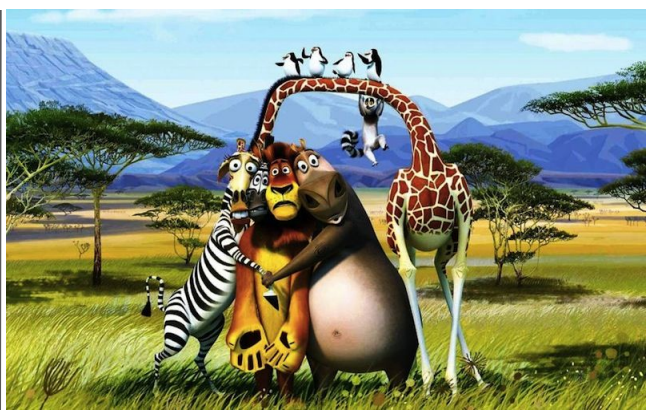
# Towards Computational Efficiency:

## The Zoo of Approximate MCs

$$\prod_{i=1}^k \text{SD}_i^{(\mathcal{R})} \rightarrow |\lambda_i| = 1$$

$\sim \text{Unif}\{-1, +1\}$

$$\mathbf{S}_4 \mathbf{D} = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$



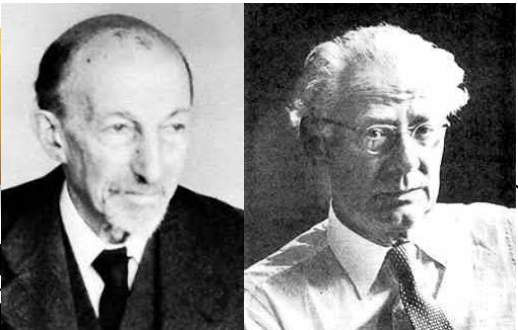
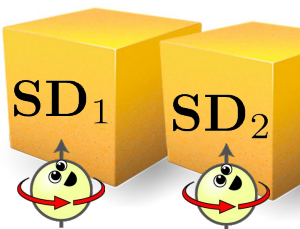
$$RC n \log(n)$$

$$j = \begin{cases} 1, & \text{if } i = j \text{ and } i \notin \{I, J\} \\ 0, & \text{if } i \neq j \text{ and } \{i, j\} \neq \{I, J\} \\ \cos \Theta, & \text{if } i = j \text{ and } i \in \{I, J\} \\ \sin \Theta, & \text{if } i = J, j = I \\ -\sin \Theta, & \text{if } i = I, j = J \end{cases}$$



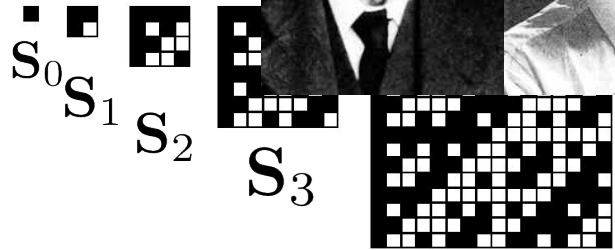
# Towards Computational Efficiency:

## The Zoo of Approximate MCs

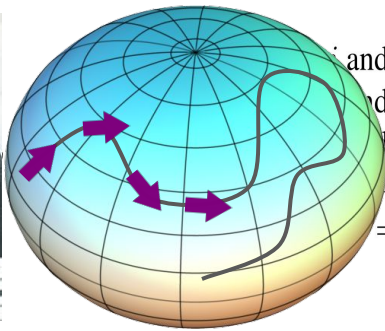
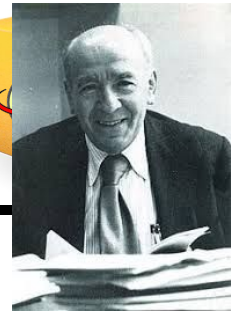
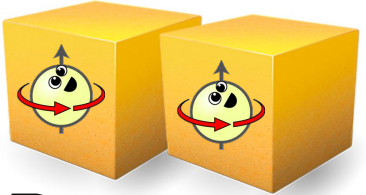
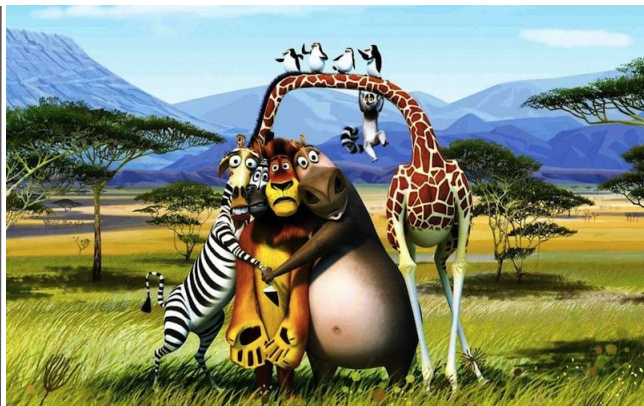


