

# Data-Driven Linear Complexity Low-Rank Approximation of General Kernel Matrices

Yuanzhe Xi  
Department of Mathematics  
Emory University

Joint work with Difeng Cai (Emory), Edmond Chow (Georgia Tech)

Fast Direct Solver 21, Purdue University

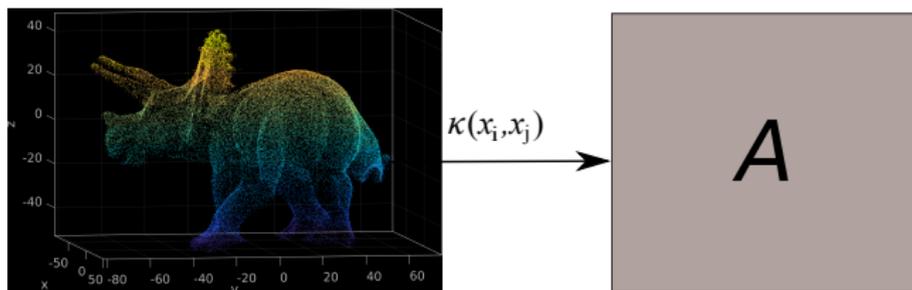
October 23, 2021

Supported by NSF OAC 2003720

# DENSE KERNEL MATRICES

- $X, Y$  are data sets in  $\mathbb{R}^d$

$$A = [\kappa(\mathbf{x}, \mathbf{y})]_{\mathbf{x} \in X, \mathbf{y} \in Y} = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{y}_1) & \kappa(\mathbf{x}_1, \mathbf{y}_2) & \dots & \kappa(\mathbf{x}_1, \mathbf{y}_n) \\ \kappa(\mathbf{x}_2, \mathbf{y}_1) & \kappa(\mathbf{x}_2, \mathbf{y}_2) & \dots & \kappa(\mathbf{x}_2, \mathbf{y}_n) \\ \vdots & \vdots & \vdots & \vdots \\ \kappa(\mathbf{x}_m, \mathbf{y}_1) & \kappa(\mathbf{x}_m, \mathbf{y}_2) & \dots & \kappa(\mathbf{x}_m, \mathbf{y}_n) \end{bmatrix}$$



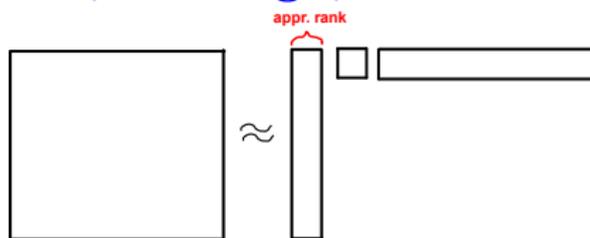
<http://www.geo.tuwien.ac.at/downloads/pg/pctools>

- Canonical kernels  $\kappa(\mathbf{x}, \mathbf{y})$ :

$$\sqrt{1 + \frac{1}{10}|\mathbf{x} - \mathbf{y}|^2}, \frac{1}{|\mathbf{x} - \mathbf{y}|}, \frac{e^{i|\mathbf{x} - \mathbf{y}|}}{|\mathbf{x} - \mathbf{y}|}, e^{-|\mathbf{x} - \mathbf{y}|/\sigma}, e^{-|\mathbf{x} - \mathbf{y}|^2/\sigma^2}, \dots$$

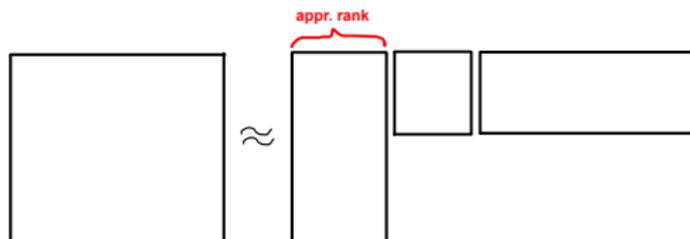
# LOW-RANK APPROXIMATION TECHNIQUES

## 1 Algebraic methods (SVD, RRQR): small rank, high cost



## 2 Analytic methods: Taylor expansion, interpolation [Hackbush, Borm, Le Borne], spherical harmonic [Greengard, Rokhlin]

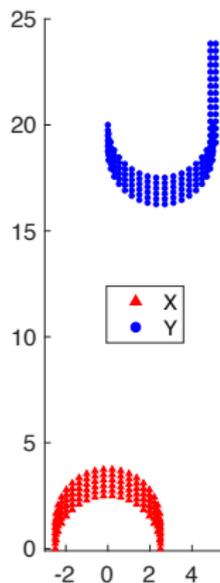
$$\kappa(\mathbf{x}, \mathbf{y}) \approx \sum_{k,l=0}^{r-1} c_{k,l} \phi_k(\mathbf{x}) \psi_l(\mathbf{y}), \quad \forall (\mathbf{x}, \mathbf{y}) \in X \times Y$$



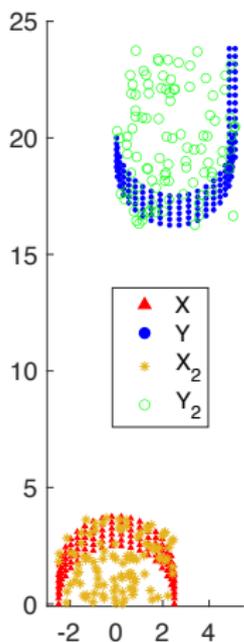
# HYBRID COMPRESSION TECHNIQUES

- 3 **Hybrid compression techniques:** proxy point method, proxy surface method [Ying, Biros, Martinsson, Gillman, Darve, Xing, Chow, et.al.]

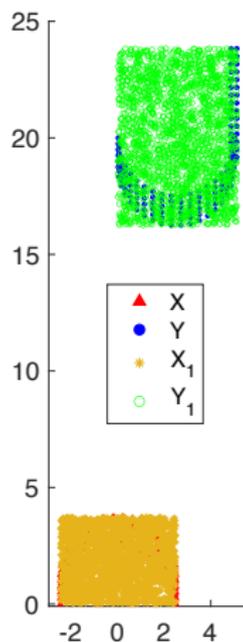
Compressing small kernel matrices from points sampled on an extended continuous domain  $\Omega_X \times \Omega_Y$



