

Near-Optimal Compression In Near-Linear Time

Raaz Dwivedi, Lester Mackey, Abhishek Shetty

raaz@mit.edu



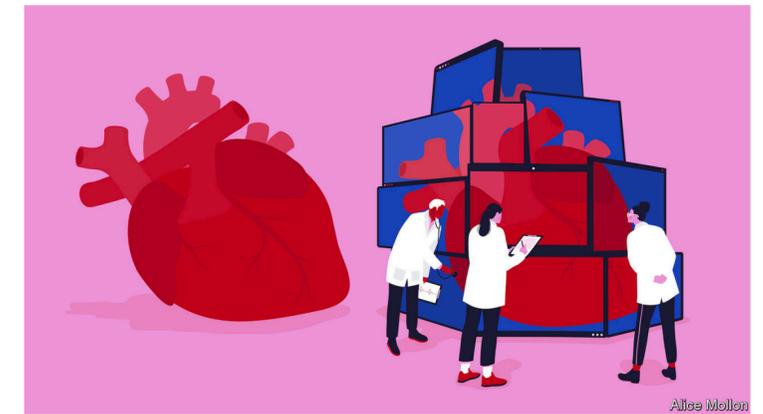
Harvard John A. Paulson
School of Engineering
and Applied Sciences



MSRI Workshop, Berkeley

Mar 8, 2022

Computational cardiology



Modeling **digital twin heart** to predict therapy response in a *non-invasive way* requires single-cell modeling.

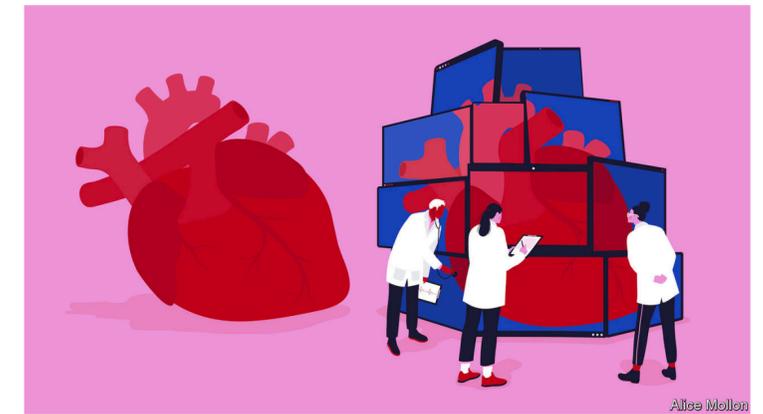
A common methodology:

- Estimate single cell model using Bayesian set-up: Use millions of Markov chain Monte Carlo (MCMC) points to approximate posterior \mathbb{P}^\star
- Propagate uncertainty at heart-level by passing these points to the whole-heart simulator

$$\mathbb{P}^\star f \triangleq \int f(x) d\mathbb{P}^\star(x) \approx \frac{1}{n} \sum_{i=1}^n f(x_i) \triangleq \mathbb{P}_n f$$

x_i = MCMC sample for single cell
model parameters
 f = heart simulator

Computational cardiology



Modeling **digital twin heart** to predict therapy response in a *non-invasive way* requires single-cell modeling.

A common methodology:

- Estimate single cell model using Bayesian set-up: Use millions of Markov chain Monte Carlo (MCMC) points to approximate posterior \mathbb{P}^\star
- Propagate uncertainty at heart-level by passing these points to the whole-heart simulator

$$\mathbb{P}^\star f \triangleq \int f(x) d\mathbb{P}^\star(x) \approx \frac{1}{n} \sum_{i=1}^n f(x_i) \triangleq \mathbb{P}_n f$$

x_i = MCMC sample for single cell model parameters
 f = heart simulator

1 Million MCMC samples ~ 2 weeks
Single evaluation of f ~ 5 weeks

Goal: Represent \mathbb{P}^\star using a few high quality points $(x_i)_{i=1}^n$

Common solutions: I.I.D. sampling, and MCMC sampling exhibit a slow root-n Monte Carlo rate $\left| \mathbb{P}^\star f - \mathbb{P}_n f \right| = \Theta(n^{-1/2})$, e.g., $\sim 10^6$ points for 0.1% error

- Prohibitive for computationally expensive f

Goal: Represent \mathbb{P}^\star using a few high quality points $(x_i)_{i=1}^n$

Common solutions: I.I.D. sampling, and MCMC sampling exhibit a slow root-n Monte Carlo rate $\left| \mathbb{P}^\star f - \mathbb{P}_n f \right| = \Theta(n^{-1/2})$, e.g., $\sim 10^6$ points for 0.1% error

- Prohibitive for computationally expensive f

Data compression: Approximate \mathbb{P}^\star by compressing given n points

- Uniform thinning, or standard thinning--choose every t -th point

Goal: Represent \mathbb{P}^\star using a few high quality points $(x_i)_{i=1}^n$

Common solutions: I.I.D. sampling, and MCMC sampling exhibit a slow root-n Monte Carlo rate $|\mathbb{P}^\star f - \mathbb{P}_n f| = \Theta(n^{-1/2})$, e.g., $\sim 10^6$ points for 0.1% error

- Prohibitive for computationally expensive f

Data compression: Approximate \mathbb{P}^\star by compressing given n points

- Uniform thinning, or standard thinning--choose every t -th point
- Accuracy degrades with such thinning-- $\Theta(\sqrt{t/n})$ worst-case error--same as the Monte Carlo rate with n/t points; e.g., $n^{-1/4}$ rate with \sqrt{n} points

Goal: Represent \mathbb{P}^\star using a few high quality points $(x_i)_{i=1}^n$

Common solutions: I.I.D. sampling, and MCMC sampling exhibit a slow root-n Monte Carlo rate $|\mathbb{P}^\star f - \mathbb{P}_n f| = \Theta(n^{-1/2})$, e.g., $\sim 10^6$ points for 0.1% error

- Prohibitive for computationally expensive f

Data compression: Approximate \mathbb{P}^\star by compressing given n points

- Uniform thinning, or standard thinning--choose every t -th point
- Accuracy degrades with such thinning-- $\Theta(\sqrt{t/n})$ worst-case error--same as the Monte Carlo rate with n/t points; e.g., $n^{-1/4}$ rate with \sqrt{n} points

Can we do better?

Minimax lower bounds

There exists some \mathbb{P}^\star such that the worst-case integration error

- Is $\Omega(n^{-1/2})$ for **any compression scheme** returning \sqrt{n} points
[Philips and Tai, 2020]

Minimax lower bounds

There exists some \mathbb{P}^\star such that the worst-case integration error

- Is $\Omega(n^{-1/2})$ for **any compression scheme** returning \sqrt{n} points
[Philips and Tai, 2020]
- Is $\Omega(n^{-1/2})$ for **any approximation** based on n i.i.d. points
[Tolstikhin, Sriperumbudur, and Muandet, 2017]

This talk: **Kernel thinning-Compress++**

- **KT-Compress++**: A practical strategy based on two new algorithms to provide near-optimal compression in near-linear time
 - **Kernel thinning (KT)**: Provides **near-optimal compression**
 - **Compress++**: Provides **significantly reduced runtime** for generic thinning algorithms with **minimal worsening of error**

This talk: **Kernel thinning-Compress++**

- **KT-Compress++**: A practical strategy based on two new algorithms to provide near-optimal compression in near-linear time
 - **Kernel thinning (KT)**: Provides **near-optimal compression**
 - **Compress++**: Provides **significantly reduced runtime** for generic thinning algorithms with **minimal worsening of error**
- Overall, KT-Compress++ a solution for finding
 - better than Monte Carlo points
 - high quality coresets
 - good prototypes

Problem set-up

- **Input:**

Points $(x_i)_{i=1}^n$ with empirical distribution $\mathbb{P}_{in} \triangleq \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Target output size s ($s = \sqrt{n}$ for heavy compression)

- **Goal:**

Return a subset of input points with size s , empirical distribution \mathbb{P}_{out} with error rate $o(s^{-1/2})$, i.e., better than Monte Carlo rate

Reproducing kernel Hilbert space (RKHS)

- RKHS of \mathbf{k} is given by $\mathbb{H}_{\mathbf{k}} \triangleq \overline{\text{span}\{\mathbf{k}(x, \cdot), x \in \mathcal{X}\}}$
- $\mathbb{H}_{\mathbf{k}}$ is dense in the space of continuous functions for universal \mathbf{k} like

$$\text{Gaussian } \mathbf{k}(x, y) = \exp\left(-\frac{1}{2}\|x - y\|^2\right); \text{IMQ } \mathbf{k}(x, y) = \frac{1}{(1 + \|x - y\|^2)^{1/2}}$$

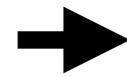
$\mathbf{k} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a reproducing kernel if the matrix $\mathbf{K} = (\mathbf{k}(x_i, x_j))_{i,j=1}^n$ is a symmetric positive definite matrix for any n and any (x_1, \dots, x_n)

Kernel Thinning

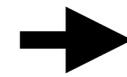
x_1, x_2, \dots, x_n

kernel \mathbf{k}

$$\mathbb{P}_{in} \triangleq \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$



Kernel
Thinning
(KT)



Non-uniform sub-sample of size s

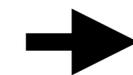
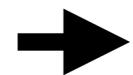
$$y_1, \dots, y_s$$
$$\mathbb{P}_{KT} \triangleq \frac{1}{s} \sum_{i=1}^s \delta_{y_i}$$

Kernel Thinning: A two-staged procedure

x_1, x_2, \dots, x_n

kernel \mathbf{k}

$$\mathbb{P}_{in} \triangleq \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$

