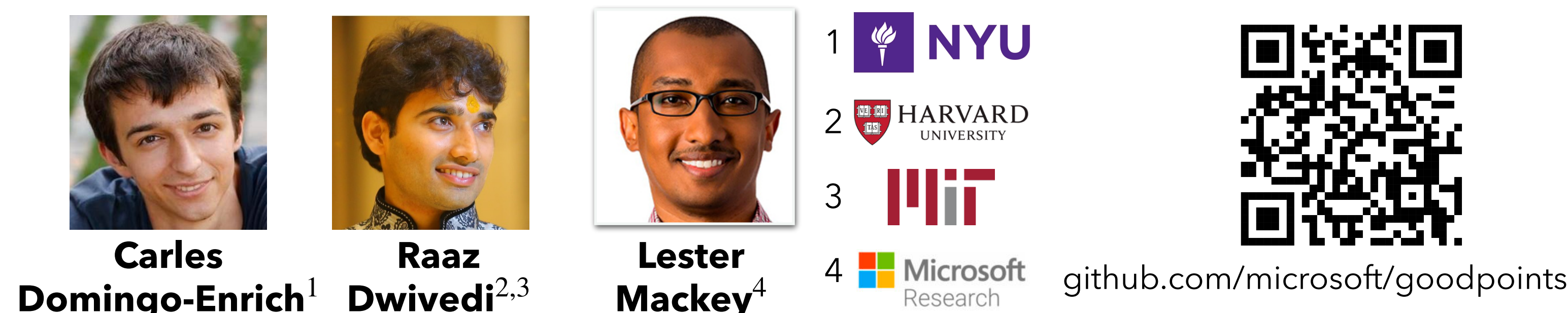


Compress Then Test: Powerful Kernel Testing in Near-linear Time



Kernel two-sample testing

- $\mathbb{X}_n = (X_i)_{1 \leq i \leq n}$ i.i.d. sample of \mathbb{P} , $\mathbb{Y}_n = (Y_i)_{1 \leq i \leq n}$ i.i.d. sample of \mathbb{Q} .
- Null hypothesis: $\mathcal{H}_0 : \mathbb{P} = \mathbb{Q}$
- Non-parametric form via kernel maximum mean discrepancy (MMD)

$$\mathcal{H}_0 : \text{MMD}_{\mathbf{k}}^2(\mathbb{P}, \mathbb{Q}) \stackrel{\Delta}{=} \mathbb{E}_{X, X' \sim \mathbb{P}} \mathbf{k}(X, X') + \mathbb{E}_{Y, Y' \sim \mathbb{P}} \mathbf{k}(Y, Y') - 2\mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \mathbf{k}(X, Y) = 0$$

with Test statistic: $\Delta(\mathbb{X}_n, \mathbb{Y}_n) = \begin{cases} 0 & \text{if } \text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n) < t_\alpha \text{ (accept } \mathcal{H}_0) \\ 1 & \text{if } \text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n) \geq t_\alpha \text{ (reject } \mathcal{H}_0) \end{cases}$

Prior MMD strategies are slow

- Computing MMD takes $O(n^2)$ time

$$\text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n) = \frac{\sum_{i,i'=1}^n \mathbf{k}(X_i, X_{i'})}{n^2} + \frac{\sum_{j,j'=1}^n \mathbf{k}(Y_j, Y_{j'})}{n^2} - \frac{2 \sum_{i=1}^n \sum_{j=1}^n \mathbf{k}(X_i, Y_j)}{n^2}$$