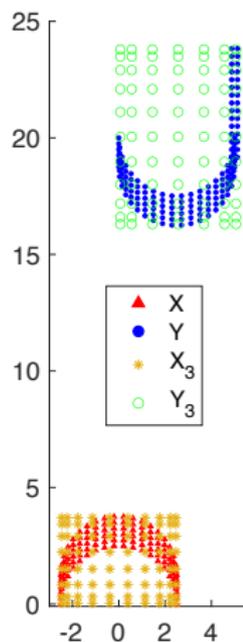
# THREE TYPES OF SAMPLES



(b) 100 rand  
pts



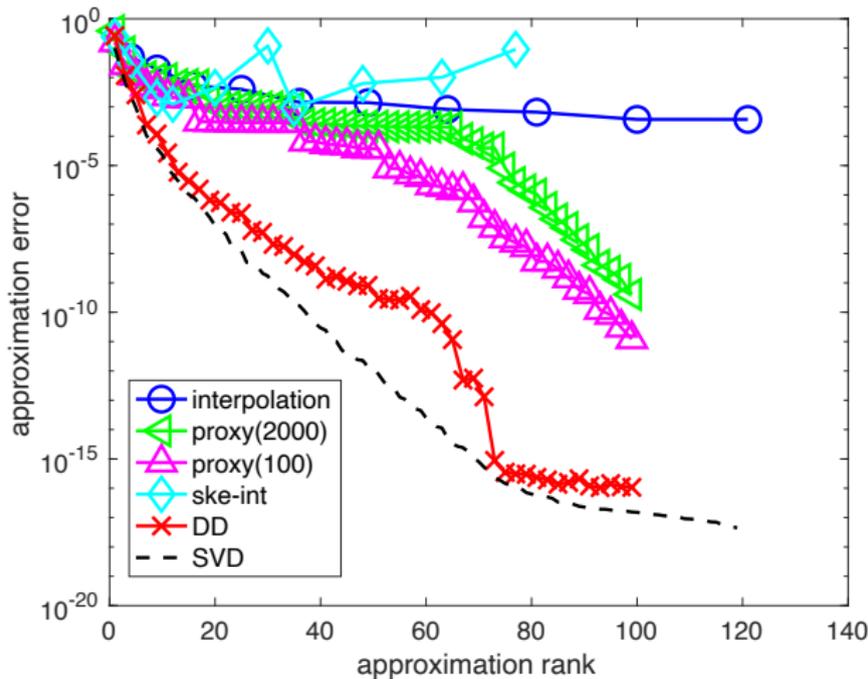
(c) 2000 rand  
pts (proxy  
point, Xing,  
Chow, 2020)



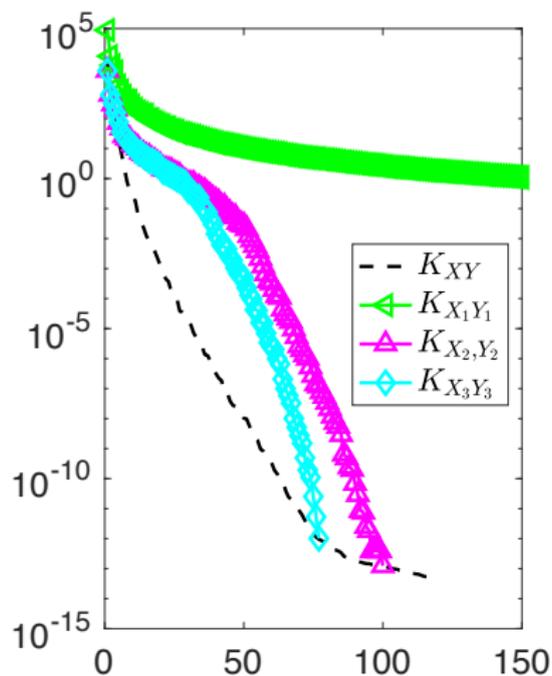
(d)  $10 \times 10$   
Cheby pts  
(Skeletonized  
interpolation,  
Darve, 2020)

# APPROXIMATION ERROR

$$\kappa(\mathbf{x}, \mathbf{y}) = \sqrt{1 + 100|\mathbf{x} - \mathbf{y} + \mathbf{a}|^2} \quad \text{with } \mathbf{a}=[0, 20]^T$$



# NEGATIVE EFFECT OF IGNORING DATA GEOMETRY



Singular values of  $K_{XY}$  (120-by-150),  $K_{X_1Y_1}$  (2000-by-2000),  $K_{X_2Y_2}$  (100-by-100), and  $K_{X_3Y_3}$  (100-by-100)

# SUBSET SELECTION-BASED METHODS

Given  $X \in \mathbb{R}^d$ ,  $Y \in \mathbb{R}^d$  and  $\kappa$ , find a subset  $S_X \in X$  and a subset  $S_Y \in Y$  such that

$$K_{XY} \approx K_{X S_Y} K_{S_X S_Y}^+ K_{S_Y Y}$$

where  $K_{XS} = [\kappa(x, y)]_{x \in X, y \in S}$ .

Similar formats to some existing factorization:

- Randomized CUR/Nyström factorization via uniform, leverage score and k-means based sampling [M. W. Mahoney; P. Drineas, Musco]
- Nearest neighbor-based importance sampling [G. Biros]
- Pseudoskeleton [E. E. Tyrtysnikov, D. Kressner]
- Interpolative decomposition [P.G. Martinsson]
- DEIM [D. Sorensen, M. Embree]

Geometrical understanding of  $S_X$  and  $S_Y$  on approx. accuracy

# ERROR ESTIMATE FOR TWO-SIDED SCHEME

## Theorem (Cai, Chow, XY, 2021)

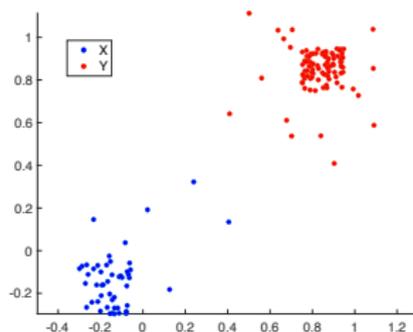
Given  $X$  and  $Y$ ,  $\kappa(x, y) \in C(\mathbb{R}^d \times \mathbb{R}^d)$  is Lipschitz continuous with Lipschitz constant  $L$ , then

$$\|K_{XY} - K_{XS_Y} K_{S_X S_Y}^+ K_{S_X Y}\|_{\max} \leq \sqrt{2}L \max\{\delta_{X, S_X}, \delta_{Y, S_Y}\} + \sqrt{r_2}L\delta_{X, S_X} + \sqrt{r_1}L\delta_{Y, S_Y} + \|K_{S_1 S_2}^+\| \sqrt{r_1 r_2} L^2 \delta_{X, S_X} \delta_{Y, S_Y}$$

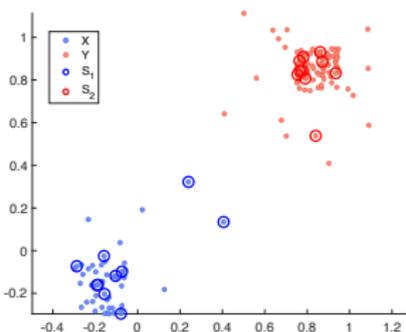
where  $\delta_{X, S_X} = \max_{x \in X} \text{dist}(x, S_X)$  and  $\delta_{Y, S_Y} = \max_{y \in Y} \text{dist}(y, S_Y)$ .