$$\prod_{i=1}^k \text{SD}_i^{(\mathcal{R})} \rightarrow |\lambda_i| = 1$$

$\sim \text{Unif}\{-1, +1\}$



$$\mathbf{S}_4 \mathbf{D} = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

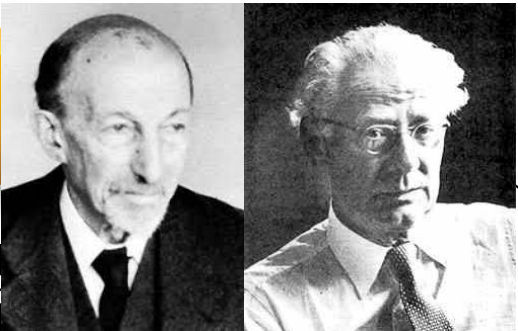


and  $i \notin \{I, J\}$   
 and  $\{i, j\} \neq \{I, J\}$   
 $i \in \{I, J\}$   
 $= I$   
 $= J$

$$RCn \log(n)$$

# Towards Computational Efficiency:

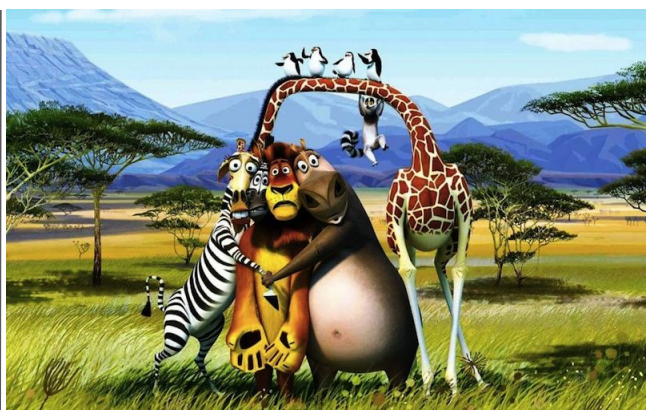
## The Zoo of Approximate MCs



$\prod_{i=1}^k \text{SD}_i^{(\mathcal{R})} \rightarrow |\lambda_i| = 1$   
 $\sim \text{Unif}\{-1, +1\}$

$\text{SD}_1$   $\text{SD}_2$   
 $S_0$   $S_1$   $S_2$   $S_3$   $S_4$

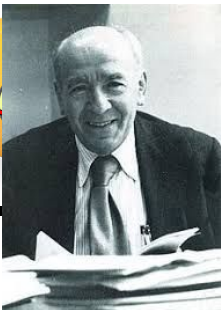
$\mathbf{D} = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$



$B^{(N)} = \begin{bmatrix} A_1^{(N/2)} C_{n-1} & A_2^{(N/2)} S_{n-1} \\ -A_1^{(N/2)} S_{n-1} & A_2^{(N/2)} C_{n-1} \end{bmatrix}$

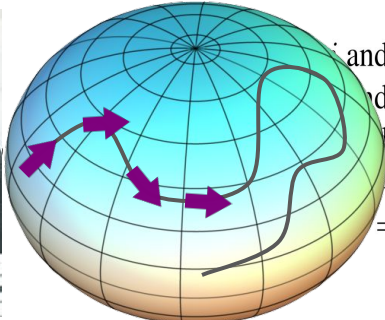
$\Downarrow$   
**size**  
**N x N**

$\Downarrow$   
**size**  
**N/2 x N/2**



$R C n \log(n)$

and  $i \notin \{I, J\}$   
 and  $\{i, j\} \neq \{I, J\}$   
 $i \in \{I, J\}$   
 $= I$   
 $= J$



$N = 2^n$

**Constraints:**

- $C_{n-1}^2 + S_{n-1}^2 = I$
- $C_{n-1} S_{n-1} = S_{n-1} C_{n-1}$

# Towards Computational Efficiency:

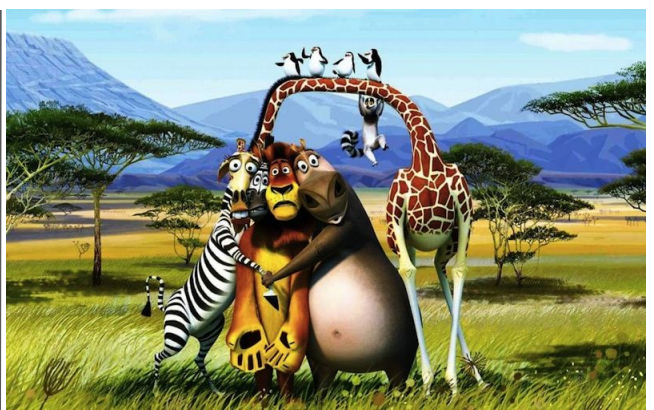
## The Zoo of Approximate MCs

$SD_1$   $SD_2$   
 $S_0$   $S_1$   $S_2$   $S_3$   $S_4$

$$\prod_{i=1}^k SD_i^{(\mathcal{R})} \rightarrow |\lambda_i| = 1$$

$\sim Unif\{-1, +1\}$

$$D = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots \\ 0 & \lambda_2 & 0 & \dots \\ 0 & 0 & \lambda_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$



$$B^{(N)} = \begin{bmatrix} A_1^{(N/2)} C_{n-1} & A_2^{(N/2)} S_{n-1} \\ -A_1^{(N/2)} S_{n-1} & A_2^{(N/2)} C_{n-1} \end{bmatrix}$$