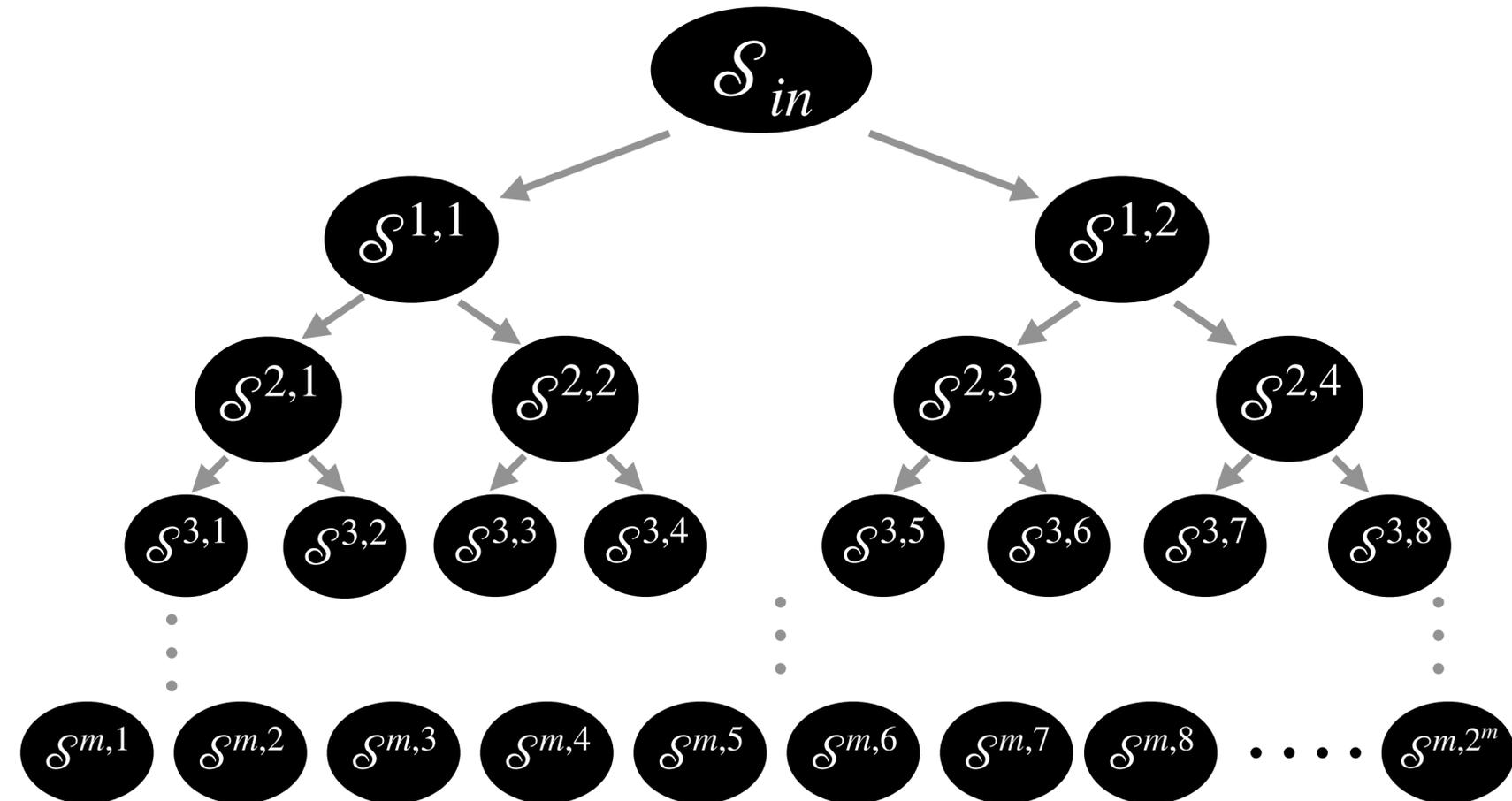


Non-uniform sub-sample of size s

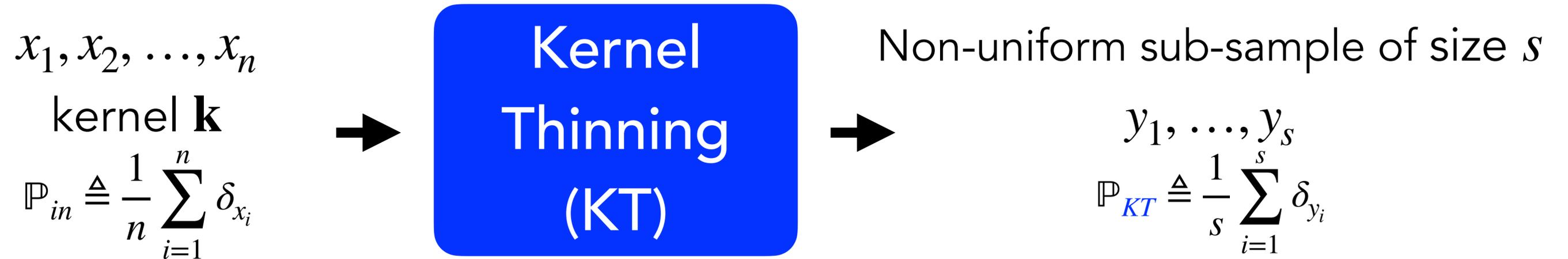
y_1, \dots, y_s

$$\mathbb{P}_{KT} \triangleq \frac{1}{s} \sum_{i=1}^s \delta_{y_i}$$

- Stage 1:** $m = \frac{1}{2} \log_2(n/s)$ recursive rounds of **non-uniform splitting** the parent coreset in two equal-sized children coresets

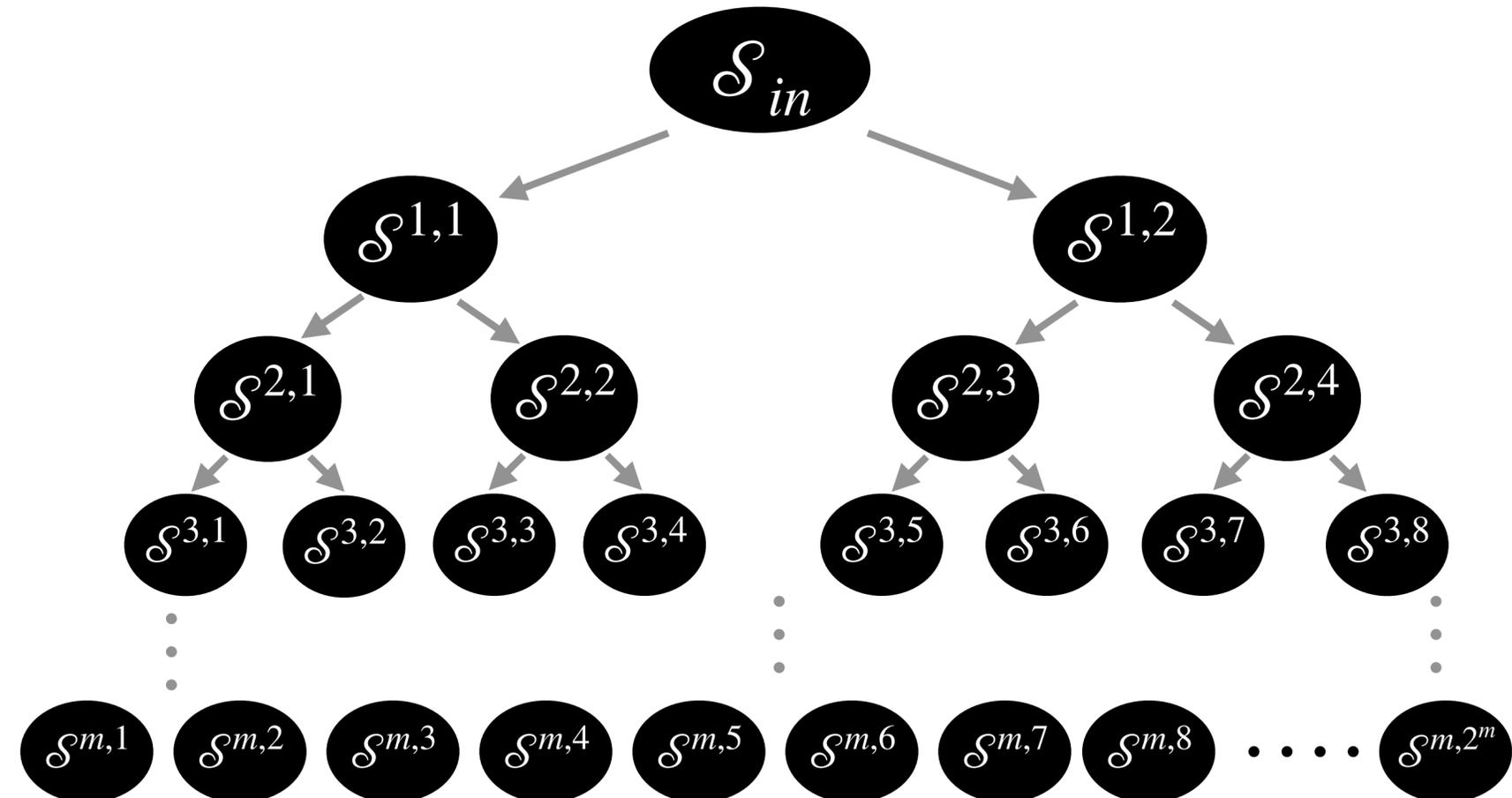


Kernel Thinning: A two-staged procedure



- Stage 1:** $m = \frac{1}{2} \log_2(n/s)$ recursive rounds of **non-uniform splitting** the parent coresets in two equal-sized children coresets

- Stage 2:** Point-by-point **refinement** of the best child coresets



Kernel Thinning: Better than Monte Carlo rate for \mathbb{P}_{in}

With n input points, s output points, with high probability over the randomness in KT, **for any fixed** $g \in \mathbb{H}_{\mathbf{k}}$ we have

$$|\mathbb{P}_{in}g - \mathbb{P}_{KT}g| \lesssim \frac{1}{s} \cdot \|g\|_{\mathbf{k}} \sqrt{\|\mathbf{k}\|_{\infty} (\log s + \log \log(n/s))} \ll \mathcal{O}\left(\frac{1}{\sqrt{s}}\right) \begin{array}{l} \text{Monte Carlo rate} \\ \text{(standard thinning rate)} \end{array}$$

for any kernel on any space!

Kernel Thinning: Better than Monte Carlo rate for \mathbb{P}_{in}

With n input points, \sqrt{n} output points, with high probability over the randomness in KT, **for any fixed** $g \in \mathbb{H}_{\mathbf{k}}$ we have

$$|\mathbb{P}_{in}g - \mathbb{P}_{KT}g| \lesssim \frac{1}{\sqrt{n}} \cdot \|g\|_{\mathbf{k}} \sqrt{\|\mathbf{k}\|_{\infty} \log n} \ll \mathcal{O}\left(\frac{1}{n^{1/4}}\right) \quad \begin{array}{l} \text{Monte Carlo rate} \\ \text{(standard thinning rate)} \end{array}$$

for any kernel on any space!

Kernel Thinning: Better than Monte Carlo rate for \mathbb{P}_{in}

With n input points, \sqrt{n} output points, with high probability over the randomness in KT, **for any fixed** $g \in \mathbb{H}_{\mathbf{k}}$ we have

$$|\mathbb{P}_{in}g - \mathbb{P}_{KT}g| \lesssim \frac{1}{\sqrt{n}} \cdot \|g\|_{\mathbf{k}} \sqrt{\|\mathbf{k}\|_{\infty} \log n} \ll \mathcal{O}\left(\frac{1}{n^{1/4}}\right) \quad \begin{array}{l} \text{Monte Carlo rate} \\ \text{(standard thinning rate)} \end{array}$$

$$\text{If } |\mathbb{P}^{\star}g - \mathbb{P}_{in}g| = \mathcal{O}\left(\sqrt{\frac{\log n}{n}}\right) \text{ then } |\mathbb{P}^{\star}g - \mathbb{P}_{KT}g| = \mathcal{O}\left(\sqrt{\frac{\log n}{n}}\right)$$

Monte Carlo input + KT \Rightarrow Better than Monte Carlo output for \mathbb{P}^\star