# SUBSET QUALITY INDICATORS

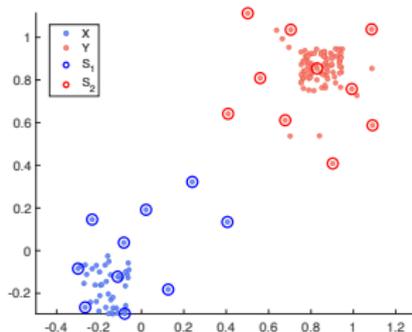
indicator 3 =  $\delta_{X,S_1}$ , indicator 4 =  $\delta_{Y,S_2}$ , indicator 5 =  $\|K_{S_1 S_2}^+\|$



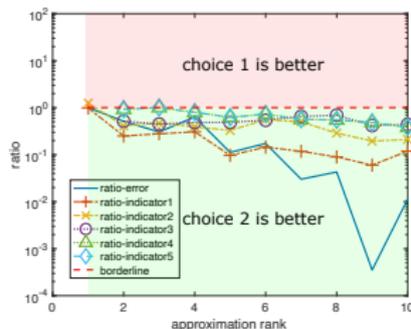
(a) Datasets X, Y



(b) Choice 1



(c) Choice 2



Given  $X \in \mathbb{R}^d$ ,  $Y \in \mathbb{R}^d$  and  $\kappa$ ,

- 1 Select a subset  $S_Y \in Y$
- 2 Perform interpolative decomposition [P.G. Martinsson] to obtain

$$K_{XS_Y} = UK_{IS_Y}$$

- 3 Form  $K_{IY}$  such that

$$K_{XY} \approx UK_{IY}$$

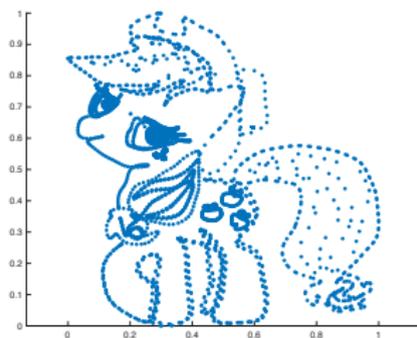
## Theorem

Given  $X$  and  $Y$ ,  $\kappa(x, y) \in C(\mathbb{R}^d \times \mathbb{R}^d)$  is Lipschitz continuous with Lipschitz constant  $L$ , then

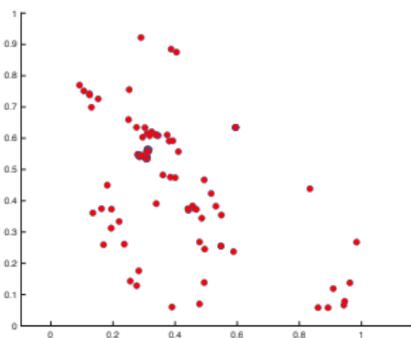
$$\begin{aligned}\|K_{XY} - UK_{\mathcal{I}Y}\|_{\max} &\leq \sqrt{2}L\delta_{Y,S} + (1 + 2r)\sqrt{m}L\delta_{Y,S} \\ &\quad + 2r\sqrt{2}L \max\{\delta_{X,\mathcal{I}}, \delta_{Y,S}\},\end{aligned}$$

where  $m = \text{card}(X)$ ,  $r = \text{card}(\mathcal{I})$ .

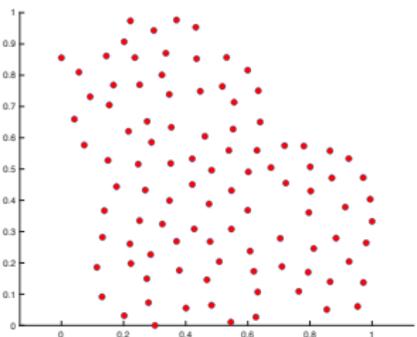
# LINEAR COMPLEXITY SUBSET SELECTION SCHEMES



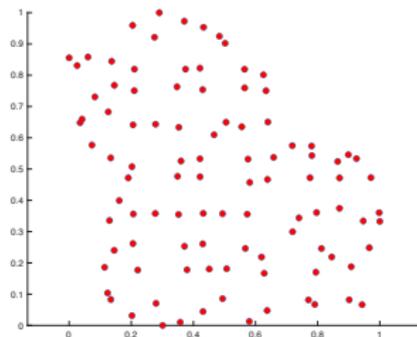
(e) Pony dataset



(f) Uniform sampling



(g) Farthest point sampling



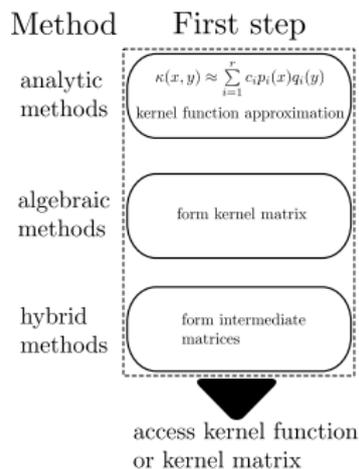
(h) Anchor net method (Cai, Nagy, Xi, 2021)

Select  $r$  points from  $n$  points in  $\mathbb{R}^d$

- 1 Farthest point sampling (**FPS**):  $O(dr^2n)$
- 2 Anchor net sampling (**ANC**):  $O(drn)$

Computational complexity for data-driven methods

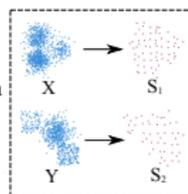
- 1 Two-sided scheme:  $O(dr(m+n))$  (**ANC**), or  $O(dr^2(m+n))$  (**FPS**)
- 2 One-sided scheme:  $O(dr^2(m+n))$



matrix  
compression

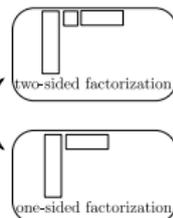
Method First step

data-driven  
method



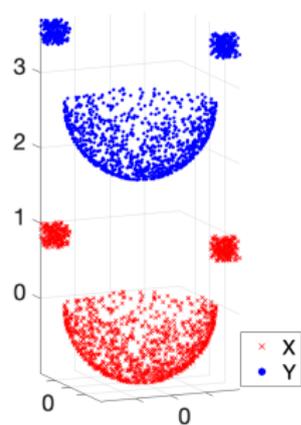
data compression  
(no function or matrix)