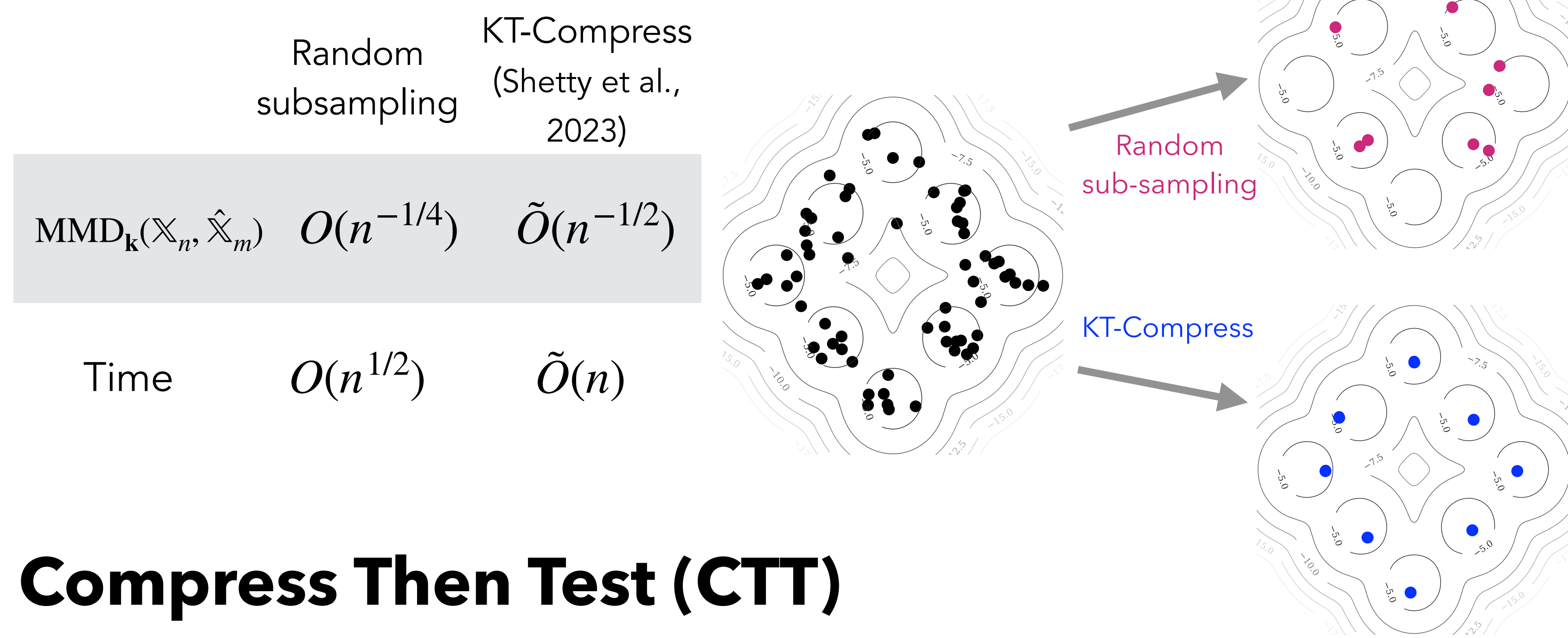
- Can we speed up the test while

- respecting Type I error: reject null rarely when $\mathbb{P} = \mathbb{Q}$, and
- keeping our test powerful: reject null often when $\mathbb{P} \neq \mathbb{Q}$?

- Prior speedup strategies **sacrifice power!**

Speed up testing by compressing

- Compress \mathbb{X}_n to $\hat{\mathbb{X}}_m$ of size $m = \sqrt{n}$ with small $\text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \hat{\mathbb{X}}_m)$



Compress Then Test (CTT)

1. Run KT-Compress to get $\mathbb{X}_n \rightarrow \hat{\mathbb{X}}_m$ & $\mathbb{Y}_n \rightarrow \hat{\mathbb{Y}}_m$ with $m = 2^g \sqrt{n}$
2. Use $\text{MMD}_{\mathbf{k}}^2(\hat{\mathbb{X}}_m, \hat{\mathbb{Y}}_m)$ instead of $\text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n)$
3. Compute threshold t_α via *cheap permutations*: Group data into $s \ll \sqrt{n}$ bins, sample \mathcal{B} permutations of $[s]$, and permute the s bins

CTT runtime: $\tilde{O}(n) + O(s^2 B)$ if $g = \log \log n$

Original test runtime: $O(n^2 B)$

Compress (the Data and) Then Test

1. Turns quadratic time tests to near-linear time tests
2. Provides up to 200x speed-up (1 hour \rightarrow 20 sec)
3. Maintains level and power provably
4. Works with kernel approximations (Low-rank CTT)
5. Applies for kernel selection (Aggregated CTT)

CTT guarantees

A. Exact Type 1 error:

$$\text{If } \text{MMD}_{\mathbf{k}}(\mathbb{P}, \mathbb{Q}) = 0, \text{ then } \Pr[\Delta(\mathbb{X}_n, \mathbb{Y}_n) = 1] = \alpha$$

B. Power:

$$\text{If } \text{MMD}_{\mathbf{k}}(\mathbb{P}, \mathbb{Q}) \geq \text{Separation}(\beta), \text{ then } \Pr[\Delta(\mathbb{X}_n, \mathbb{Y}_n) = 1] \geq 1 - \beta$$

Test name	MMD separation	Runtime	Tails of \mathbb{P} on \mathbb{R}^d	Choice of \mathbf{k}'	$R_{\mathbf{k},\mathbf{k}}(\mathbb{P}, n, \delta, g)$
CTT (ours, Thm. 1)	$\frac{R_{\mathbf{k},\mathbf{k}}(\mathbb{P}, n, \delta, g)}{2^g \sqrt{n}} + n^{-\frac{1}{2}}$	$4^g n \log n$	Compact	Compact \mathbf{k}_{rt}	$(\log \frac{n}{\delta})^2$
Complete MMD (Gretton et al., 2012a)	$n^{-\frac{1}{2}}$	n^2	Subexp	Analytic \mathbf{k}	$(\log \frac{n}{\delta})^{\frac{3d+5}{2}}$
Block MMD (Zaremba et al., 2013)	$(Bn)^{-\frac{1}{4}}$	Bn	Subexp	Subexp \mathbf{k}_{rt}	$c_{m,\delta} (\log \frac{n}{\delta})^{\frac{d+5}{2}}$
Incomplete MMD (Yamada et al., 2019)	$\ell^{-\frac{1}{4}}$	ℓ	ρ -Heavy-tail	ρ -Heavy-tail \mathbf{k}_{rt}	$(\frac{n}{\delta})^{\frac{d}{2\rho}} (\log \frac{n}{\delta})^{\frac{d+5}{2}}$

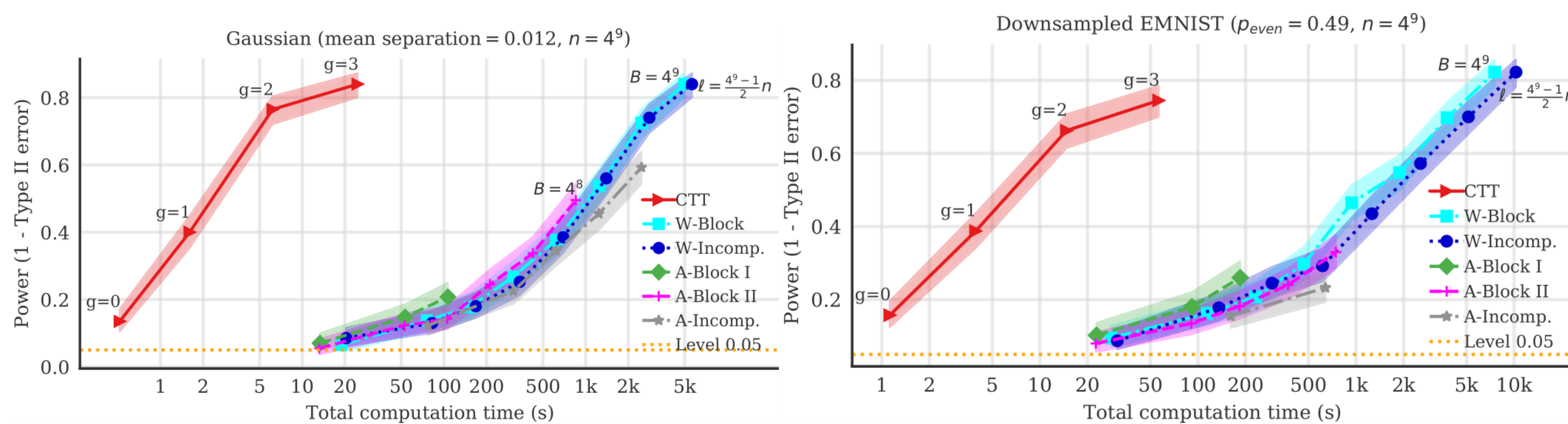
B = Number of blocks used in Block MMD test