$\Downarrow$   
**size**  
 $N \times N$

$\Downarrow$   
**size**  
 $N/2 \times N/2$

$RCn \log(n)$

and  $i \notin \{I, J\}$   
 and  $\{i, j\} \neq \{I, J\}$   
 $i \in \{I, J\}$   
 $= I$   
 $= J$

$$N = 2^n$$

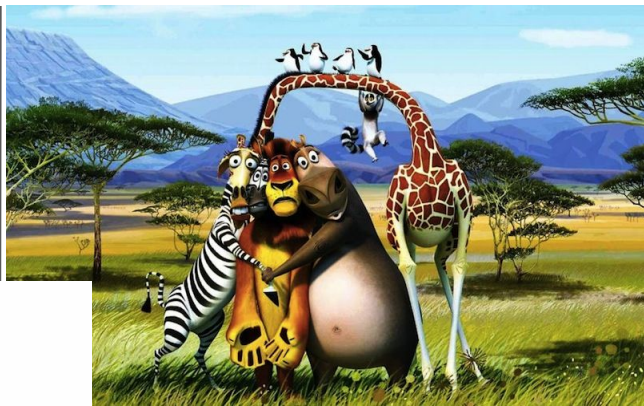
### Constraints:

- $C_{n-1}^2 + S_{n-1}^2 = I$
- $C_{n-1} S_{n-1} = S_{n-1} C_{n-1}$

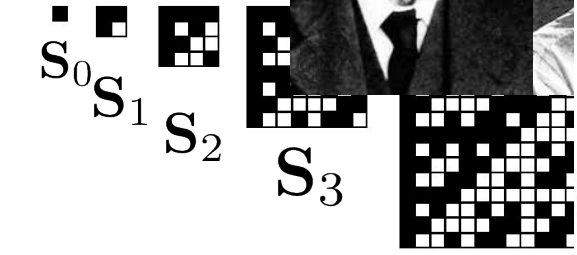
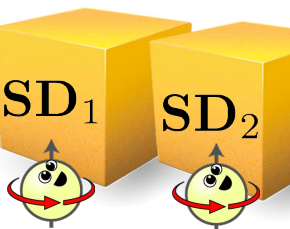


# Towards Computational Efficiency:

## The Zoo of Approximate MCs



$$\prod_{i=1}^k \text{SD}(\mathcal{R}) \rightarrow |\lambda_i| = 1$$



$$) = \begin{bmatrix} A_1^{(N/2)} C_{n-1} & A_2^{(N/2)} S_{n-1} \\ -A_1^{(N/2)} S_{n-1} & A_2^{(N/2)} C_{n-1} \end{bmatrix}$$

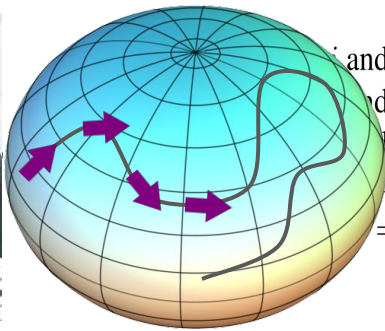
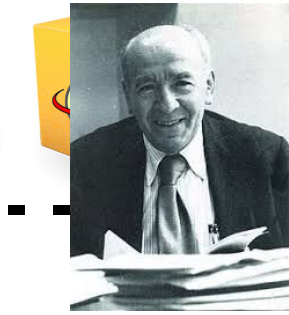
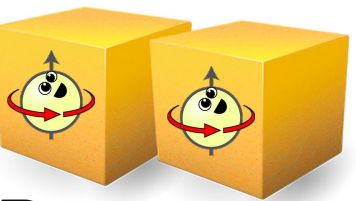


size  
N/2 x N/2

$$N = 2^n$$

### Constraints:

- $C_{n-1}^2 + S_{n-1}^2 = I$
- $C_{n-1} S_{n-1} = S_{n-1} C_{n-1}$



and  $i \notin \{I, J\}$   
and  $\{i, j\} \neq \{I, J\}$   
 $i \in \{I, J\}$   
 $= I$   
 $= J$