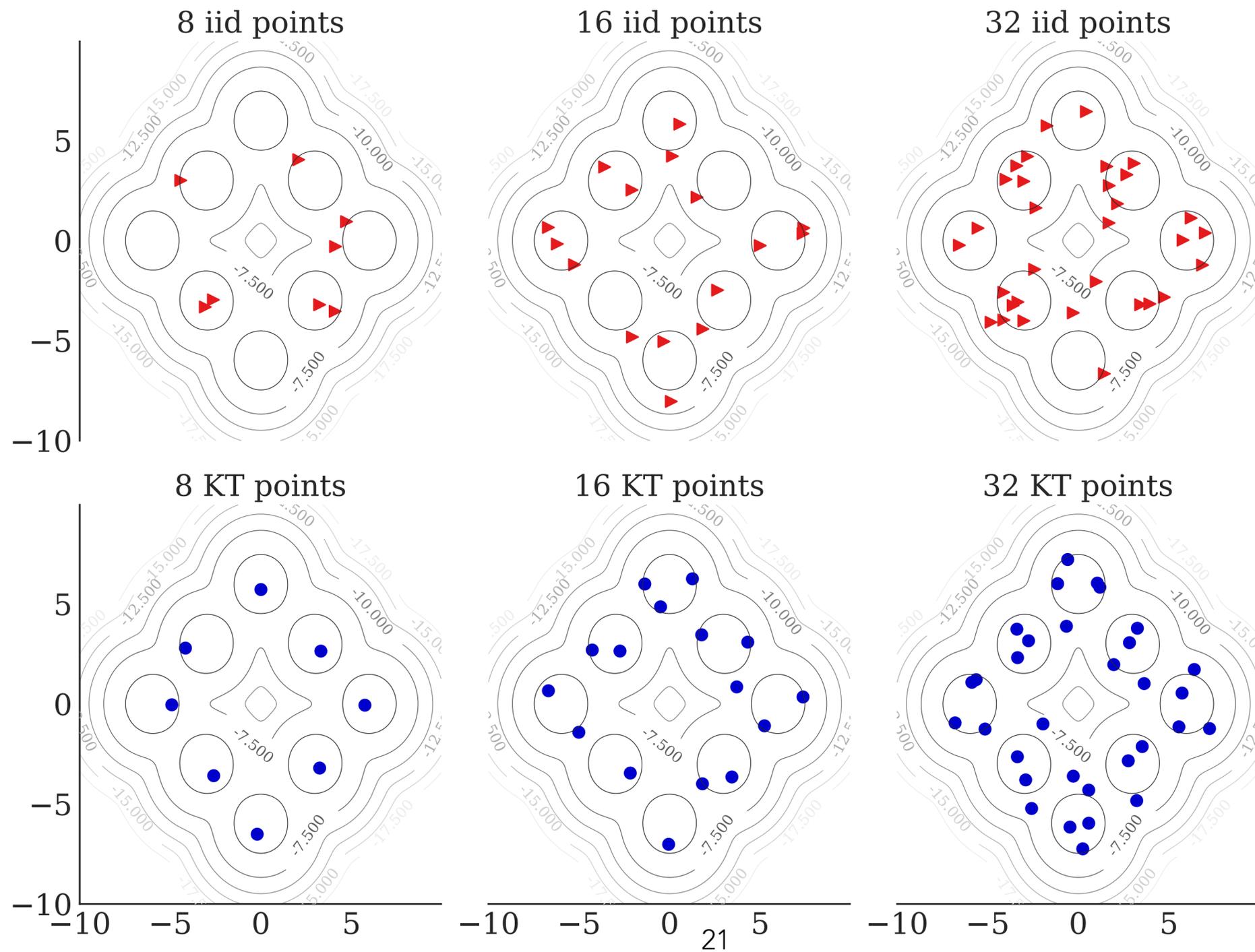
With n input points, \sqrt{n} output points, with high probability over the randomness in KT, **for any fixed** $g \in \mathbb{H}_k$ we have

$$|\mathbb{P}_{in}g - \mathbb{P}_{KT}g| \lesssim \frac{1}{\sqrt{n}} \cdot \|g\|_k \sqrt{\|\mathbf{k}\|_\infty \log n} \ll \mathcal{O}\left(\frac{1}{n^{1/4}}\right) \quad \begin{array}{l} \text{Monte Carlo rate} \\ \text{(standard thinning rate)} \end{array}$$

$$\text{If } |\mathbb{P}^\star g - \mathbb{P}_{in}g| = \mathcal{O}\left(\sqrt{\frac{\log n}{n}}\right) \text{ then } |\mathbb{P}^\star g - \mathbb{P}_{KT}g| = \mathcal{O}\left(\sqrt{\frac{\log n}{n}}\right)$$

- This is the Monte Carlo rate for input!
- Easily satisfied for input from iid sampling, MCMC, quadrature methods ...

Intuition: KT finds “diverse and representative” points



Worst-case error

- For points in \mathbb{R}^d , **the worst-case error**—Maximum Mean Discrepancy (MMD) error—in the reproducing kernel Hilbert space (RKHS) satisfies

$$\sup_{\|g\|_k \leq 1} |\mathbb{P}_{in} g - \mathbb{P}_{KT} g|$$

Worst-case error: $\widetilde{O}(n^{-1/2})$ with \sqrt{n} points for decaying \mathbf{k}

- For points in \mathbb{R}^d , **the worst-case error**—Maximum Mean Discrepancy (MMD) error—in the reproducing kernel Hilbert space (RKHS) satisfies

$$\sup_{\|g\|_{\mathbf{k}} \leq 1} |\mathbb{P}_{in}g - \mathbb{P}_{KT}g| \lesssim_d \begin{cases} n^{-1/2} \sqrt{\log n} & \text{(Compactly supported; e.g., B-spline } \mathbf{k}) \\ n^{-1/2} \sqrt{\log^{d/2+1} n \log \log n} & \text{(Sub-Gaussian tails; e.g., Gaussian } \mathbf{k}) \\ n^{-1/2} \sqrt{\log^{d+1} n \log \log n} & \text{(Sub-exponential tails; e.g., Matern } \mathbf{k}) \end{cases}$$

Assuming similar tails for \mathbf{P}_{in}

- For output size s , the MMD error is $\widetilde{O}(1/s)$

Monte Carlo input + KT \Rightarrow Better than Monte Carlo output for \mathbb{P}^\star

- For points in \mathbb{R}^d , **the worst-case error**—Maximum Mean Discrepancy (MMD) error—in the reproducing kernel Hilbert space (RKHS) satisfies

$$\sup_{\|g\|_{\mathbf{k}} \leq 1} |\mathbb{P}^\star g - \mathbb{P}_{KT} g| \lesssim_d \begin{cases} n^{-1/2} \sqrt{\log n} & \text{(Compactly supported; e.g., B-spline } \mathbf{k}) \\ n^{-1/2} \sqrt{\log^{d/2+1} n \log \log n} & \text{(Sub-Gaussian tails; e.g., Gaussian } \mathbf{k}) \\ n^{-1/2} \sqrt{\log^{d+1} n \log \log n} & \text{(Sub-exponential tails; e.g., Matern } \mathbf{k}) \end{cases}$$

Assuming similar tails for \mathbf{P}_{in}

whenever $\sup_{\|g\|_{\mathbf{k}} \leq 1} |\mathbb{P}^\star g - \mathbb{P}_{in} g| \lesssim n^{-1/2}$ —holds for iid / fast mixing MCMC input

Comparison with related work

- finding good approximations to \mathbb{P}^\star by thinning, reweighting or directly

Related work: \sqrt{n} points with $\mathcal{O}(n^{-1/4})$ MMD

- **Known guarantees *no better than Monte Carlo rate*:**

Standard thinning iid points [Tolstikhin-Sriperumbudur-Muandet, 2017]

Standard thinning geometrically ergodic MCMC [Dwivedi-Mackey 2021]

Kernel herding for infinite-dimensional kernels [Chen-Welling-Smola 2010, Lacoste-Julien-Lindsten-Bach 2015]

Stein Points MCMC [Chen-Barp-Briol-Gorham-Girolami-Mackey-Oates, 2019]

Greedy sign selection [Karnin-Liberty 2019]

- ***Unknown guarantees*:**

Support points [Mak-Joseph 2018]

Supersampling from a reservoir [Paige-Sejdinovic-Wood, 2016]:

Related work: \sqrt{n} points with $o(n^{-1/4})$ MMD

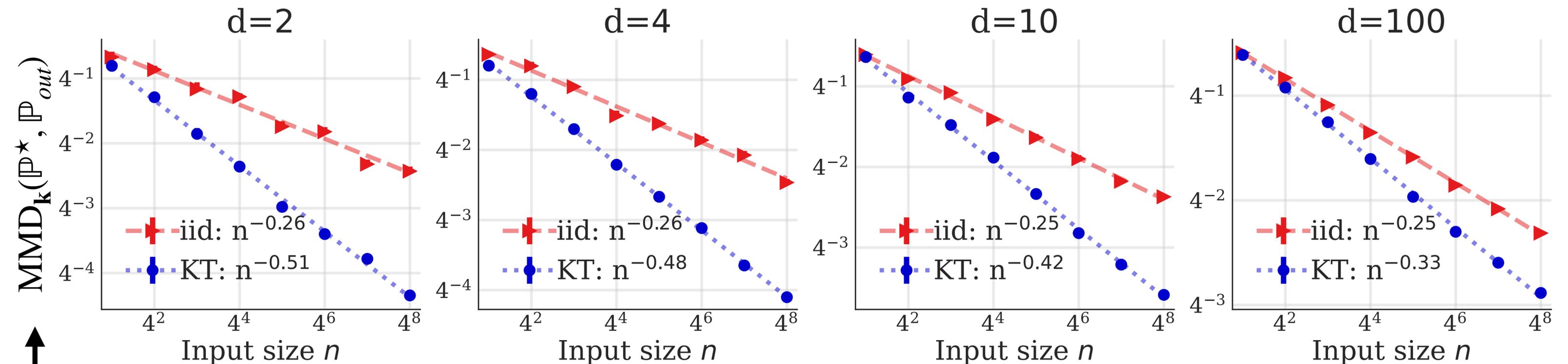
- **Finite-dimensional linear kernels:** Discrepancy construction [Harvey and Samadi, 2014]
- **Uniform \mathbb{P}^* on $[0,1]^d$:**
Quasi Monte Carlo [Hickernell 1998, Novak-Wozniakowski 2010],
Haar thinning [Dwivedi-Feldheim-Gurel-Gurevich-Ramdas 2019]
- **\mathbb{P}^* with *bounded support* with known $\mathbb{P}^* \mathbf{k}$:**
Bayesian quadrature [O'Hagan 1991]
Bayes' Sard cubature [Karvonen et al. 2018]
Determinantal point processes [Belhadji et al. 2020]
- **$(\mathbf{k}, \mathbb{P}^*)$ with *known/bounded eigenfunctions*:**
Determinantal point process kernel quadrature [Belhadji et al. 2019]
Black-box importance sampling [Liu et al. 2018]

Kernel thinning advantages

1. \sqrt{n} points with $O(\sqrt{\log n/n})$ integration-error for any fixed function in the RKHS for any kernel on any space (iid sampling gives $\Omega(n^{-1/4})$ error)
2. \sqrt{n} points with $\tilde{O}(n^{-1/2})$ -worst-case/MMD error for decaying kernels
3. Valid for **non-uniform** target distributions with **unbounded support**
4. Valid for **infinite-dimensional** smooth/decaying kernels
5. Valid for **generic input points** including iid/MCMC/quadrature etc with mild conditions
6. Requires **only kernel evaluations** to implement
7. **Matches MMD lower bounds** up to log factors
8. **Matches L^∞ -error lower bounds** up to log factors