matrix  
compression

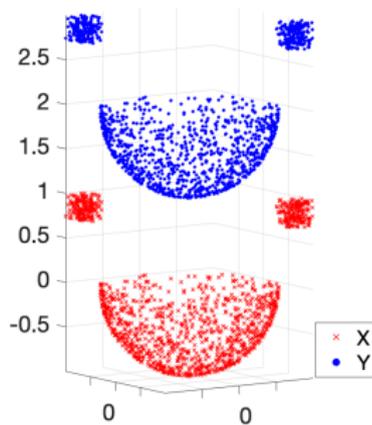


low-rank  
approximation

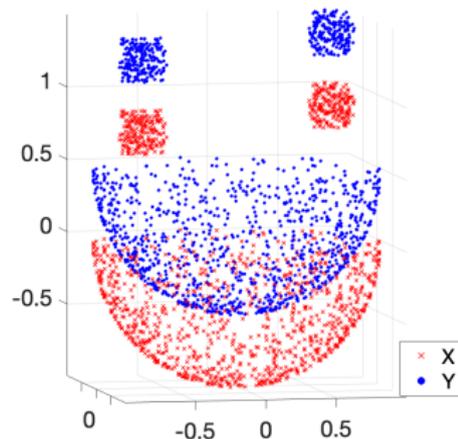
# DATASETS IN THREE DIMENSIONS



(j) Dataset 1



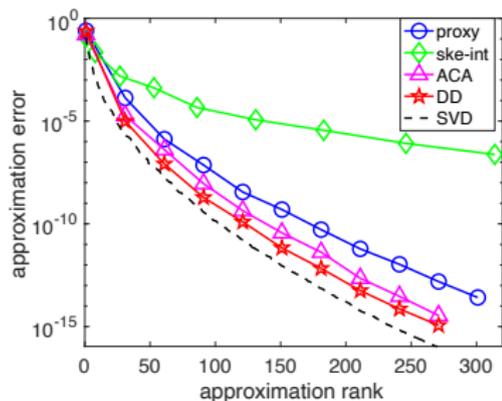
(k) Dataset 2



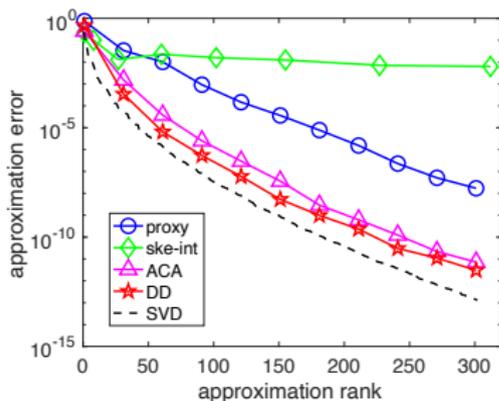
(l) Dataset 3

$Y$ : vertical shift of  $X$  by 2.7, 2.0, 0.5 with  $\text{dist}(X, Y) = 1.00, 0.43, 0.12$

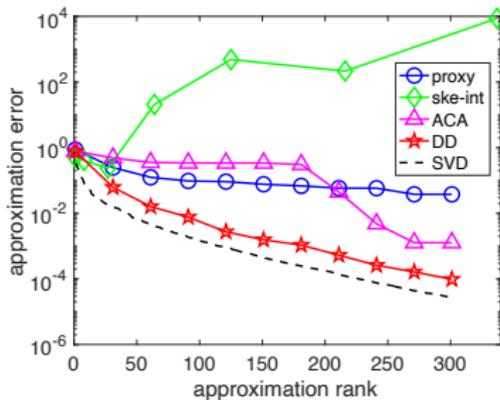
# KERNEL FUNCTION $\kappa(\mathbf{x}, \mathbf{y}) = 1/|\mathbf{x} - \mathbf{y}|$



(a) Dataset 1

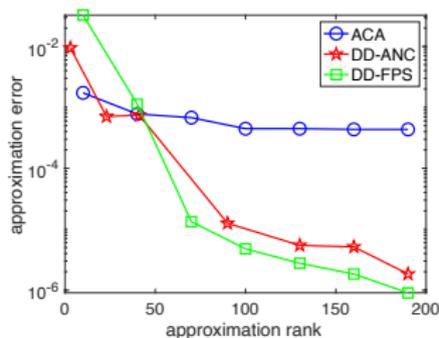


(b) Dataset 2

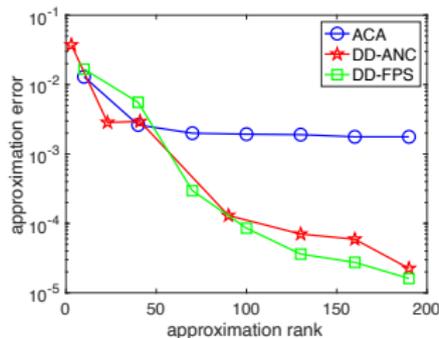


# COVERTYPE ( $d = 54$ )<sup>1</sup>

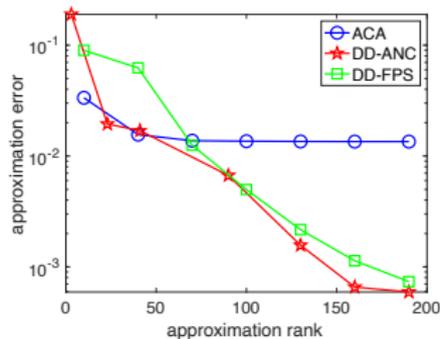
Consider Gaussian kernel  $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-|\mathbf{x} - \mathbf{y}|^2/\sigma^2)$