$$RCn \log(n)$$

# On the Hunt for the Unifying Theory:

## The World of Givens Reflections and Rotations

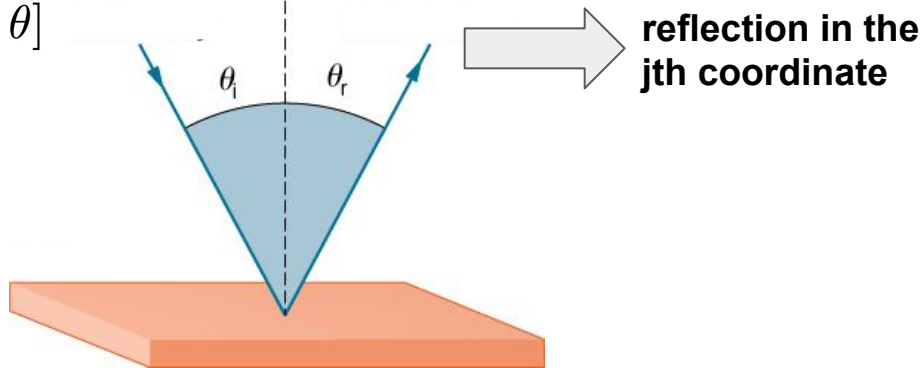


**Givens rotations**

$$\mathbf{G}[i, j, \theta]_{k,l} = \begin{cases} \cos(\theta) & \text{if } k = l \in \{i, j\} \\ -\sin(\theta) & \text{if } k = i, l = j \\ \sin(\theta) & \text{if } k = j, l = i \\ 1 & \text{if } k = l \notin \{i, j\} \\ 0 & \text{otherwise.} \end{cases}$$

**Givens reflections**

$$\tilde{\mathbf{G}}[i, j, \theta]$$



**Kac's random walk matrices**

$$\mathbf{K}_T = \prod_{t=1}^T \mathbf{G}[I_t, J_t, \theta_t]$$

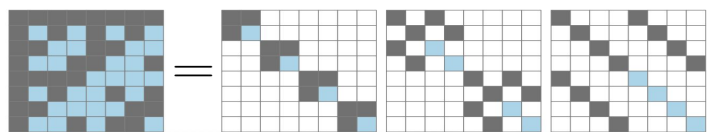
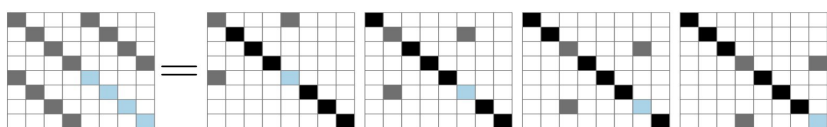
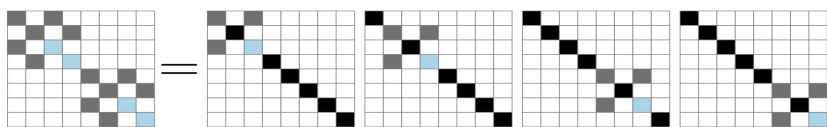
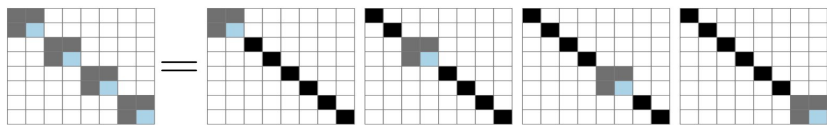
**Hadamard-Rademacher Chains**

$$\mathbf{X}_T = \prod_{t=1}^T \mathbf{HD}_t$$



# On the Hunt for the Unifying Theory:

## The World of Givens Reflections and Rotations



Kac's random walk matrices

$$\mathbf{K}_T = \prod_{t=1}^T \mathbf{G}[I_t, J_t, \theta_t]$$

Hadamard-Rademacher Chains

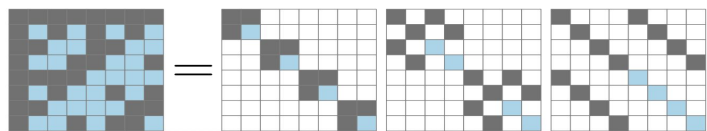
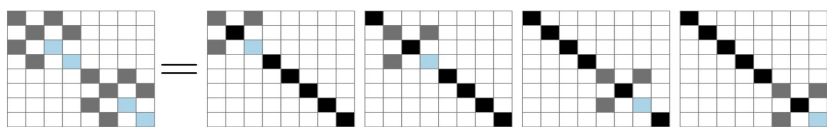
$$\mathbf{X}_T = \prod_{t=1}^T \mathbf{HD}_t$$

$$\mathbf{HD}_t = \left( \prod_{i=1}^{L-1} \tilde{\mathbf{F}}^{i,L} \right) (\tilde{\mathbf{F}}^{L,L} \mathbf{D}_t)$$

$$\tilde{\mathbf{F}}^{j,L} = \prod_{\lambda \in \mathbb{F}_2^L, \lambda_j=0} \tilde{\mathbf{G}}[\lambda, \lambda + \mathbf{e}_j, \pi/4] \in \mathcal{O}(2^L)$$