KT vs iid: Gaussian \mathbb{P}^\star in \mathbb{R}^d

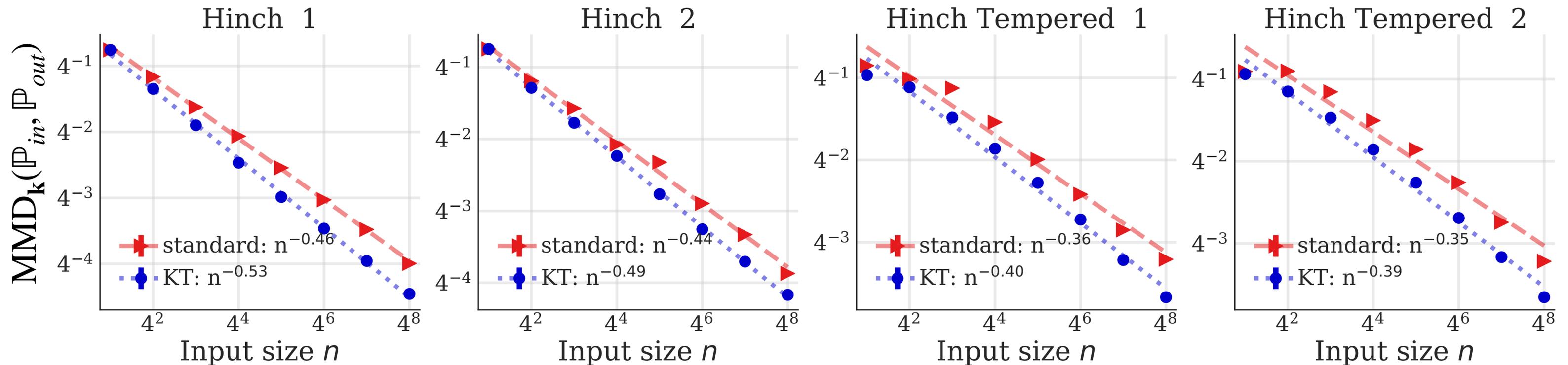
(Gaussian kernel with $\sigma^2 = 2d$)



In practice, significant gains even in dimension $d = 100$

(Worst-case error in the unit ball of Gaussian RKHS)

KT on MCMC samples from computational cardiology



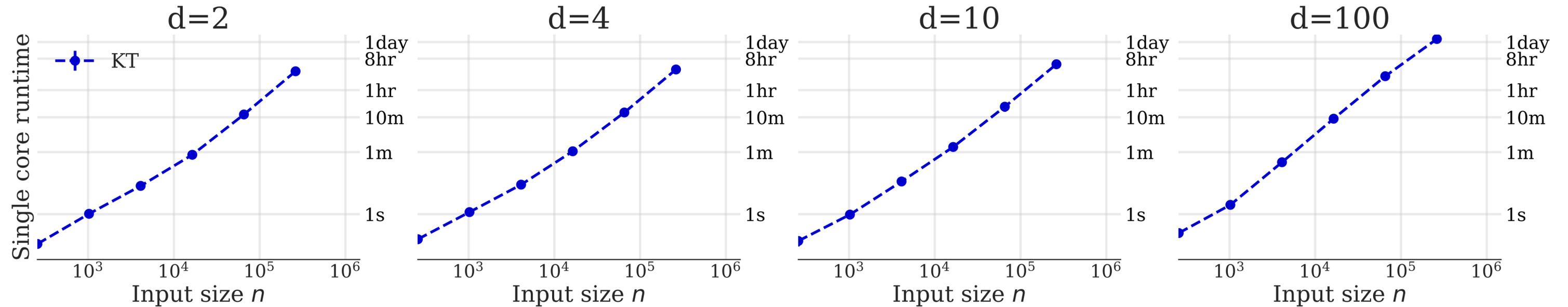
In this setting with $d = 38$, standard thinning is already good (the chain is mixing slowly), but **KT provides further improvement!** Each point saves 1000s of CPU hours!!

*MCMC samples taken from Riabiz-Chen-Cockayne-Swietach-Niederer-Mackey-Oates, 2021

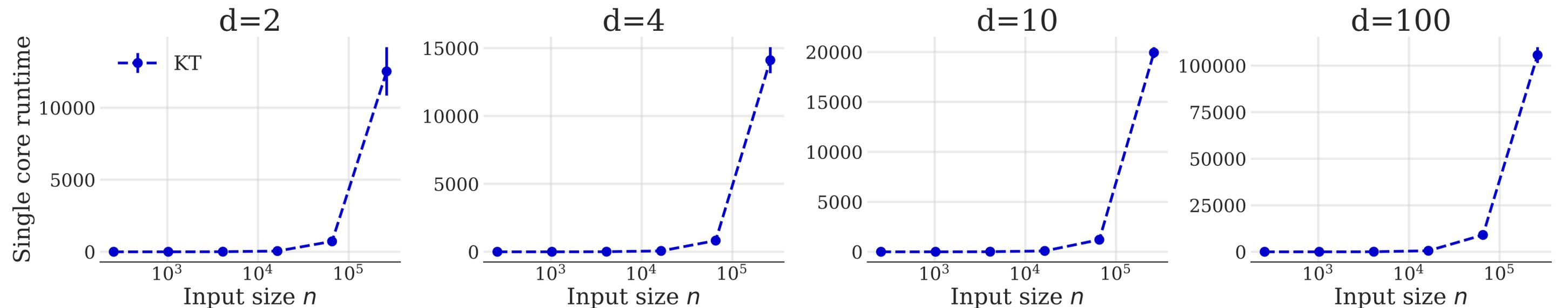
KT drawback: n^2 runtime with n input points

Y-axis = Runtime in log-scale

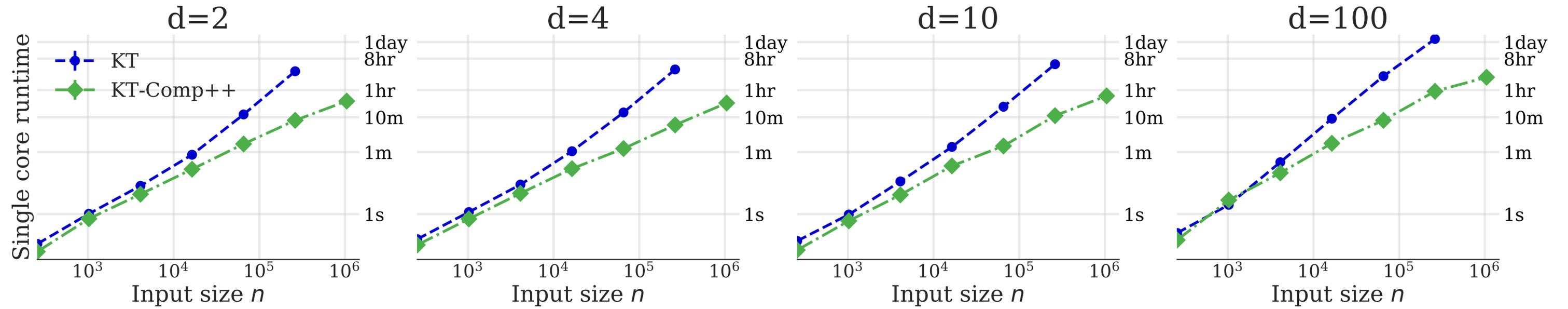
(Runtime dominated by kernel evaluations)



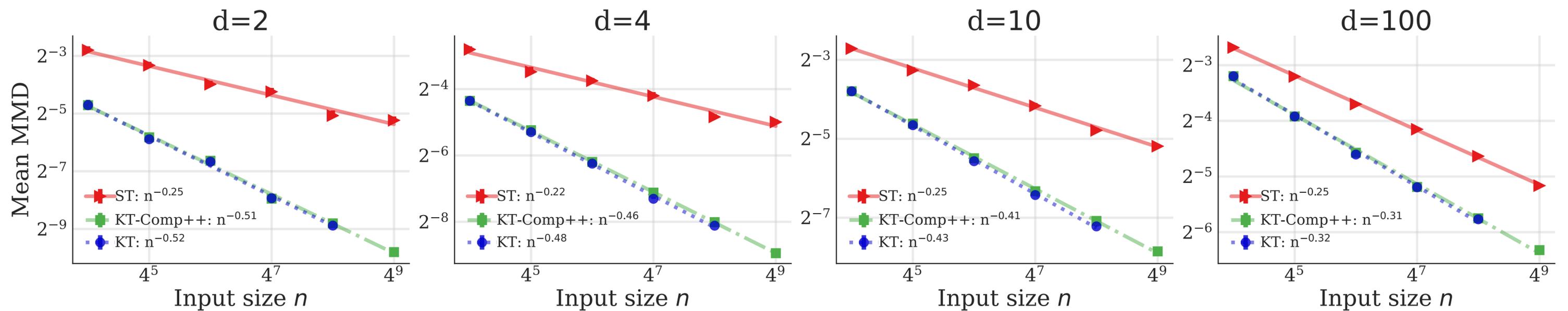
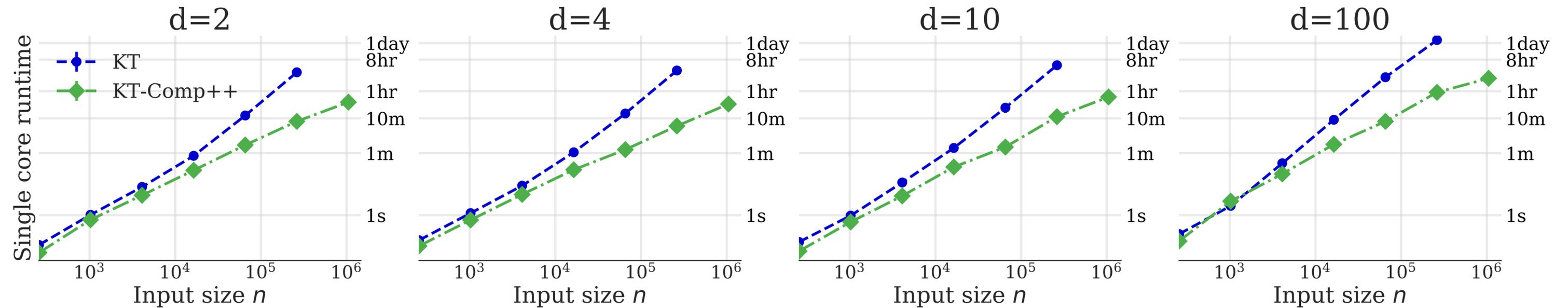
Y-axis = Runtime in linear scale (seconds)



Compress++: Reducing runtime



Compress++: Reducing runtime with minimal loss in accuracy!!



Compress++: A recursive strategy to reduce runtime for generic thinning algorithms

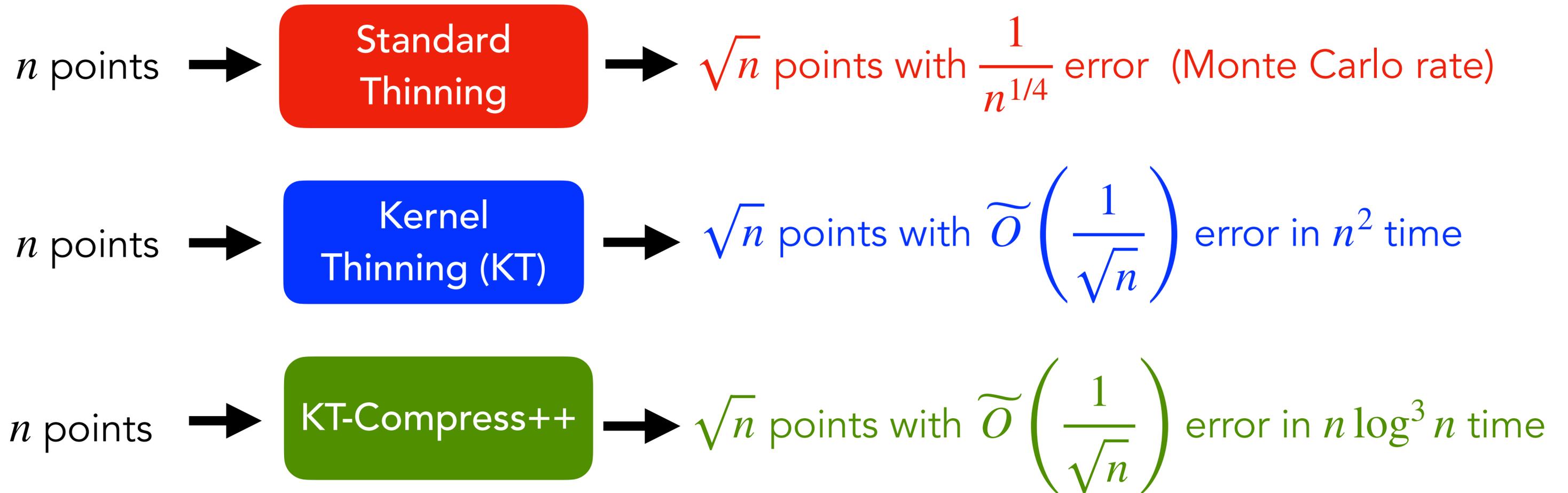


Compress(\mathcal{S} , \mathbf{g} , ALG):