ℓ = Number of ordered index pairs in Incomplete MMD test

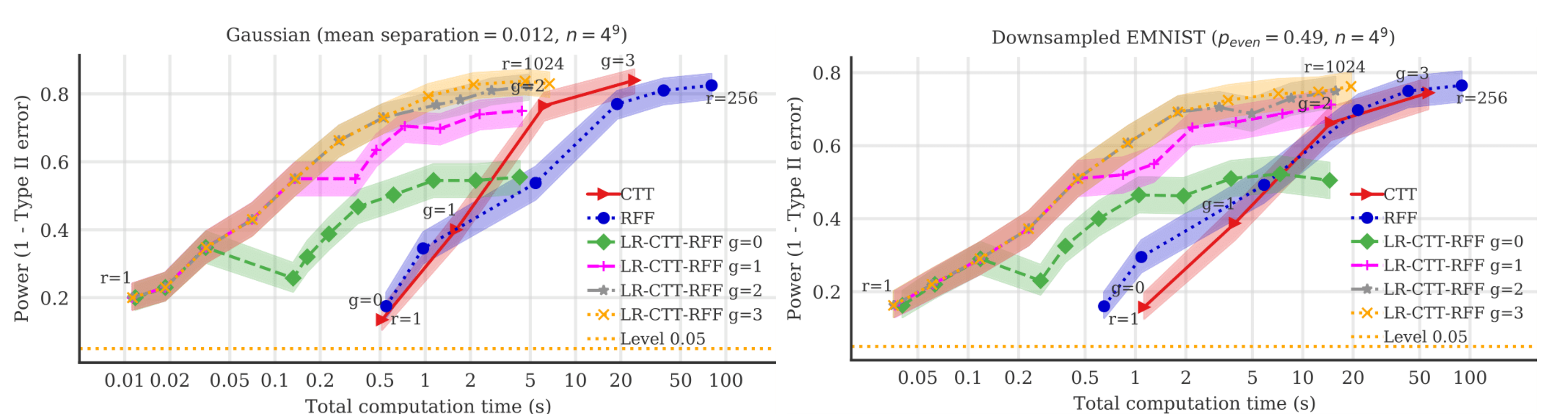
δ = Failure probability for CTT guarantees

\mathbf{k}' = auxiliary kernel used by KT-Compress & $\mathbf{k}(x, y) = \int \mathbf{k}_{\text{rt}}(x, z) \mathbf{k}_{\text{rt}}(z, y) dz$

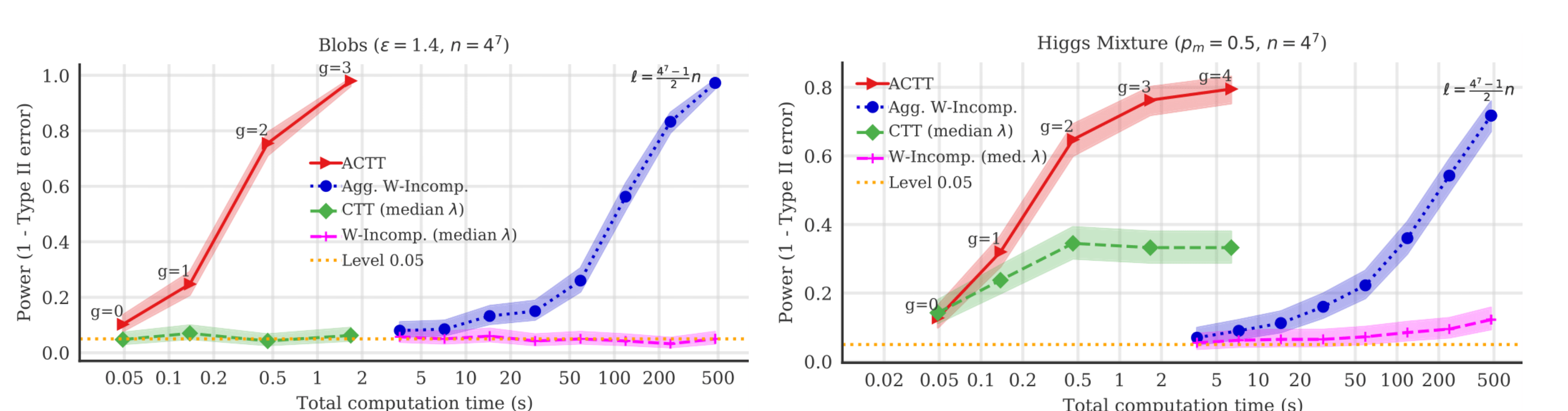
CTT's dominance in time-power tradeoff on Gaussian data and EMNIST data