# On the Hunt for the Unifying Theory:

## The World of Givens Reflections and Rotations



Hadamard-MultiRademachers

$$\mathbf{M}_L = \prod_{i=1}^L \left( \tilde{\mathbf{F}}^{i,L} \mathbf{D}_i \right)$$

Butterfly Matrices



$$\mathbf{F}^{j,L}[(\theta_j, \mu)_{\mu \in \mathbb{F}_2^{L-j}}] = \prod_{\substack{\lambda \in \mathbb{F}_2^L \\ \lambda_j = 0}} \mathbf{G}[\lambda, \lambda + \mathbf{e}_j, \theta_j, \lambda_{j+1:L}] \in \mathcal{O}(2^L)$$

$$((\theta_i, \mu)_{\mu \in \mathbb{F}_2^{L-i}})_{i=1}^L \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0, 2\pi))$$

$$\mathbf{B}_L = \prod_{i=1}^L \mathbf{F}^{i,L}[(\theta_i, \mu)_{\mu \in \mathbb{F}_2^{L-i}}]$$

$$\tilde{\mathbf{F}}^{j,L} = \prod_{\lambda \in \mathbb{F}_2^L, \lambda_j = 0} \tilde{\mathbf{G}}[\lambda, \lambda + \mathbf{e}_j, \pi/4] \in \mathcal{O}(2^L)$$

# First Theoretical Results for Free-Lunch Phenomenon in the Nonlinear Regime



**Theorem** (Kac's random walk estimators of RBF kernels). *Let  $K_d : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be the Gaussian kernel and let  $\epsilon > 0$ . Let  $\mathcal{B}$  be a set satisfying  $\text{diam}(\mathcal{B}) \leq B$  for some universal constant  $B$  that does not depend on  $d$  ( $\mathcal{B}$  might be for instance a unit sphere). Then there exists a constant  $C = C(B, \epsilon) > 0$  such that for every  $\mathbf{x}, \mathbf{y} \in \mathcal{B} \setminus S(\epsilon)$  and  $d$  large enough we have:*

$$\text{MSE}(\widehat{K}_{\text{kac}}^{\phi, m, k}(\mathbf{x}, \mathbf{y})) < \text{MSE}(\widehat{K}_{\text{base}}^{\phi, m}(\mathbf{x}, \mathbf{y})),$$

where  $k = C \cdot d \log d$  and  $m = ld$  for some  $l \in \mathbb{N}$ .

# First Theoretical Results for Free-Lunch Phenomenon in the Nonlinear Regime



**Theorem** (Kac's kernels). Let  $K_d : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be a kernel and let  $\epsilon > 0$ . Let  $B$  be a universal constant  $B$  be for instance a unit ball in  $\mathbb{R}^d$ . Let  $C = C(B, \epsilon) > 0$  such that for  $d$  large enough we have

Still more accurate estimator than unstructured MC baseline

for RBF kernel and  $\exists$  for some  $l$  ( $B$  might be a constant  $S(\epsilon)$  and  $d$

$$\text{MSE}(\widehat{K}_{\text{kac}}^{\phi, m, k}(\mathbf{x}, \mathbf{y})) < \text{MSE}(\widehat{K}_{\text{base}}^{\phi, m}(\mathbf{x}, \mathbf{y})),$$

where  $k = C \cdot d \log d$  and  $m = ld$  for some  $l \in \mathbb{N}$ .

# First Theoretical Results for Free-Lunch Phenomenon in the Nonlinear Regime



**Theorem** (Kac's random walk estimators of RBF kernels). Let  $\mathcal{B}$  be a compact set with  $\text{diam}(\mathcal{B}) \leq B$  for some  $B > 0$ . Let  $\epsilon > 0$  and let  $S(\epsilon)$  be a subset of  $\mathcal{B}$  with  $|S(\epsilon)| \geq \epsilon^d$ . Then there exists a constant  $C = C(d, \epsilon, B)$  such that for any  $d \geq 1$  and  $n \geq C \cdot d \log d$ , there exists a constant  $k = k(d, \epsilon, B)$  such that for any  $m \geq k$  and any  $\mathbf{x}, \mathbf{y} \in \mathcal{B} \setminus S(\epsilon)$ , we have