(a)  $\sigma = \sigma_1 \approx 74.46$

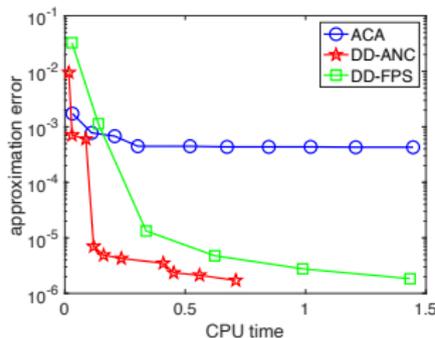


(b)  $\sigma = \sigma_2 \approx 37.23$

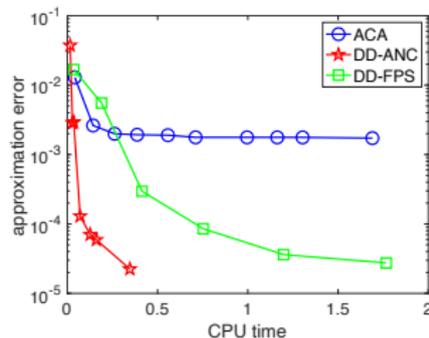


(c)  $\sigma = \sigma_3 \approx 14.89$

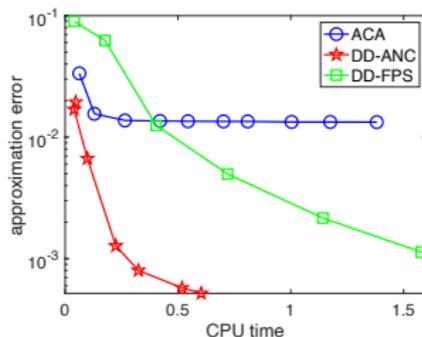
# COVERTYPE ( $d = 54$ )<sup>2</sup>



(d)  $\sigma = \sigma_1 \approx 74.46$

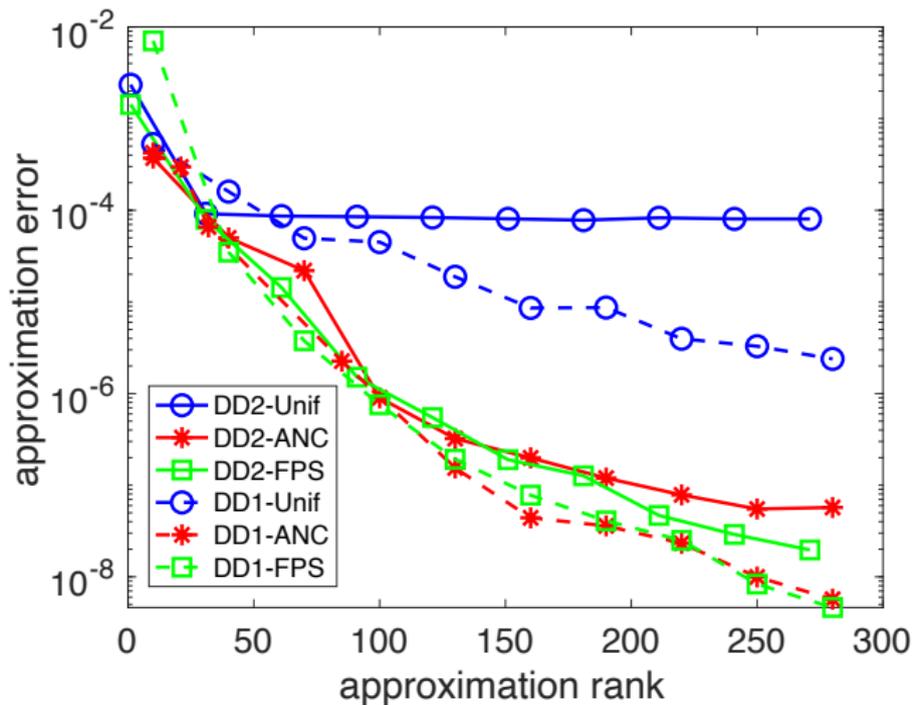


(e)  $\sigma = \sigma_2 \approx 37.23$



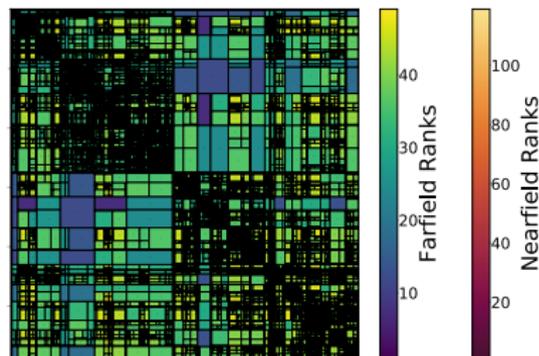
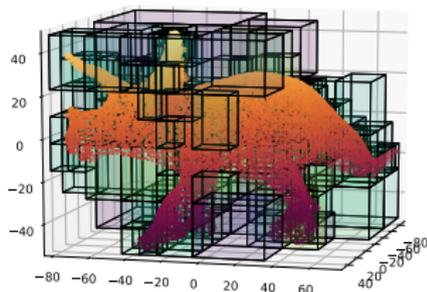
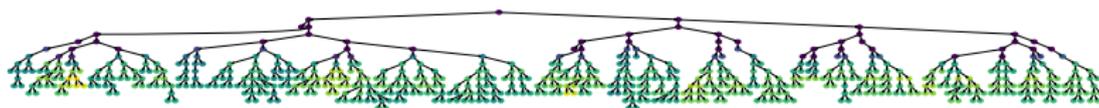
(f)  $\sigma = \sigma_3 \approx 14.89$

# TWO-SIDED V.S. ONE-SIDED



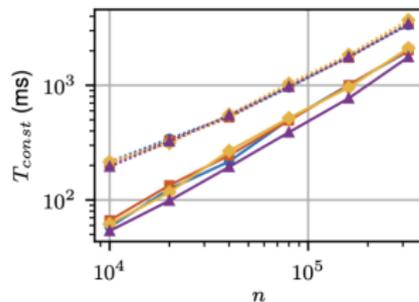
# HIERARCHICAL MATRICES

When kernel has singularity, low-rank properties are associated with certain off-diagonal blocks

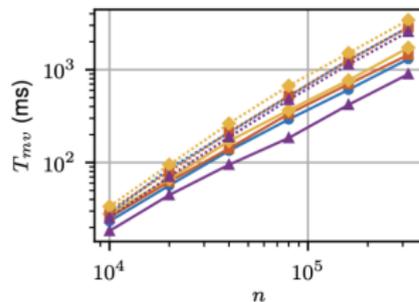


Examples: Fast Multipole Method (Rokhlin 1985, Greengard, Rokhlin 1987),  $\mathcal{H}^2$  matrix (Hackbusch et al. 2001-)

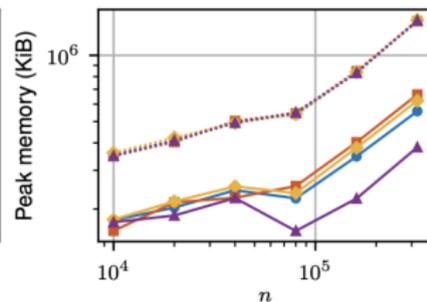
# DATA-DRIVEN VS INTERPOLATION



(a) Construction time (ms)



(b) Matrix-vector product time (ms)