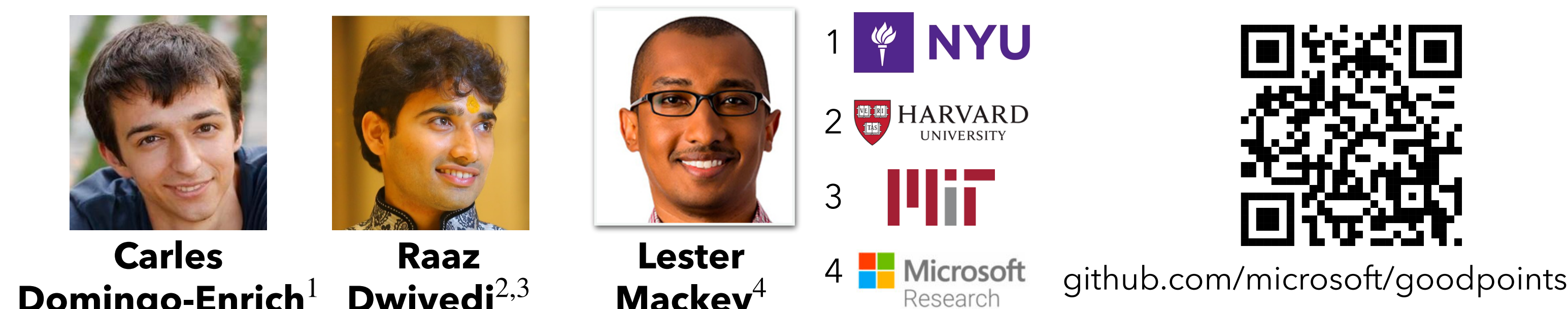
Low-rank CTT's dominance over random Fourier features based test



Aggregated CTT's dominance for kernel selection on Blobs and Higgs data



Compress Then Test: Powerful Kernel Testing in Near-linear Time



Kernel two-sample testing

- $\mathbb{X}_n = (X_i)_{1 \leq i \leq n}$ i.i.d. sample of \mathbb{P} , $\mathbb{Y}_n = (Y_i)_{1 \leq i \leq n}$ i.i.d. sample of \mathbb{Q} .
- Null hypothesis: $\mathcal{H}_0 : \mathbb{P} = \mathbb{Q}$
- Non-parametric form via kernel maximum mean discrepancy (MMD)

$$\mathcal{H}_0 : \text{MMD}_{\mathbf{k}}^2(\mathbb{P}, \mathbb{Q}) \stackrel{\Delta}{=} \mathbb{E}_{X, X' \sim \mathbb{P}} \mathbf{k}(X, X') + \mathbb{E}_{Y, Y' \sim \mathbb{P}} \mathbf{k}(Y, Y') - \mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \mathbf{k}(X, Y) = 0$$

with Test statistic: $\Delta(\mathbb{X}_n, \mathbb{Y}_n) = \begin{cases} 0 & \text{if } \text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n) < t_\alpha \text{ (accept } \mathcal{H}_0) \\ 1 & \text{if } \text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n) \geq t_\alpha \text{ (reject } \mathcal{H}_0) \end{cases}$

Prior MMD strategies are slow

- Computing MMD takes $O(n^2)$ time

$$\text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n) = \frac{\sum_{i, i'=1}^n \mathbf{k}(X_i, X_{i'})}{n^2} + \frac{\sum_{j, j'=1}^n \mathbf{k}(Y_j, Y_{j'})}{n^2} - \frac{2 \sum_{i=1}^n \sum_{j=1}^n \mathbf{k}(X_i, Y_j)}{n^2}$$

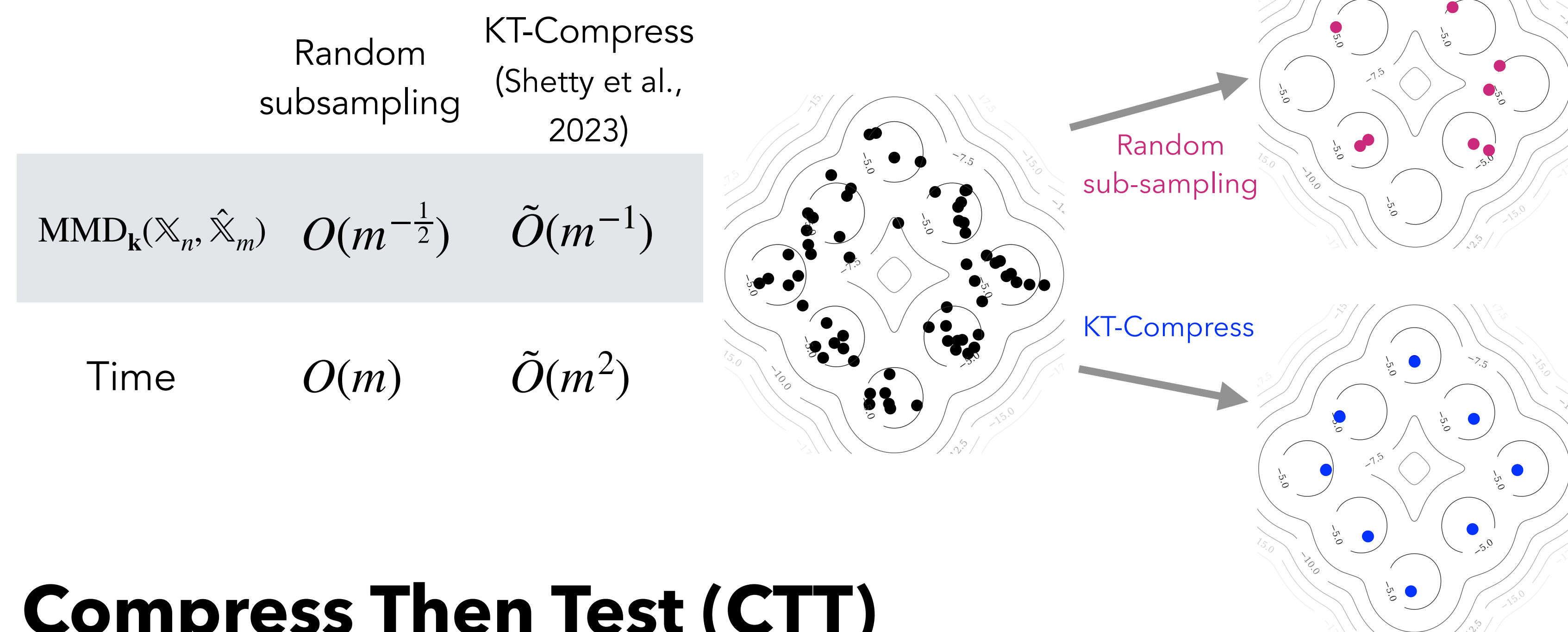
- Can we speed up the test while

- respecting Type I error: reject null rarely when $\mathbb{P} = \mathbb{Q}$, and
- keeping our test powerful: reject null often when $\mathbb{P} \neq \mathbb{Q}$?