**Log-Linear Time Complexity**  
(unstructured MC baseline has quadratic)

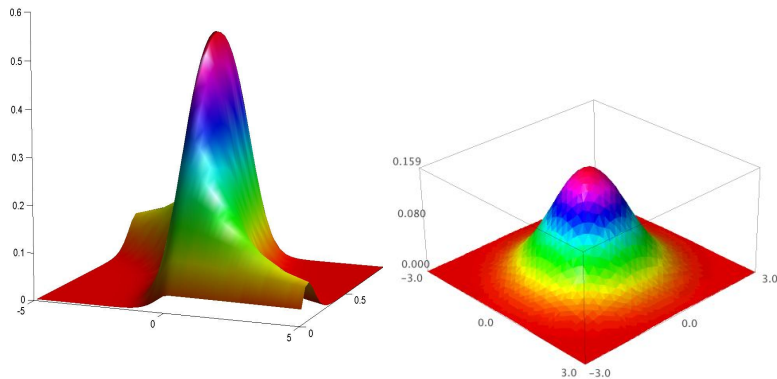
(Kac's random walk estimators of RBF kernels). Let  $\mathcal{B}$  be a compact set with  $\text{diam}(\mathcal{B}) \leq B$  for some  $B > 0$ . Let  $\epsilon > 0$  and let  $S(\epsilon)$  be a subset of  $\mathcal{B}$  with  $|S(\epsilon)| \geq \epsilon^d$ . Then there exists a constant  $C = C(d, \epsilon, B)$  such that for any  $d \geq 1$  and  $n \geq C \cdot d \log d$ , there exists a constant  $k = k(d, \epsilon, B)$  such that for any  $m \geq k$  and any  $\mathbf{x}, \mathbf{y} \in \mathcal{B} \setminus S(\epsilon)$  and  $d$

$$\mathbb{E}(\widehat{K}_{\text{base}}^{\phi, m}(\mathbf{x}, \mathbf{y})),$$

where  $k = C \cdot d \log d$  and  $m = ld$  for some  $l \in \mathbb{N}$ .



# First Theoretical Results for Free-Lunch Phenomenon in the Nonlinear Regime



Analysis of the Total Variation Distance between Haar measure on d-sphere and measure induced by standard Kac's random walk on d-sphere

$$\text{MSE}(\underbrace{Y}_{\text{estimator}}) = \mathbb{E}[(Y - \underbrace{\mu}_{\text{estimated value}})^2] = \int_0^\infty \mathbb{P}[|Y - \mu| > \sqrt{t}] dt$$

**Pillai, Smith 2016**

**Kac's random walk on d-sphere mixes in  $O(d \log d)$  steps**

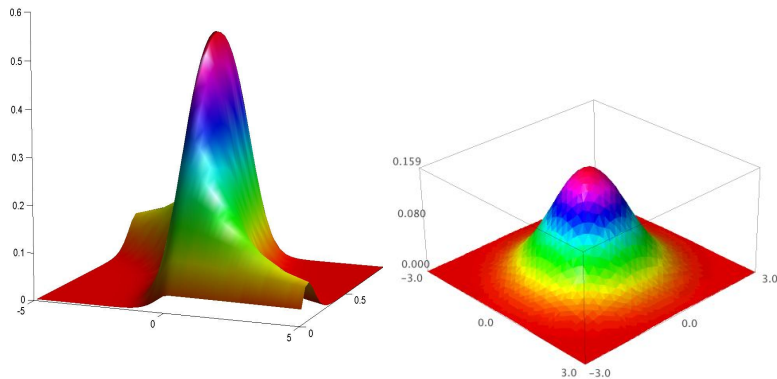
**Theorem** Fix  $C_1 < \frac{1}{2}$  and  $C_2 > 200$ . If the sequence of times  $\{T_1(n)\}_{n \in \mathbb{N}}$  satisfies  $T_1(n) < C_1 n \log(n)$  for all  $n$ , then

$$\lim_{n \rightarrow \infty} \inf_{X_0 \in S^{n-1}} \|\mathcal{L}(X_{T_1(n)}) - \mu\|_{\text{TV}} = 1.$$

If the sequence of times  $\{T_2(n)\}_{n \in \mathbb{N}}$  satisfies  $T_2(n) > C_2 n \log(n)$  for all  $n$ , then

$$\lim_{n \rightarrow \infty} \sup_{X_0 \in S^{n-1}} \|\mathcal{L}(X_{T_2(n)}) - \mu\|_{\text{TV}} = 0.$$

# First Theoretical Results for Free-Lunch Phenomenon in the Nonlinear Regime



Analysis of the Total Variation Distance between Haar measure on d-sphere and measure induced by standard Kac's random walk on d-sphere

$$\text{MSE}(\underbrace{Y}_{\text{estimator}}) = \mathbb{E}[(Y - \underbrace{\mu}_{\text{estimated value}})^2] = \int_0^\infty \mathbb{P}[|Y - \mu| > \sqrt{t}] dt$$

**Pillai, Smith 2016**

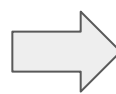
**Kac's random walk on d-sphere mixes in  $O(d \log d)$  steps**

**Theorem** Fix  $C_1 < \frac{1}{2}$  and  $C_2 > 200$ . If the sequence of times  $\{T_1(n)\}_{n \in \mathbb{N}}$  satisfies  $T_1(n) < C_1 n \log(n)$  for all  $n$ , then

$$\lim_{n \rightarrow \infty} \inf_{X_0 \in S^{n-1}} \|\mathcal{L}(X_{T_1(n)}) - \mu\|_{\text{TV}} = 1.$$

If the sequence of times  $\{T_2(n)\}_{n \in \mathbb{N}}$  satisfies  $T_2(n) > C_2 n \log(n)$  for all  $n$ , then