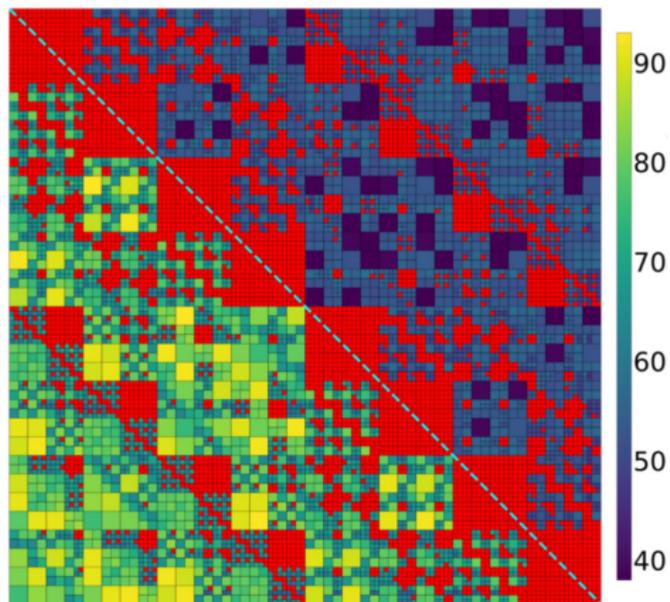


(c) Peak memory usage

DD (solid line) vs Interpolation (dotted line)

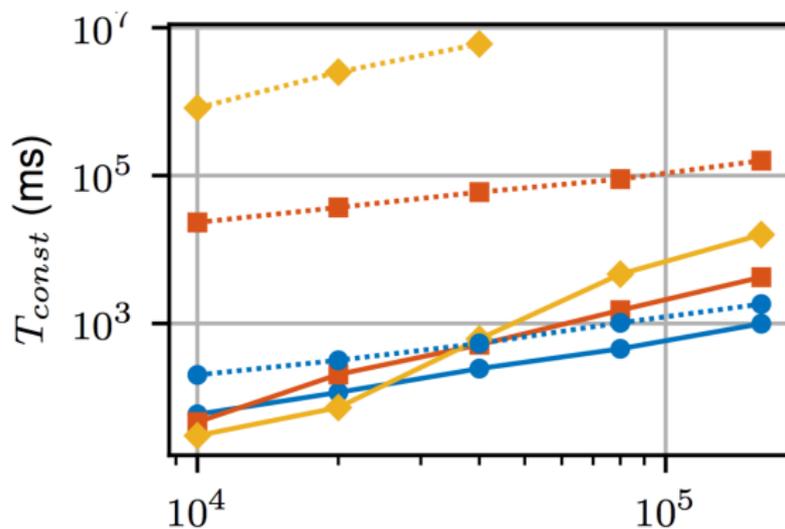
# DATA-DRIVEN VS INTERPOLATION

Coulomb potential, 3D cube, fixed tolerance:  $1E - 7$



Rank comparison. Lower triangular: interpolation; upper triangular: DD

# DATA-DRIVEN SAMPLING VS INTERPOLATION



Matvec in high dimensions

## Data-driven hierarchical matrix methods

- ① Data-driven sampling
- ② General kernel functions
- ③ Efficient hierarchical matrix construction

## Future work

- ① Efficient preconditioning techniques
- ② Tensor kernel functions
- ③ High dimensional problems

- 1 D. Cai, L. Erlandson, E. Chow and YX, Accelerating Parallel Hierarchical Matrix-Vector Products via Data-Driven Sampling, 34th IEEE International Parallel and Distributed Processing Symposium (IPDPS), 2020
- 2 D. Cai, J. Nagy, and YX, Fast deterministic approximation of general symmetric kernel matrices in high dimensions, 2021, arXiv:2102.05215
- 3 D. Cai, E. Chow and YX, Data-driven linear complexity low-rank approximation of general kernel matrices: A Geometric Approach, submitted, 2021

## Definition

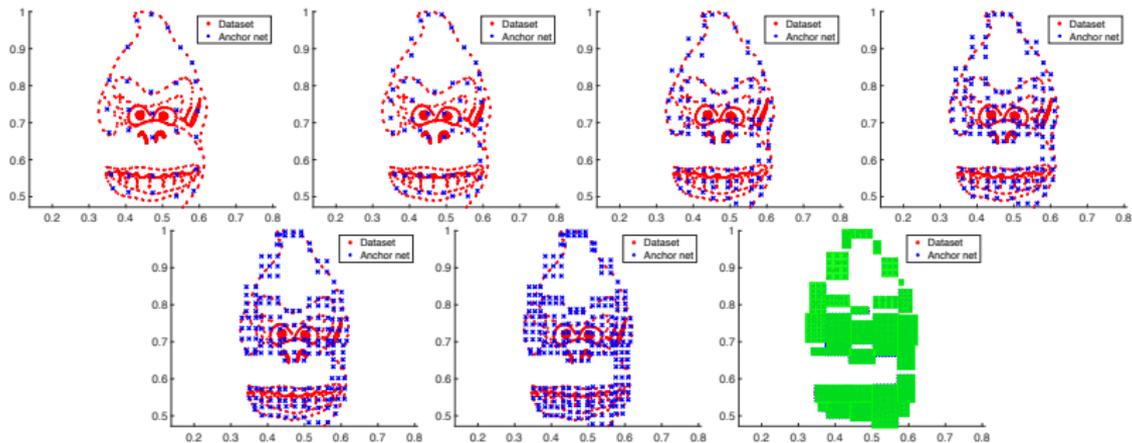
The star discrepancy  $D_N^*(\mathcal{A})$  of  $\mathcal{A} = \{x_1, \dots, x_N\} \subset [0, 1]^d$  is defined by

$$D_N^*(\mathcal{A}) := \sup_{J \in \mathcal{J}_1} |\#(\mathcal{A} \cap J)/N - \lambda(J)|,$$

where  $\mathcal{J}_1$  is the family of all boxes in  $[0, 1]^d$  of the form  $\prod_{i=1}^d [0, \alpha_i)$ .

Examples: Halton sequences, digital nets with  $D_N^*(\mathcal{A}_N) = O(N^{-1}(\log N)^d)$  with  $\mathcal{A}_N$  being the first  $N$  terms.