- If $\text{size}(\mathcal{S}) == 1$, **Return** \mathcal{S}
- Else:
 - (i) Call Compress separately on 4 equal splits of \mathcal{S}
 - (ii) Concatenate the 4 outputs from step (i)
 - (iii) **Return** Halved output of step (ii) using ALG

Summary:

KT-Compress++ provides near-optimal compression in near-linear time



 python™ **pip install goodpoints**

[arXiv.org](https://arxiv.org)

<https://arxiv.org/abs/2105.05842>

[Kernel Thinning, COLT 2021]

<https://arxiv.org/abs/2110.01593>

[Generalized Kernel Thinning, ICLR 2022]

<https://arxiv.org/pdf/2111.07941.pdf>

[Distribution Compression In Near-linear Time, ICLR 2022]

Additional slides

Details of kernel thinning

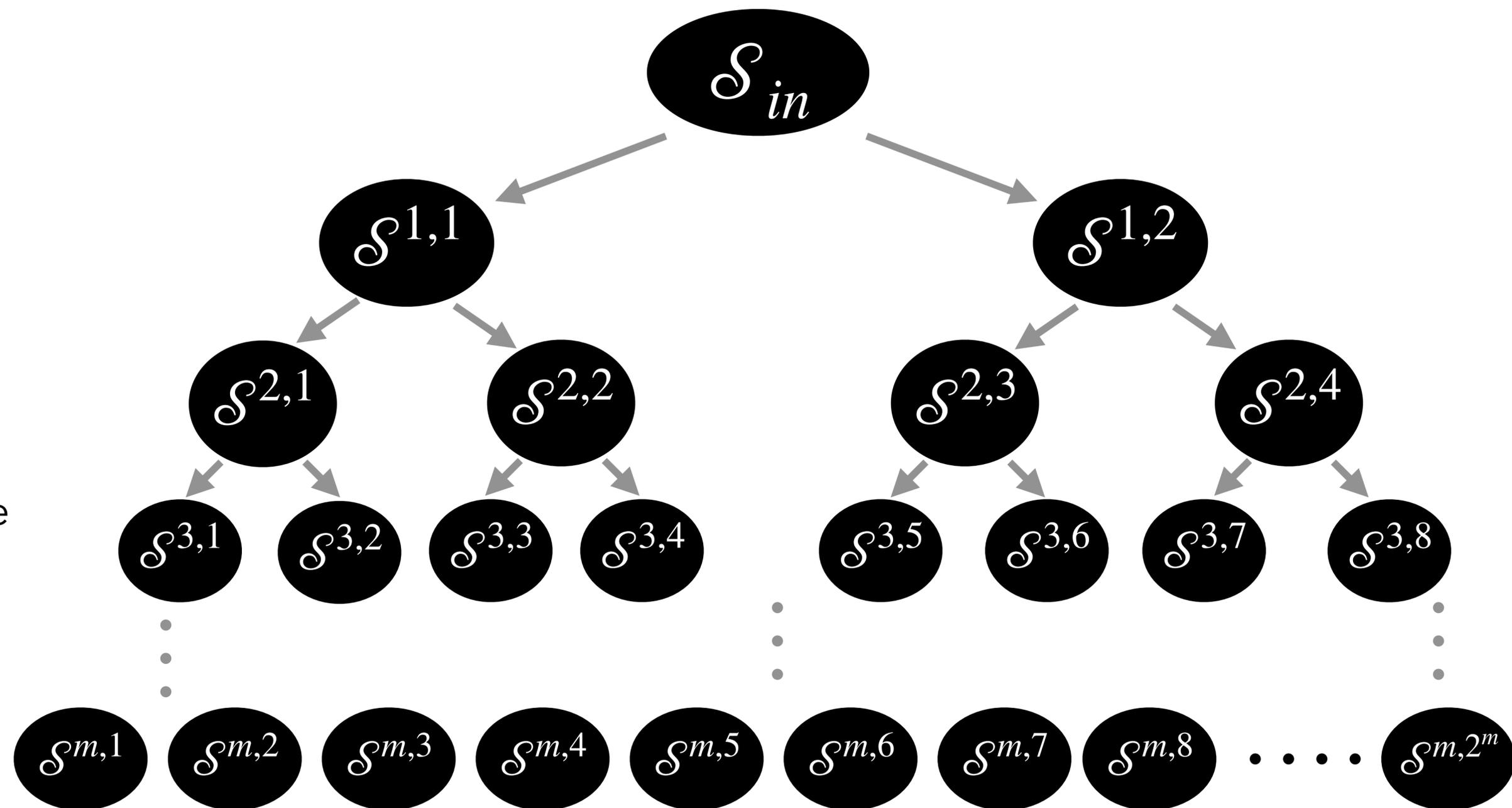
KT: A two-staged algorithm

- **Input:** Kernel \mathbf{k} , input points \mathcal{S}_{in} of size n , thinning factor m
- **KT-Split:**
 - Split \mathcal{S}_{in} into 2^m balanced candidate coresets each of size $\frac{n}{2^m}$
 - When $m = \frac{1}{2} \log_2 n$, we have \sqrt{n} coresets each of size \sqrt{n}
- **KT-Swap:**
 - Pick the best candidate coreset that minimized $\text{MMD}_{\mathbf{k}}$ to input
 - Iteratively refine each point in the selected coreset by swapping with the best alternative \mathcal{S}_{in} if it improves the MMD error

Computation: $\mathcal{O}(n^2)$ kernel evaluations
Storage: $n \min(n, d)$

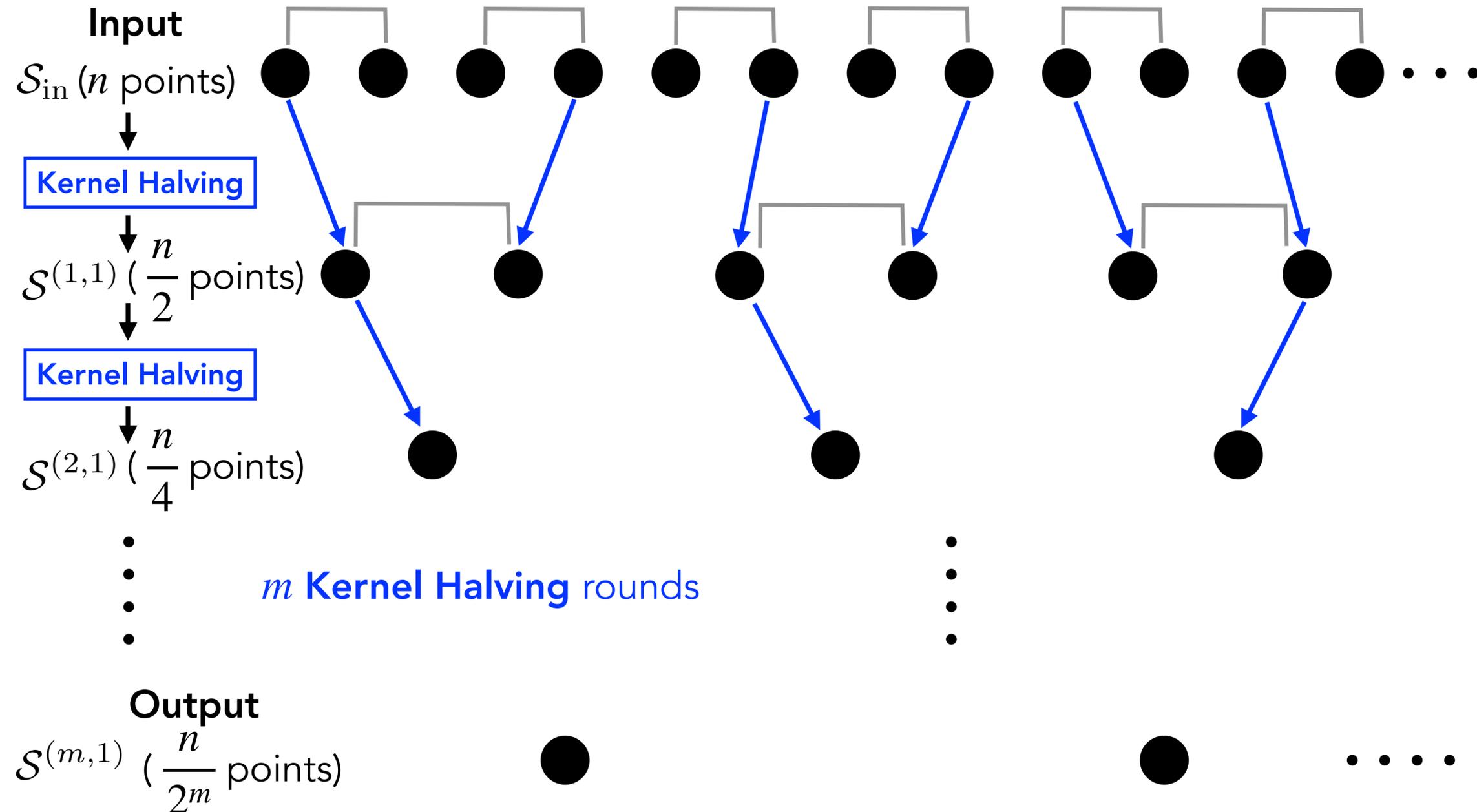
KT-Split

- Repeated rounds of splitting the parent coreset in two equal-sized children coresets
- Runs online, after seeing t input points, the bottom nodes have $t/2^m$ points



KT-Split

- One path on the tree is obtained by repeated **kernel halving**
- At each halving round, remaining points are paired, and one point is selected **non-uniformly** from each pair using a **new Hilbert space generalization of the self-balancing walk** of [Alweiss-Liu-Sawhney 2020]



Algorithm 2: Self-balancing Hilbert Walk

Input: sequence of functions $(f_i)_{i=1}^n$ in Hilbert space \mathcal{H} , threshold sequence $(\mathbf{a}_i)_{i=1}^n$

$\psi_0 \leftarrow \mathbf{0} \in \mathcal{H}$

for $i = 1, 2, \dots, n$ **do**

$\alpha_i \leftarrow \langle \psi_{i-1}, f_i \rangle_{\mathcal{H}}$ // Compute Hilbert space inner product

if $|\alpha_i| > \mathbf{a}_i$:

$\psi_i \leftarrow \psi_{i-1} - f_i \cdot \alpha_i / \mathbf{a}_i$

else:

$\eta_i \leftarrow 1$ with probability $\frac{1}{2}(1 - \alpha_i / \mathbf{a}_i)$ and $\eta_i \leftarrow -1$ otherwise

$\psi_i \leftarrow \psi_{i-1} + \eta_i f_i$

end

return ψ_n , combination of signed input functions

Algorithm 2: Self-balancing Hilbert Walk

Input: sequence of functions $(f_i)_{i=1}^n$ in Hilbert space \mathcal{H} , threshold sequence $(\mathbf{a}_i)_{i=1}^n$