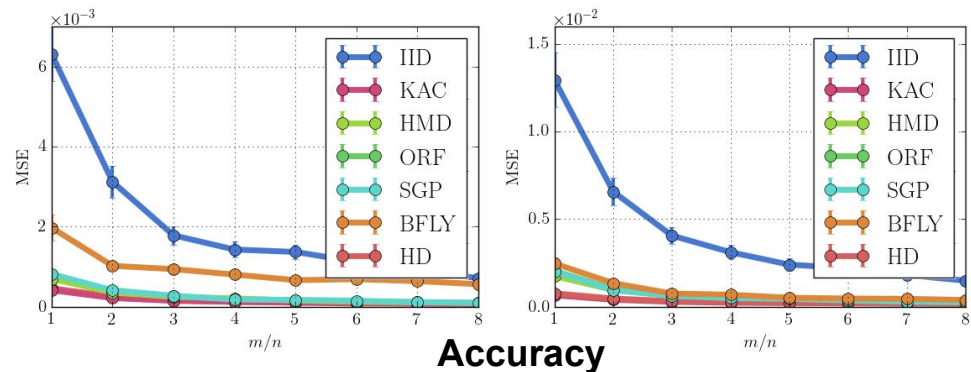
$$\lim_{n \rightarrow \infty} \sup_{X_0 \in S^{n-1}} \|\mathcal{L}(X_{T_2(n)}) - \mu\|_{\text{TV}} = 0.$$