- Prior speedup strategies **sacrifice power!**

Speed up testing by compressing

- Compress \mathbb{X}_n to $\hat{\mathbb{X}}_m$ of size $m \ll n$ with small $\text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \hat{\mathbb{X}}_m)$



Compress Then Test (CTT)

1. Run KT-Compress to get $\mathbb{X}_n \rightarrow \hat{\mathbb{X}}_m$ & $\mathbb{Y}_n \rightarrow \hat{\mathbb{Y}}_m$ with $m = 2^g \sqrt{n}$
2. Use $\text{MMD}_{\mathbf{k}}^2(\hat{\mathbb{X}}_m, \hat{\mathbb{Y}}_m)$ instead of $\text{MMD}_{\mathbf{k}}^2(\mathbb{X}_n, \mathbb{Y}_n)$
3. Compute threshold t_α via *cheap permutations*: Group data into $s \ll m$ bins, sample \mathcal{B} permutations of $[s]$, and permute the s bins

Total runtime: $\tilde{O}(4^g n) = \tilde{O}(n)$ if $g = \log \log n$

Compress (the Data and) Then Test

1. Turns quadratic time tests to near-linear time tests
2. Provides up to 200x speed-up (1 hour \rightarrow 20 sec)
3. Maintains level and power provably
4. Works with kernel approximations (Low-rank CTT)
5. Applies for kernel selection (Aggregated CTT)

CTT guarantees

A. Exact Type 1 error:

$$\text{If } \text{MMD}_{\mathbf{k}}(\mathbb{P}, \mathbb{Q}) = 0, \text{ then } \Pr[\Delta(\mathbb{X}_n, \mathbb{Y}_n) = 1] = \alpha$$

B. Power:

$$\text{If } \text{MMD}_{\mathbf{k}}(\mathbb{P}, \mathbb{Q}) \geq \text{Separation}(\beta), \text{ then } \Pr[\Delta(\mathbb{X}_n, \mathbb{Y}_n) = 1] \geq 1 - \beta$$

Test name	MMD separation	Runtime	Tails of \mathbb{P} on \mathbb{R}^d	Choice of \mathbf{k}'	$R_{\mathbf{k}, \mathbf{k}'}(\mathbb{P}, n, \delta, g)$
CTT (ours, Thm. 1)	$\frac{R_{\mathbf{k}, \mathbf{k}'}(\mathbb{P}, n, \delta, g)}{2^g \sqrt{n}} + n^{-\frac{1}{2}}$	$4^g n \log n$	Compact	Compact \mathbf{k}_{rt}	$(\log \frac{n}{\delta})^2$
Complete MMD (Gretton et al., 2012a)	$n^{-\frac{1}{2}}$	n^2	Subexp	Analytic \mathbf{k}	$(\log \frac{n}{\delta})^{\frac{3d+5}{2}}$
Block MMD (Zaremba et al., 2013)	$(Bn)^{-\frac{1}{4}}$	Bn	Subexp	Subexp \mathbf{k}_{rt}	$c_{m, \delta} (\log \frac{n}{\delta})^{\frac{d+5}{2}}$
Incomplete MMD (Yamada et al., 2019)	$\ell^{-\frac{1}{4}}$	ℓ	ρ -Heavy-tail	ρ -Heavy-tail \mathbf{k}_{rt}	$(\frac{n}{\delta})^{\frac{d}{2\rho}} (\log \frac{n}{\delta})^{\frac{d+5}{2}}$

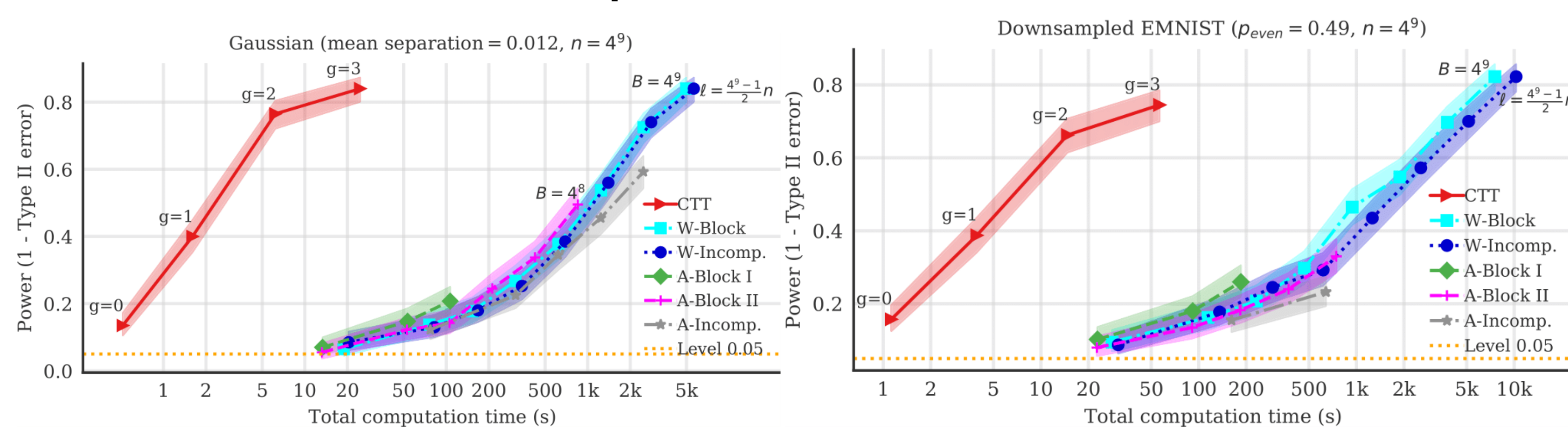
B = Number of blocks used in Block MMD test