# ILLUSTRATION OF ANCHOR NETS



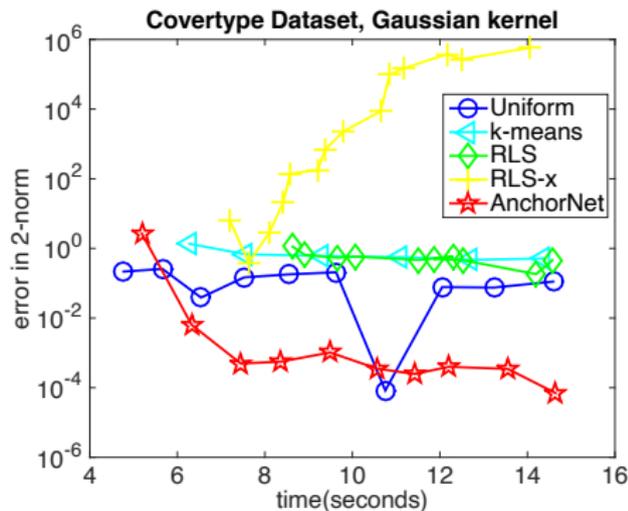
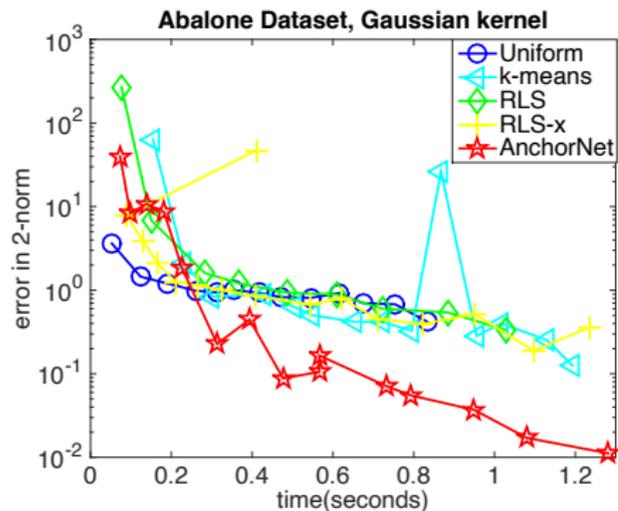
## Property (Anchor net $\mathcal{A}_X$ )

The Anchor net  $\mathcal{A}_X \subset [0, \infty)^d$  parametrized by  $p = 0, 1, 2, \dots$  satisfies

- ① For  $\Omega = \bigcap_{\epsilon > 0} \limsup_{N \rightarrow \infty} \{x \in \mathbb{R}^d : \text{dist}_\infty(x, \mathcal{A}_X) \leq \epsilon\}$ . Then  $\lambda(\Omega) > 0$  and  $X \subset \Omega$ .
- ②  $D_N^*(\mathcal{A}) := \sup_{J \in \mathcal{J}_1} |\#(\mathcal{A} \cap J)/N - \lambda(J)|$ .

1.  $\mathcal{A}_X$  densely surrounds  $X$  as  $p$  increases
2. Uniformity of points in  $\mathcal{A}_X$  is guaranteed

# DATA: ABALONE ( $d = 8$ ), COVERTYPE ( $d = 54$ )<sup>3</sup>



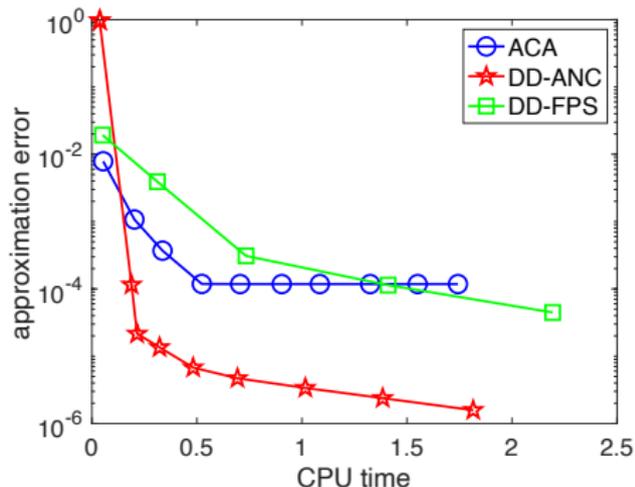
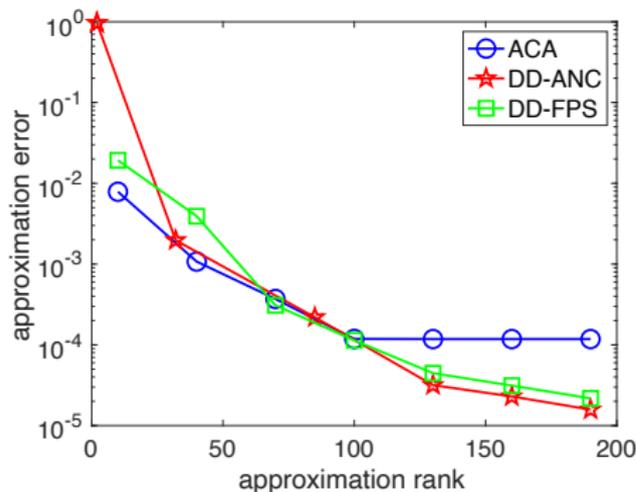
Error-Time plot for approximating Gaussian kernel matrix