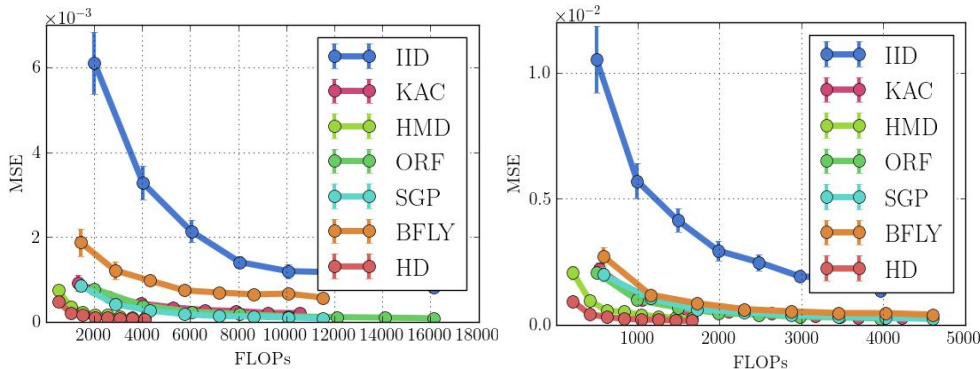
**More careful analysis of the LHS**

# How Does It Work In Practice ?

## Kernel Approximation via Random Features

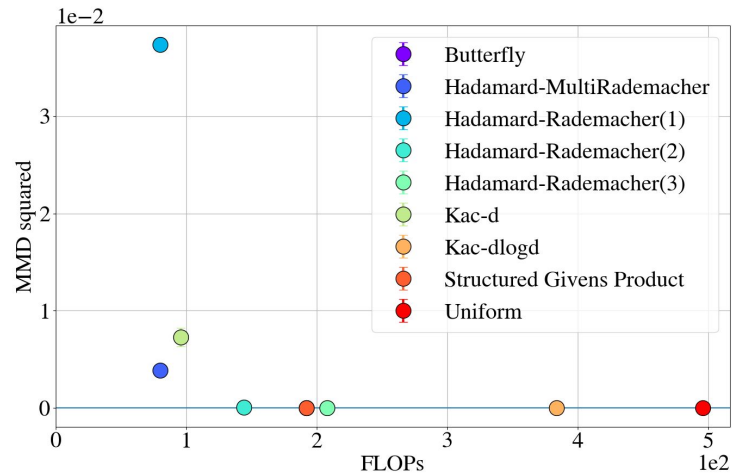


Accuracy

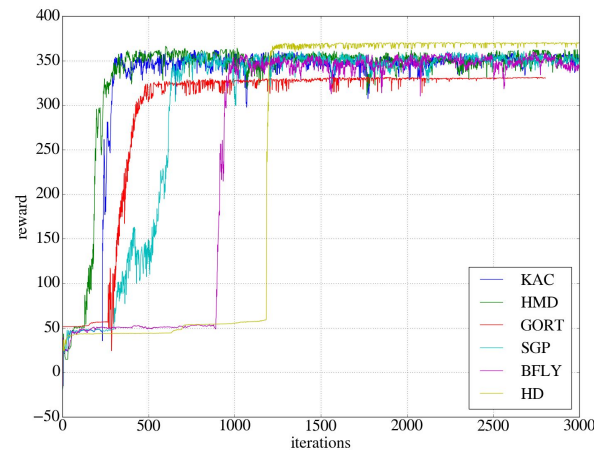


Computational Efficiency

## Maximum Mean Discrepancy Experiment

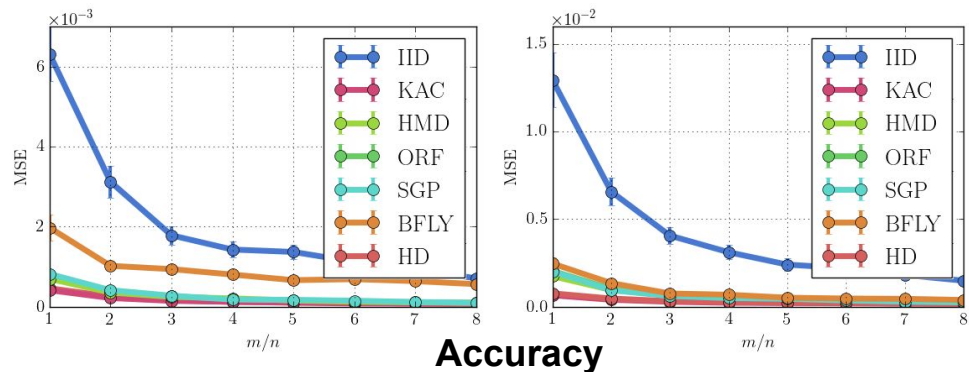


## Reinforcement Learning via ES-methods



# How Does It Work In Practice ?

## Kernel Approximation via Random Features



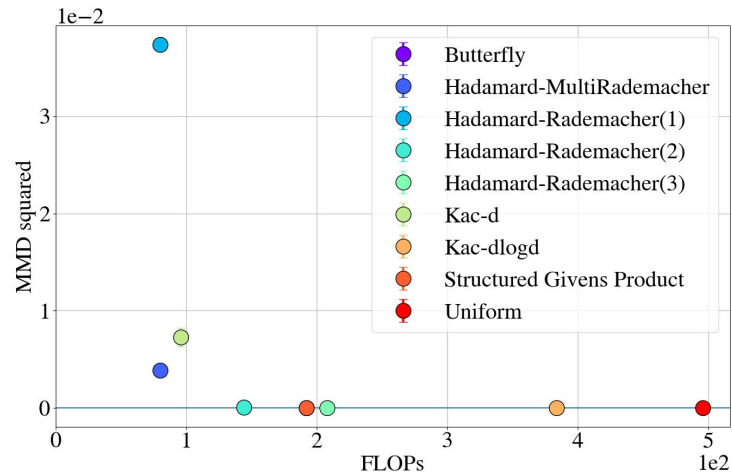
Accuracy



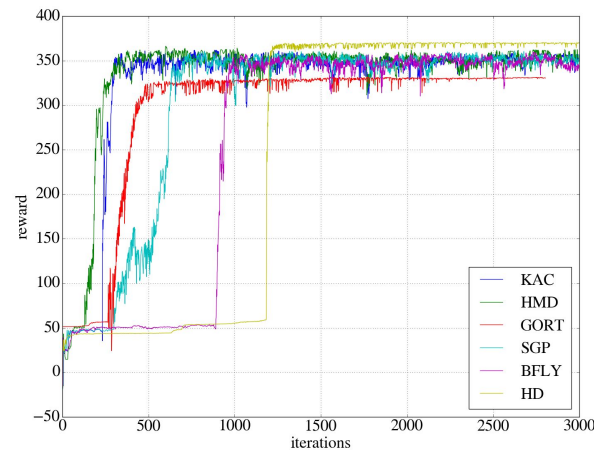
Computational Efficiency

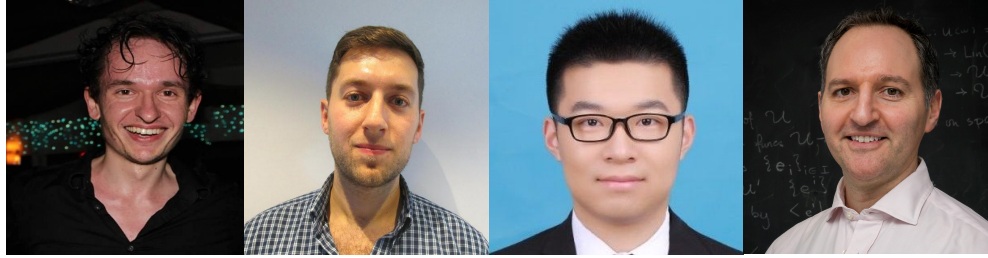
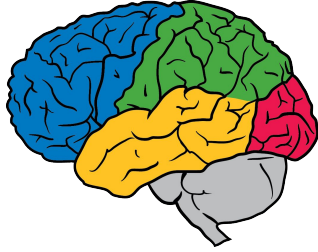


## Maximum Mean Discrepancy Experiment



## Reinforcement Learning via ES-methods





**The  
Alan Turing  
Institute**

**Thank you for your attention !**