$\psi_0 \leftarrow \mathbf{0} \in \mathcal{H}$

for $i = 1, 2, \dots, n$ **do**

$\alpha_i \leftarrow \langle \psi_{i-1}, f_i \rangle_{\mathcal{H}}$ // Compute Hilbert space inner product

if $|\alpha_i| > \mathbf{a}_i$:

$\psi_i \leftarrow \psi_{i-1} - f_i \cdot \alpha_i / \mathbf{a}_i$  We choose \mathbf{a}_i such that this step **does not occur** with high probability

else:

$\eta_i \leftarrow 1$ with probability $\frac{1}{2}(1 - \alpha_i / \mathbf{a}_i)$ and $\eta_i \leftarrow -1$ otherwise

$\psi_i \leftarrow \psi_{i-1} + \eta_i f_i$

end

return ψ_n , combination of signed input functions

- Exact **Kernel halving**: When $f_i = \mathbf{k}(x_{2i}, \cdot) - \mathbf{k}(x_{2i-1}, \cdot)$, exactly half of input points (\mathcal{S}_{out}) given -1 sign after $n/2$ steps

$$\frac{1}{n}\psi = \frac{1}{n} \sum_{x \in \mathcal{S}_{in}} \mathbf{k}(x, \cdot) - \frac{2}{n} \sum_{x \in \mathcal{S}_{out}} \mathbf{k}(x, \cdot) = \mathbb{P}_{in} \mathbf{k} - \mathbb{P}_{out} \mathbf{k}$$

Algorithm 2: Self-balancing Hilbert Walk

Input: sequence of functions $(f_i)_{i=1}^n$ in Hilbert space \mathcal{H} , threshold sequence $(\mathbf{a}_i)_{i=1}^n$

$\psi_0 \leftarrow \mathbf{0} \in \mathcal{H}$

for $i = 1, 2, \dots, n$ **do**

$\alpha_i \leftarrow \langle \psi_{i-1}, f_i \rangle_{\mathcal{H}}$ // Compute Hilbert space inner product

if $|\alpha_i| > \mathbf{a}_i$:

$\psi_i \leftarrow \psi_{i-1} - f_i \cdot \alpha_i / \mathbf{a}_i$

else:

$\eta_i \leftarrow 1$ with probability $\frac{1}{2}(1 - \alpha_i / \mathbf{a}_i)$ and $\eta_i \leftarrow -1$ otherwise

$\psi_i \leftarrow \psi_{i-1} + \eta_i f_i$

end

return ψ_n , combination of signed input functions

- **Balance:** If \mathbf{k} is a reproducing kernel, for all $g \in \mathbb{H}_{\mathbf{k}}$,

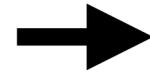
$\langle \psi_n, g \rangle_{\mathbf{k}} = \mathbb{P}_{in}g - \mathbb{P}_{out}g$ is $\mathcal{O}(n^{-1} \cdot \sqrt{\log n} \cdot \|g\|_{\mathbf{k}})$ -sub-Gaussian

If η_i were chosen i.i.d., the sub-Gaussian parameter is $\Omega(n^{-1/2})$

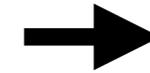
Details of Compress++

Compress++: A simple two-stage algorithm

n points, parameter g ,
halving algorithm HALVE



Compress



$2^g \sqrt{n}$ points



2^g thinning
algorithm THIN



\sqrt{n} points

For example:

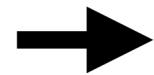
HALVE = Kernel thinning by a factor of 2

THIN = Kernel thinning by a factor of 2^g

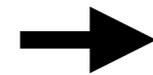
(other algorithms can be used too!)

Compress++: A simple two-stage algorithm

n points, parameter g ,
halving algorithm HALVE



Compress



$2^g \sqrt{n}$ points



2^g thinning
algorithm THIN



\sqrt{n} points

Compress(\mathcal{S} , g , ALG):

- If $\text{size}(\mathcal{S}) == 4^g$, **Return** \mathcal{S}
- Else:
 - Call Compress separately on 4 equal splits of \mathcal{S}
 - Concatenate the 4 outputs from step (i)
 - Return** Halved output of step (ii) using ALG

Compress++: Informal guarantee

Under some mild conditions, with $g = \log \log n + 1$, we have

- **Sub-Gaussian error inflation by at most 4:**

If $\text{MMD}_{\mathbf{k}}(\mathbb{P}_n, \mathbb{P}_{HALVE}) \sim e_1(n)$ and $\text{MMD}_{\mathbf{k}}(\mathbb{P}_n, \mathbb{P}_{THIN}) \sim e_2(n)$ then

$$\text{MMD}_{\mathbf{k}}(\mathbb{P}_n, \mathbb{P}_{Compress++}) \sim 4 \max(e_1(n), e_2(n))$$

$$\text{MMD}_{\mathbf{k}}(\mathbb{P}_{in}, \mathbb{P}_{out}) = \sup_{\|g\|_{\mathbf{k}} \leq 1} |\mathbb{P}_{in}g - \mathbb{P}_{out}g|$$

Compress++: Informal guarantee

Under some mild conditions, with $g = \log \log n + 1$, we have

- **Sub-Gaussian error inflation by at most 4:**

If $\text{MMD}_{\mathbf{k}}(\mathbb{P}_n, \mathbb{P}_{HALVE}) \sim e_1(n)$ and $\text{MMD}_{\mathbf{k}}(\mathbb{P}_n, \mathbb{P}_{THIN}) \sim e_2(n)$ then

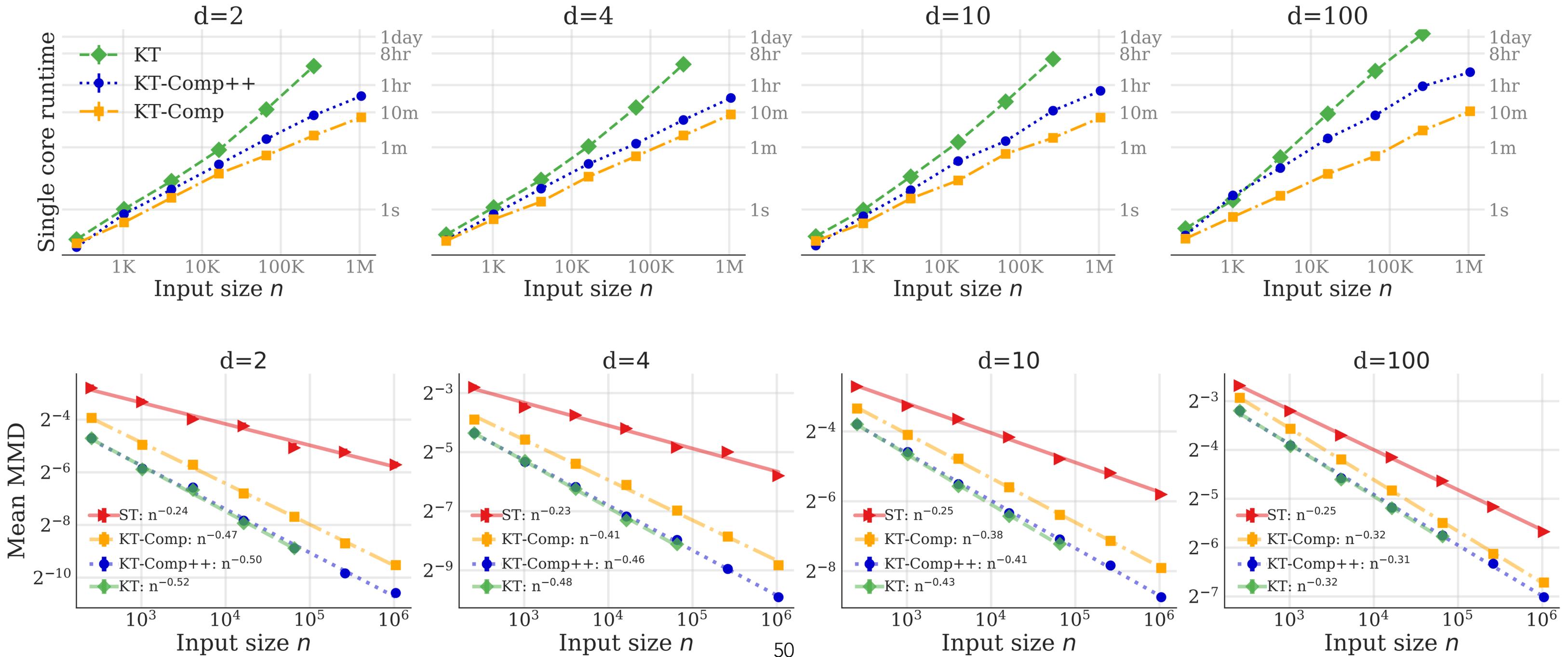
$$\text{MMD}_{\mathbf{k}}(\mathbb{P}_n, \mathbb{P}_{Compress++}) \sim 4 \max(e_1(n), e_2(n))$$

- **Quadratic reduction in runtime:**

If runtime of HALVE and THIN with n points is $\mathcal{O}(n^\tau)$ then the runtime of Compress++ with n points is $\mathcal{O}(n^{\tau/2})$ if $\tau > 2$ and $\mathcal{O}(n \log^3 n)$ if $\tau = 2$.

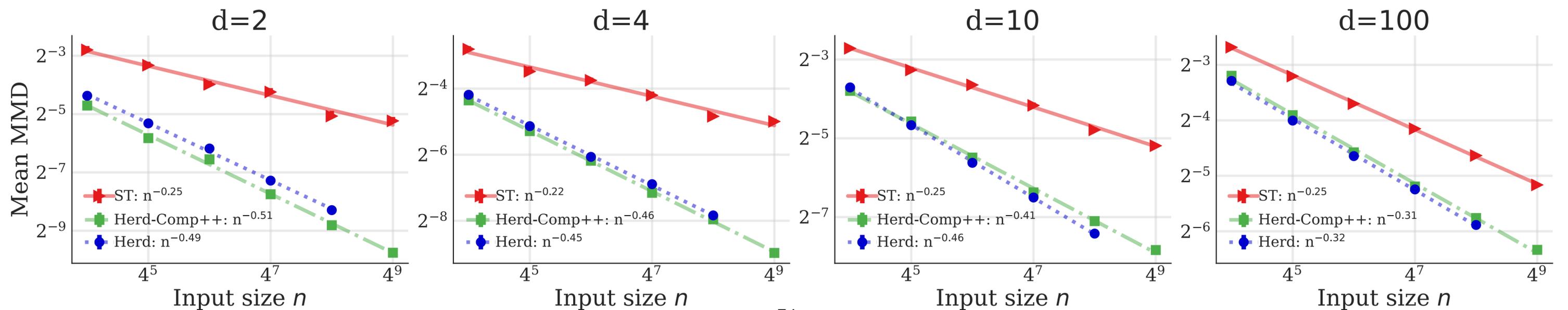
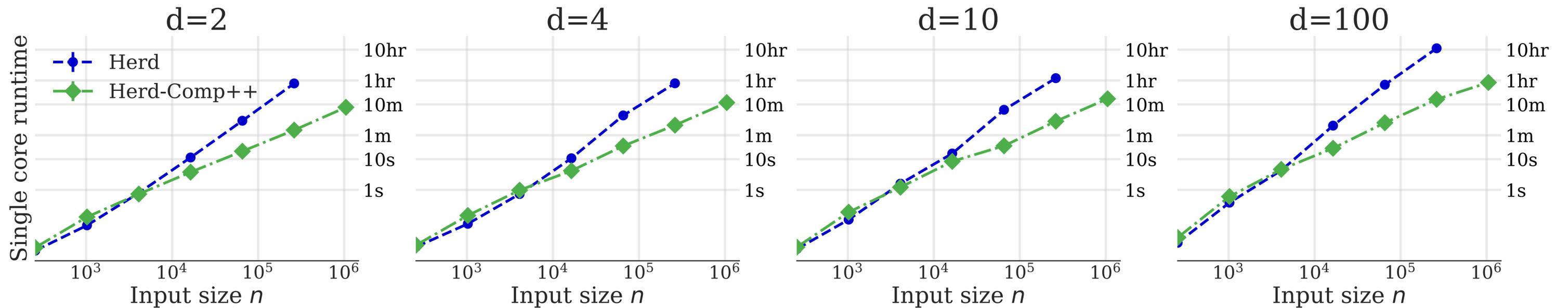
KT vs Compress ($g = 0$) vs Compress++ ($g = 4$)