**RLS**: Recursive Nystrom method (C. Musco, C. Musco, 2017)

**RLS-x**: Accelerated RLS (Neurips 2017 Spotlights)

# DATA: GAS SENSOR( $d = 128$ )<sup>4</sup>

Randomly sample **8k** points for  $X$  and **10k** points for  $Y$

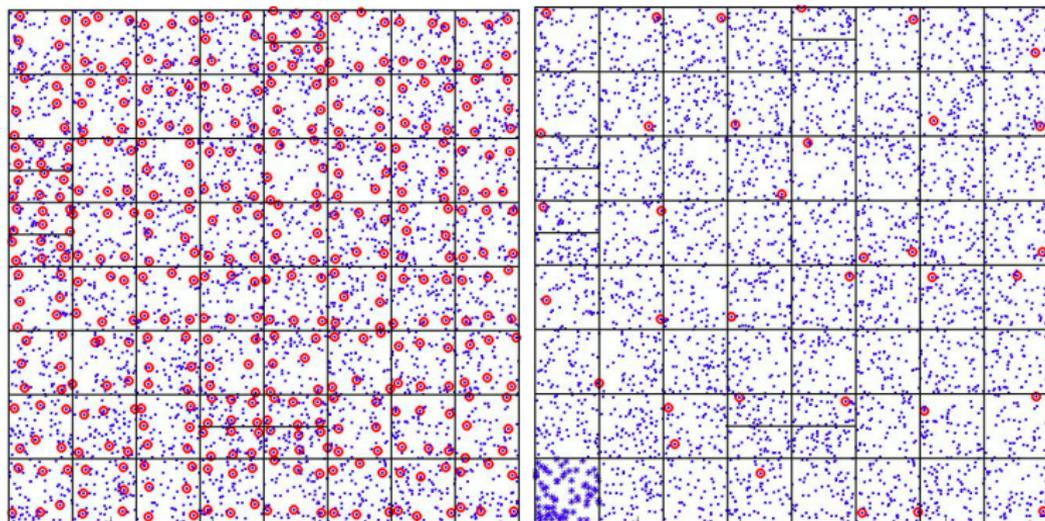


Gaussian kernel with  $\sigma = 14.59$

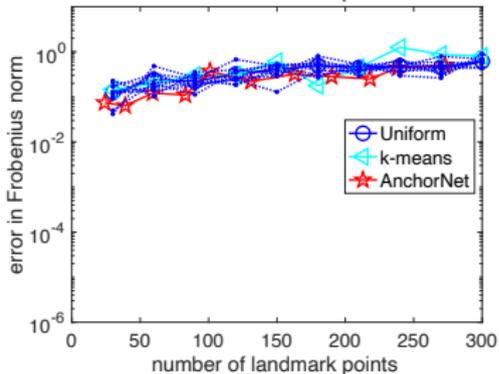
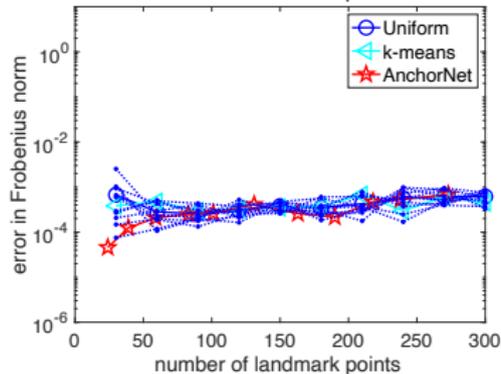
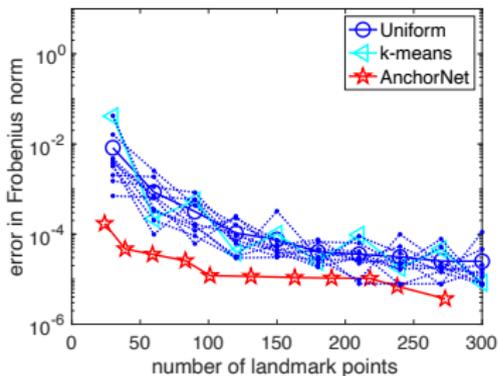
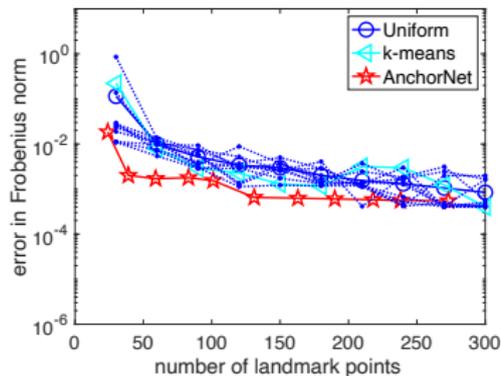
ACA: S. Borm, L. Grasedyck, 2005

# HIERARCHICAL DATA DRIVEN SAMPLING

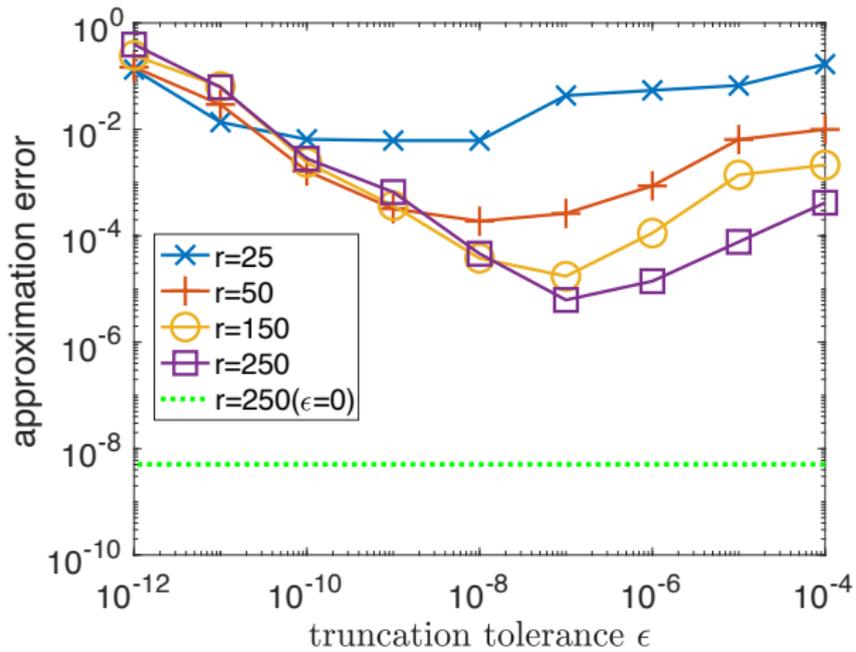
Bottom-up (left) and Top-down sampling (right)



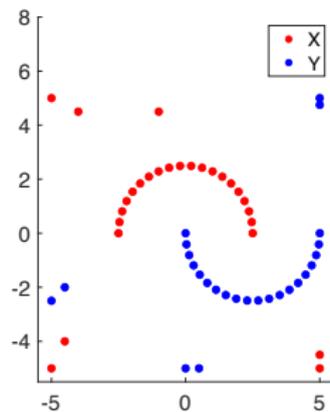
$$\bar{U}_i = [\kappa(\mathbf{x}, \mathbf{y})]_{\mathbf{x} \in X_i, \mathbf{y} \in Y_i} \text{ (Kernel Independent)}$$



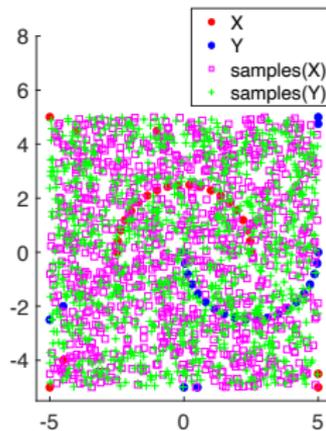
Approximation errors using  $\epsilon$ -pseudoinverse. First four figures are error-rank plots of three methods with  $\epsilon = 10^{-4}, 10^{-6}, 10^{-9}, 10^{-12}$ , respectively.



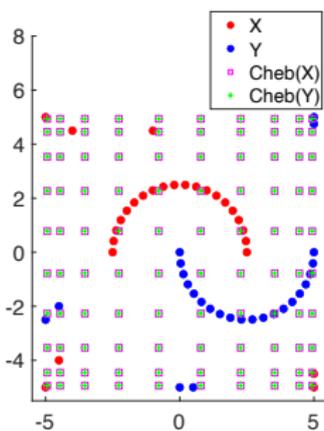
# PROXY POINT APPROXIMATION



(a) Datasets X, Y



(b) Proxy point method (Xing, Chow, 2020)



(c) Skeletonized interpolation (E. Darve 2019)

# SINGULAR VALUE DECAY OF GAUSSIAN MATRICES