ℓ = Number of ordered index pairs in Incomplete MMD test

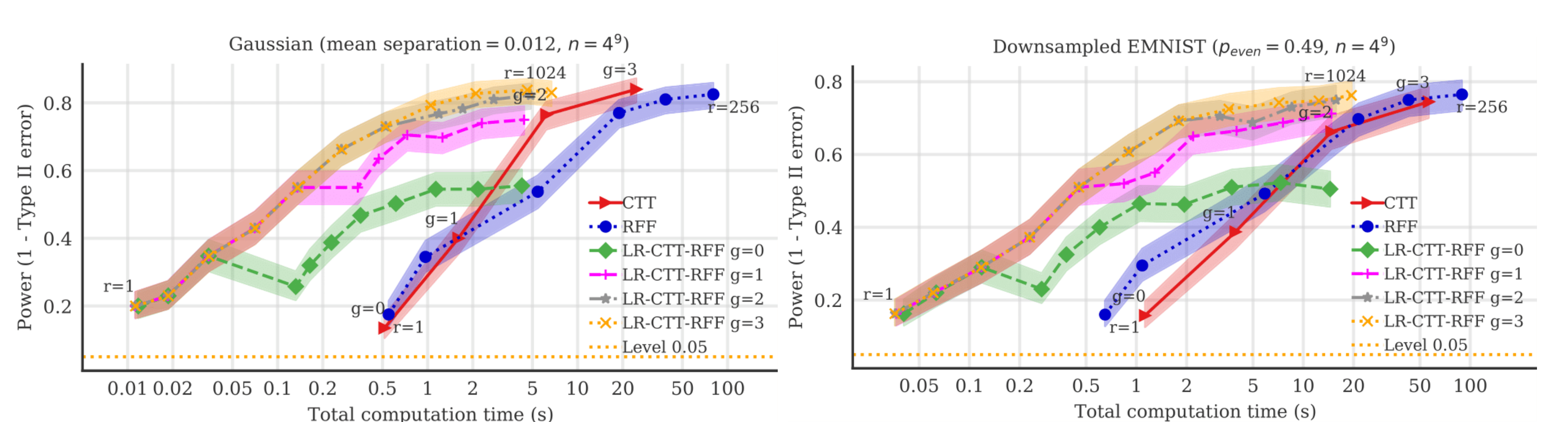
δ = Failure probability for CTT guarantees

\mathbf{k}' = auxiliary kernel used by KT-Compress & $\mathbf{k}(x, y) = \int \mathbf{k}_{\text{rt}}(x, z) \mathbf{k}_{\text{rt}}(z, y) dz$

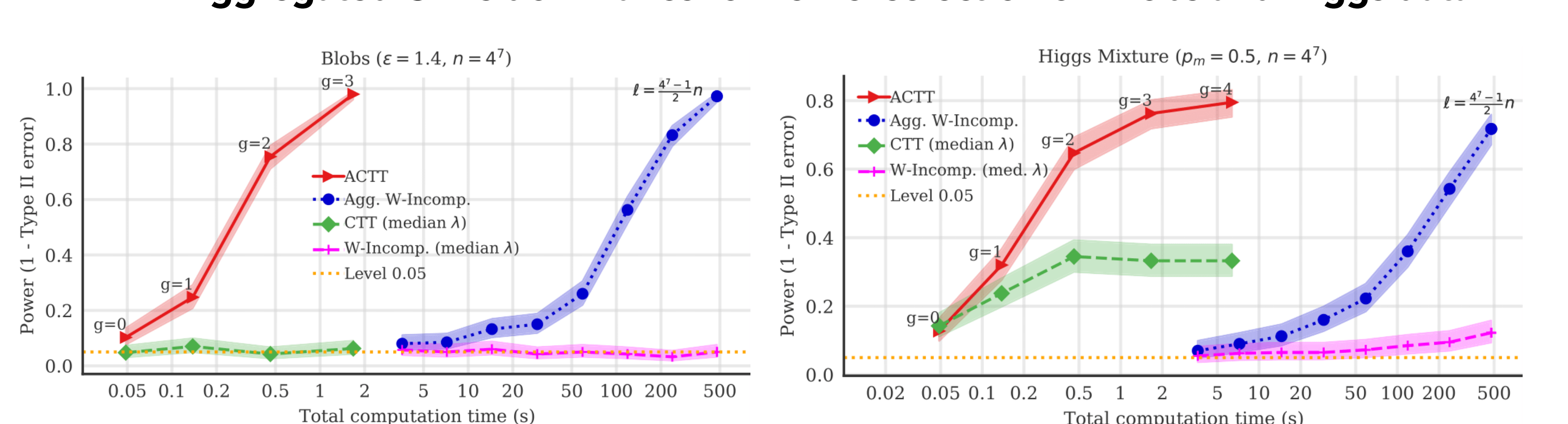
CTT's dominance in time-power tradeoff on Gaussian data and EMNIST data