The input algorithms Halve and Thin to Compress++ are derived from KT



Results for Compress++ with kernel herding (Herd)

Results for Compress++ with kernel herding (Herd)



Lower bounds

Lower bounds

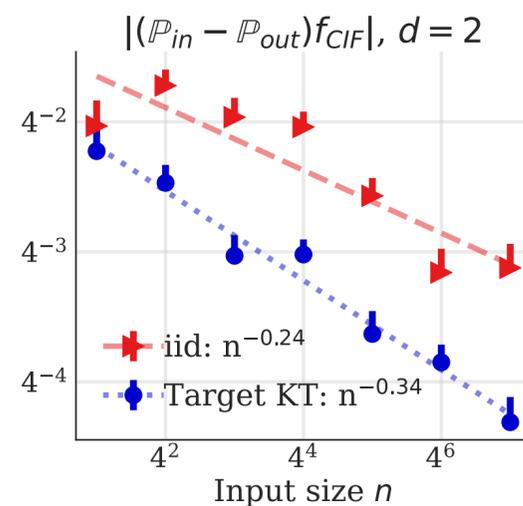
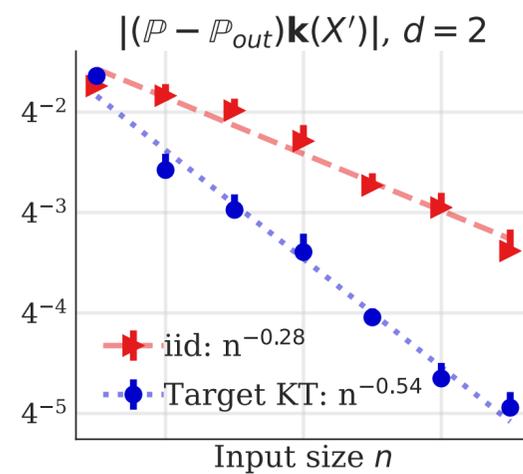
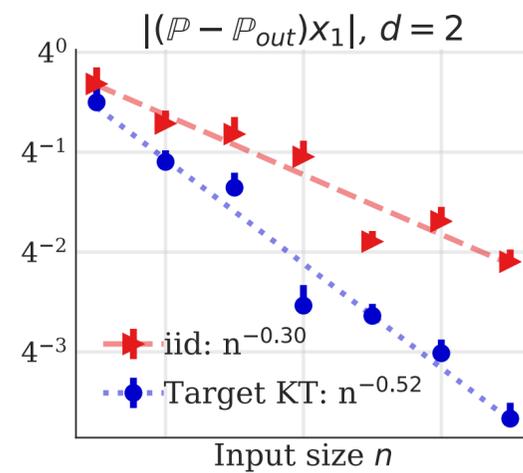
- For smooth kernels, there exists a target \mathbb{P} , such that a coresset of size \sqrt{n} suffers an MMD error of $\min(\sqrt{\frac{d}{n}}, n^{-1/4})$. [Philips and Tai 2020]
- For characteristic kernels, there exists a target \mathbb{P} , such that any estimator based on n i.i.d. input points must suffer at least $n^{-1/2}$ MMD error. [Tolstikhin et al. 2017]

Both bounds apply to Gaussian and Matérn kernels

Single function + Additional MCMC Experiments

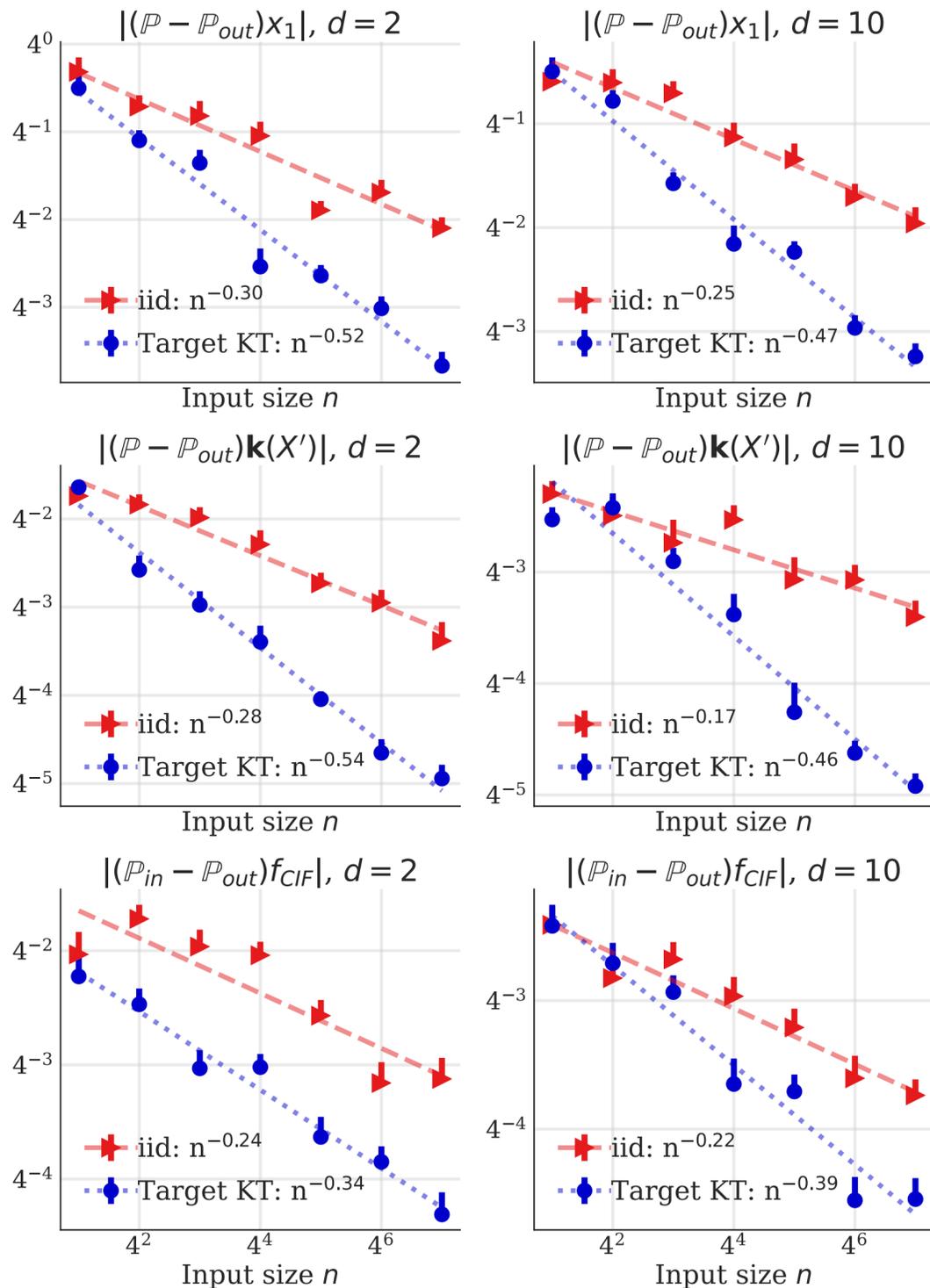
Better error for functions **inside and outside of RKHS**

(Gaussian \mathbf{k} with $\sigma^2 = 2d$ and standard Gaussian \mathbb{P}^\star)



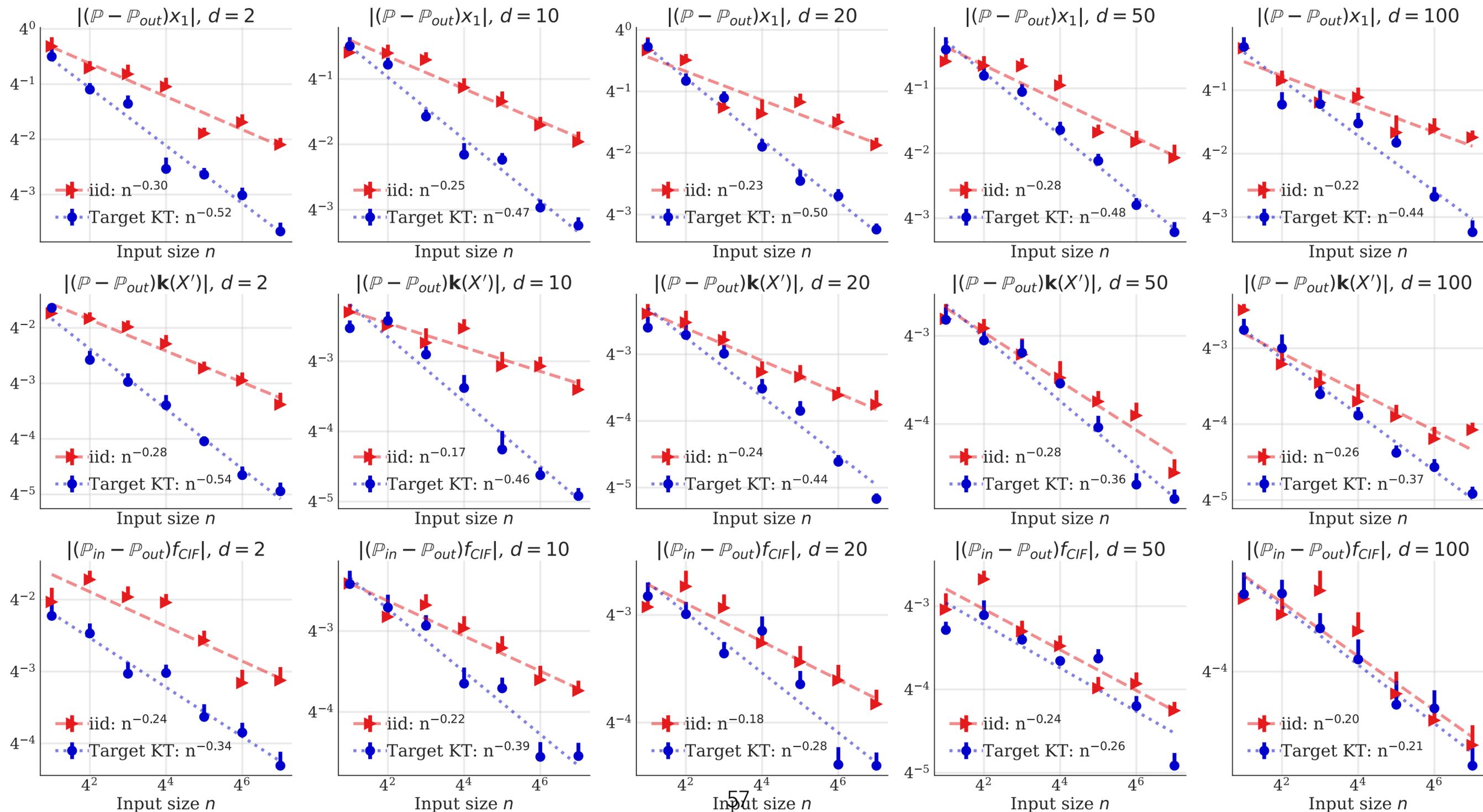
Better error for functions **inside and outside of RKHS**

(Gaussian \mathbf{k} with $\sigma^2 = 2d$ and standard Gaussian \mathbb{P}^*)



Better error for functions **inside and outside of RKHS**

(Gaussian \mathbf{k} with $\sigma^2 = 2d$ and standard Gaussian \mathbb{P}^*)



MCMC experiments: Differential equation models

Dimension $d = 4$

- 1. Lotka-Volterra model
oscillatory enzymatic control, [1925, 1926]
- 2. Goodwin model
oscillatory predator-prey evolution, [1965]

X

- 1. Posterior

X

- 1. Random walk (RW)
[Metropolis et al. 1953, Hastings 1970]
- 2. Adaptive random walk (adaRW)
[Haario et al. 1999]
- 3. Metropolis adjusted Langevin algorithm (MALA) [Roberts et al. 1996]
- 4. Preconditioned-MALA (pMALA)
[Girolami et al. 2011]

Dimension $d = 38$

- 3. Hinch calcium signal model
[Hinch-Greenstein-Tanskanen-Xu-Winslow, 2004]

X

- 1. Posterior
- 2. Tempered posterior

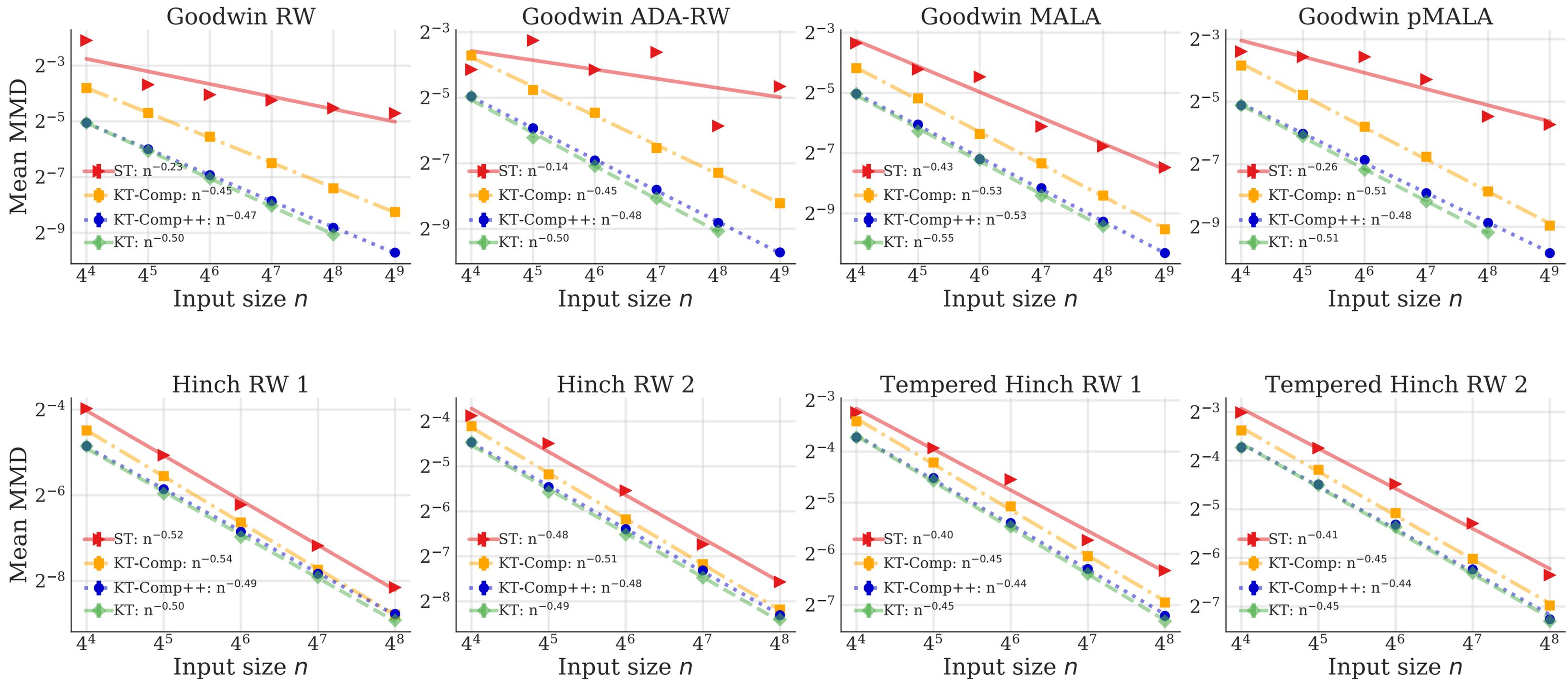
X

- 1. Random walk (RW) - run 1
- 2. Random walk (RW) - run 2

For KT, we use Gaussian kernel, and chose its bandwidth via median heuristic [Garreau et al. 2017]

MCMC samples taken from Riabiz-Chen-Cockayne-Swietach-Niederer-Mackey-Oates, 2021

Results for MCMC experiments



Details for KT result

Target KT MMD rates: \sqrt{n} points with $\widetilde{O}(n^{-1/2})$ error

- More generally

$\widetilde{O}\left(\frac{1}{\sqrt{n}}\right)$ error rate for analytic kernels

$\widetilde{O}\left(\frac{n^{d/2m}}{\sqrt{n}}\right)$ for m -times differentiable kernels

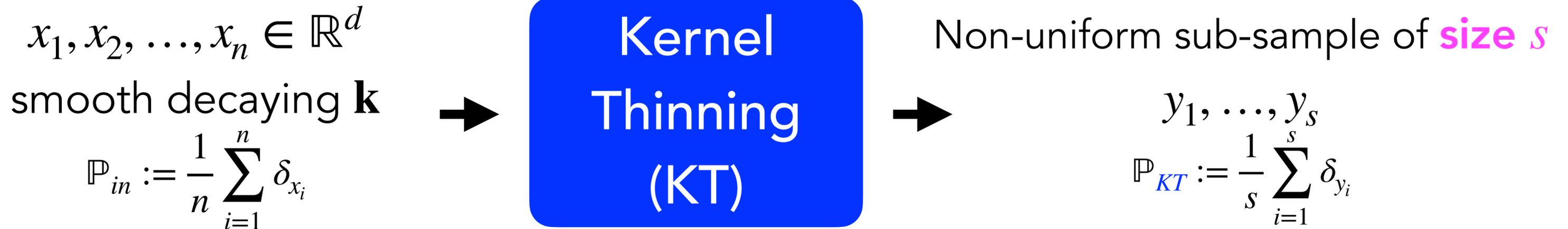
- We state explicit constants with dependence on kernel hyper-parameters in the paper

Generalized kernel thinning

	Root KT	Target KT	KT+ [Best of both worlds]
KT-Split kernel	\mathbf{k}_{rt}	\mathbf{k}	$\mathbf{k} + \mathbf{k}_{\alpha-rt}$
Single-function error	Same as MMD error	$\sqrt{\frac{\log n}{n}}$ For arbitrary \mathbf{k} on arbitrary domain	$\sqrt{\frac{\log n}{n}}$
MMD error	See <u>slide</u>	$\sqrt{\frac{\log^{ad+b} n}{n}}$ Analytic \mathbf{k} $\sqrt{\frac{n^{d/m}}{n}}$ m -times differentiable \mathbf{k}	Min(Target KT Error, α -Root KT)

$$\mathbf{k}(x, y) = \int \mathbf{k}_{rt}(x, z)\mathbf{k}_{rt}(z, y)dz \ \& \ \mathbf{k}_{\alpha-rt} = \widehat{(\hat{\mathbf{k}})^\alpha} \text{ where } \hat{\ } \text{ denotes Fourier transform}$$

Target KT or KT+: Better than Monte Carlo rate

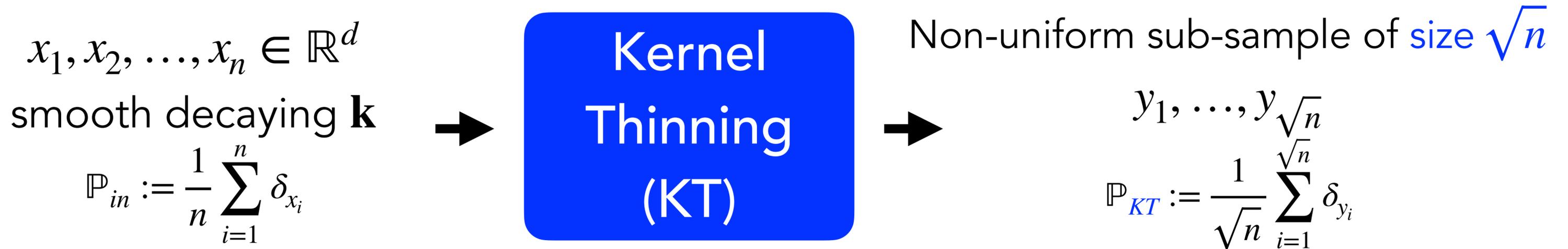


For any fixed $g \in \mathbb{H}_{\mathbf{k}}$, with probability $1 - \delta$ over the randomness in KT, we have

$$|\mathbb{P}_{in}g - \mathbb{P}_{KT}g| \leq \frac{1}{s} \cdot \|g\|_{\mathbf{k}} \sqrt{\frac{8}{3} \|\mathbf{k}\|_{\infty} \log\left(\frac{4}{\delta}\right) \log\left(\frac{6s \log(n/s)}{\delta}\right)}$$

Much faster than the Monte Carlo rate for standard/uniform thinning $\mathcal{O}\left(\frac{1}{\sqrt{s}}\right)$

Target KT or KT+: Better than Monte Carlo rate



For any fixed $g \in \mathbb{H}_{\mathbf{k}}$, with probability $1 - \delta$ over the randomness in KT, we have

$$|\mathbb{P}_{in}g - \mathbb{P}_{KT}g| \leq \frac{1}{\sqrt{n}} \cdot \|g\|_{\mathbf{k}} \sqrt{\frac{8}{3} \|\mathbf{k}\|_{\infty} \log\left(\frac{4}{\delta}\right) \log\left(\frac{6\sqrt{n} \log \sqrt{n}}{\delta}\right)}$$

Much faster than the Monte Carlo rate for standard/uniform thinning $\mathcal{O}\left(\frac{1}{n^{1/4}}\right)$

Properties of MMD

- Maximum mean discrepancy (MMD) = worst-case integration discrepancy between two distributions over a class of real-valued test functions

$$\text{MMD}_{\mathbf{k}}(\mathbb{P}_{in}, \mathbb{P}_{out}) = \sup_{\|g\|_{\mathbf{k}} \leq 1} |\mathbb{P}_{in}g - \mathbb{P}_{out}g|$$

[Gretton-Borgwardt-Rasch-Schölkopf-Smola, 2012]

- **MMD metrizes convergence in distribution** for popular infinite-dimensional kernels like Gaussian, Matern, IMQ, B-spline

[Simon-Gabriel-Barp-Schölkopf-Mackey, 2020]