Low-rank CTT's dominance over random Fourier features based test



Aggregated CTT's dominance for kernel selection on Blobs and Higgs data



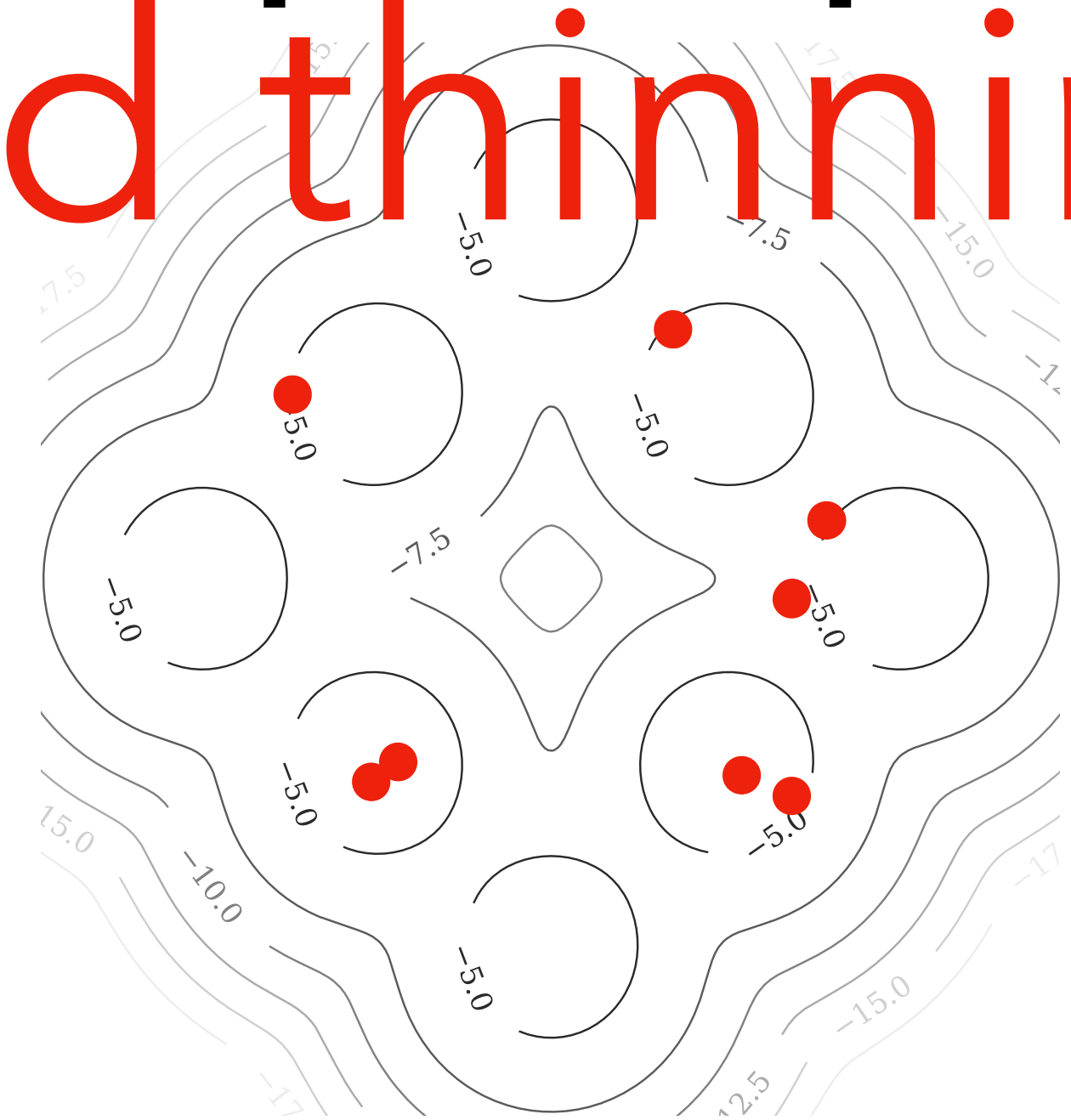
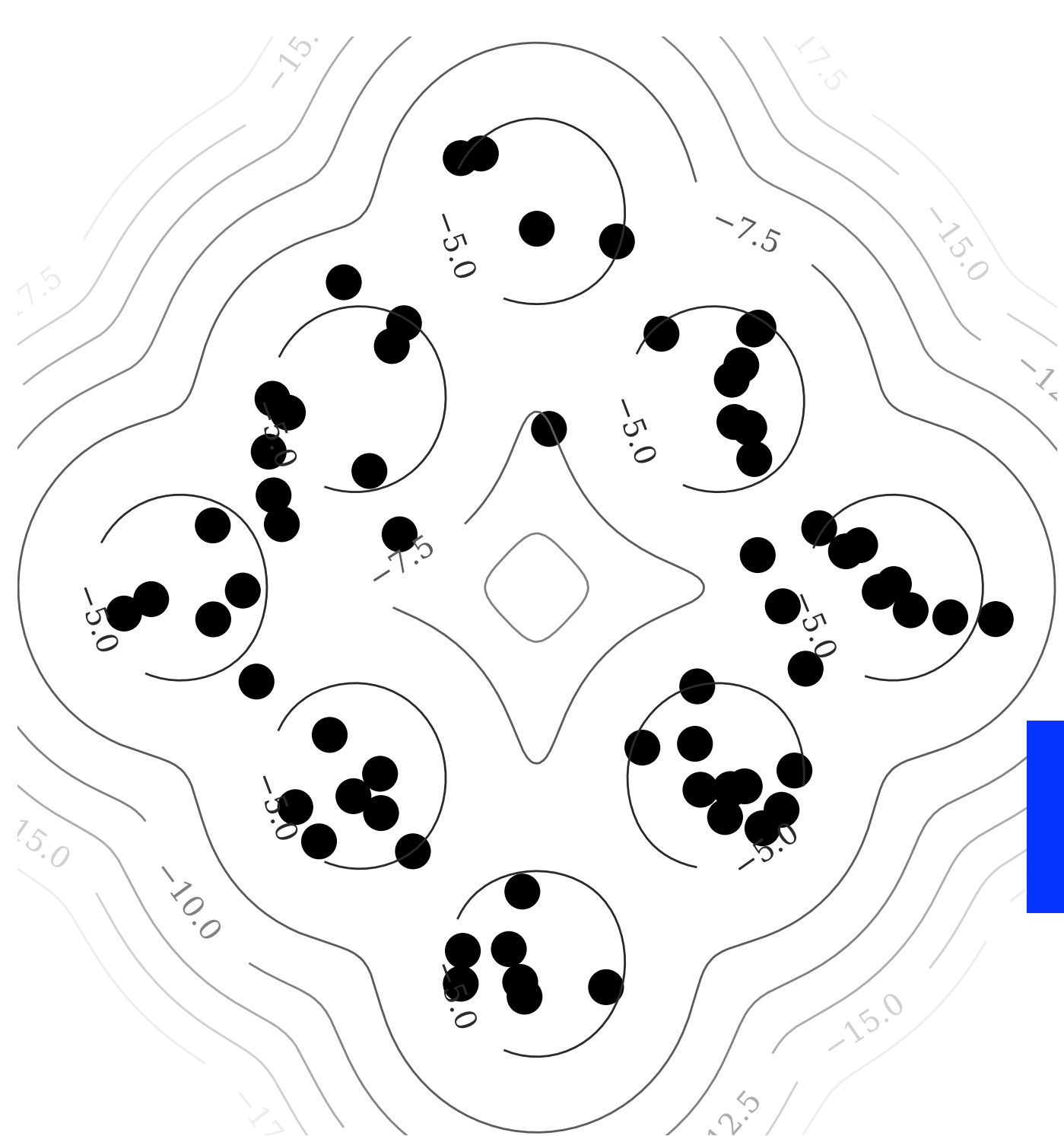
1. Compress $\mathbb{X}_m, \mathbb{Y}_m$ to $\hat{\mathbb{X}}_m, \hat{\mathbb{Y}}_m$ each of size $2^g \sqrt{m}$ using KT-Compress: $\tilde{O}(4^g m)$ time
 2. Use $\text{MMD}_k^2(\hat{\mathbb{X}}_m, \hat{\mathbb{Y}}_m)$ as the test statistic: $\tilde{O}(4^g m)$ time
- Total runtime: $\tilde{O}(m)$ if $g = \log \log m$

Visual comparison on $\mathbb{P}^\star = 8$

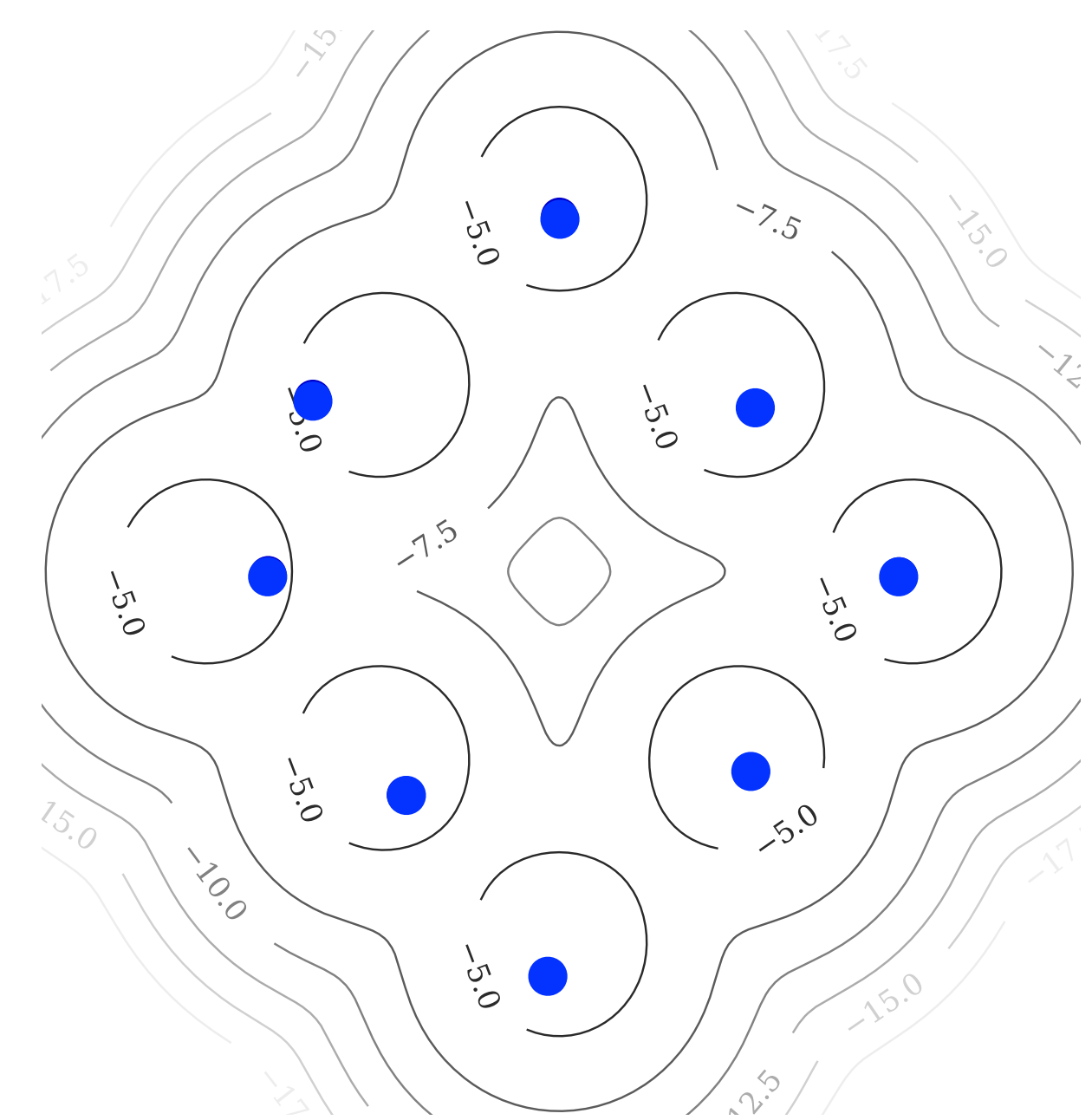
- iid input points

8 output points

Standard thinning



Kernel thinning



This one has rules marked for me



