

DATA-DRIVEN SIGNAL RECOVERY

Yury Polyanskiy (MIT)

joint with a (large) bunch of people

Talk plan

- 1.Thesis
- 2.Theory
- 3.Practice

- 1.Let's move beyond iid Gaussian noise
- 2.Likelihood-free hypothesis testing
- 3.Physics, Computer Vision, Comm

Classical detection and estimation

Adding a splash of 21st century

- ▶ How do we teach signal detection?

$$H_0 : Y_i \sim \mathcal{N}(1, \sigma^2) \quad H_1 : Y_i \sim \mathcal{N}(-1, \sigma^2)$$

and threshold the average $\frac{1}{m} \sum_i Y_i \geq 0$

- ▶ ... more generally:

$$H_0 : Y^m = s_0 + Z^m \quad H_1 : Y^m = s_1 + Z^m$$

and do matched filter: $(Y^m, s_1 - s_0) \geq 0$

- ▶ ... more generally:

$$H_0 : Y^m \sim P_{Y^m} \quad H_1 : Y^m \sim Q_{Y^m}$$

and do Neyman-Pearson

- ▶ ... more generally:

$$H_0 : Y^m \sim P, P \in \mathcal{P} \quad H_1 : Y^m \sim Q, Q \in \mathcal{Q}$$

and do what?..

- ▶ Try GLRT, otherwise search Annals of Stats
- ▶ **Problem:** if \mathcal{P}, \mathcal{Q} are realistic (i.e. large), then sample complexity is **bad** (curse of dimensionality etc)

Thesis: often we have side information (prior knowledge) about P_{Y^m}, Q_{Y^m} in the form of iid samples.

**Science:
simulations**

**Communication:
RF captures**

What is likelihood-free inference?

aka simulation-based inference

What is likelihood-free inference (LFI)?

aka simulation-based inference (SBI)

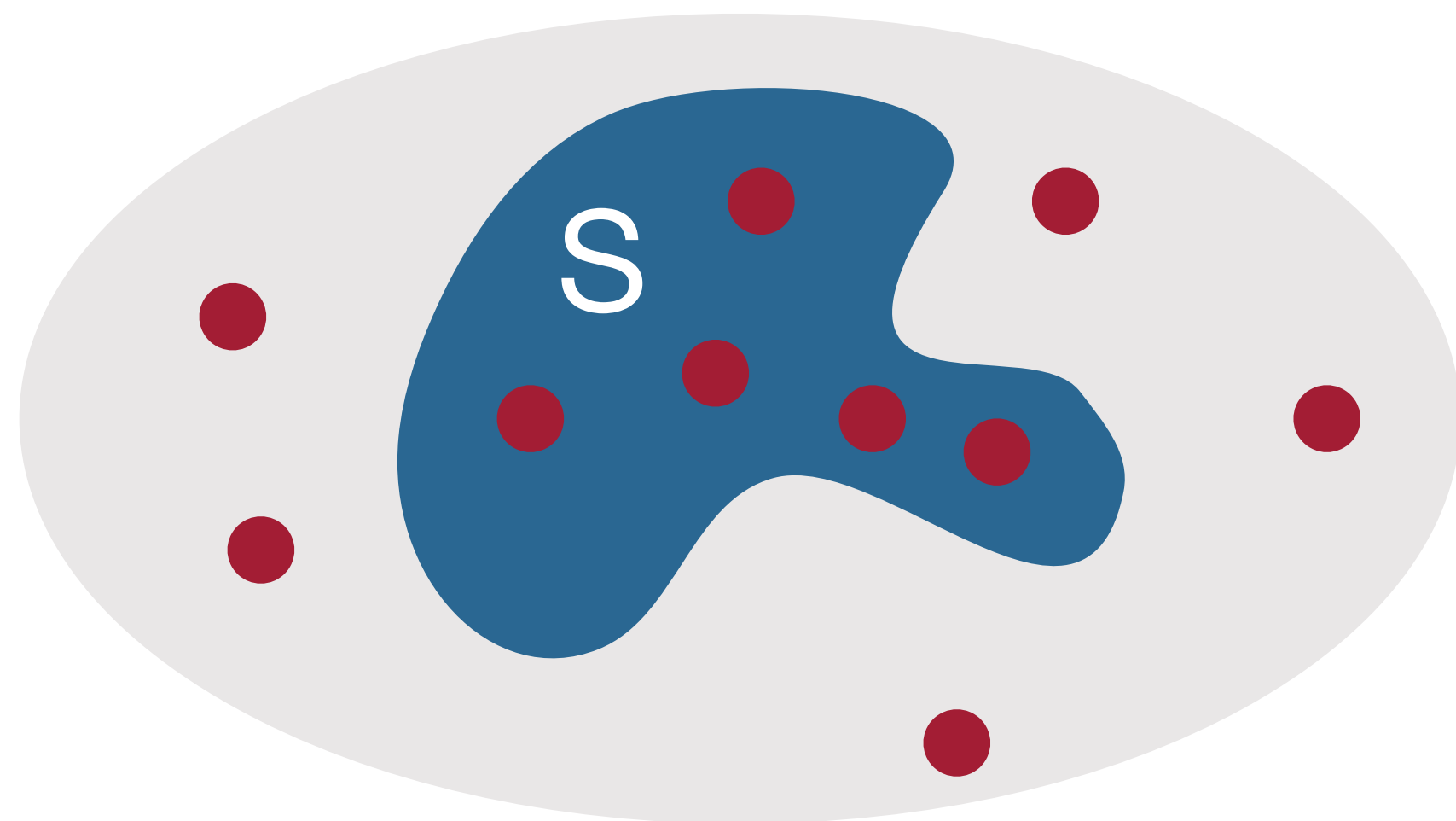
- ▶ Simulation access to **black-box** model $\theta \mapsto X \sim \mathbb{P}_\theta$
- ▶ Given true data $Z \sim \mathbb{P}_{\theta^\star}^{\otimes m}$, do **inference on θ^\star**
- ▶ Intractable likelihood: do so **without learning** the map $\theta \mapsto \mathbb{P}_\theta$
- ▶ Examples: climate modeling and particle physics

Discovery of the Higgs boson

- ▶ Observe data $Z \sim \mathbb{P}^{\otimes m}$

$$H_0 : \mathbb{P} = \mathbb{P}_{\text{noHiggs}} \quad \text{versus} \quad H_1 : \mathbb{P} = \mathbb{P}_{\text{Higgs}}$$

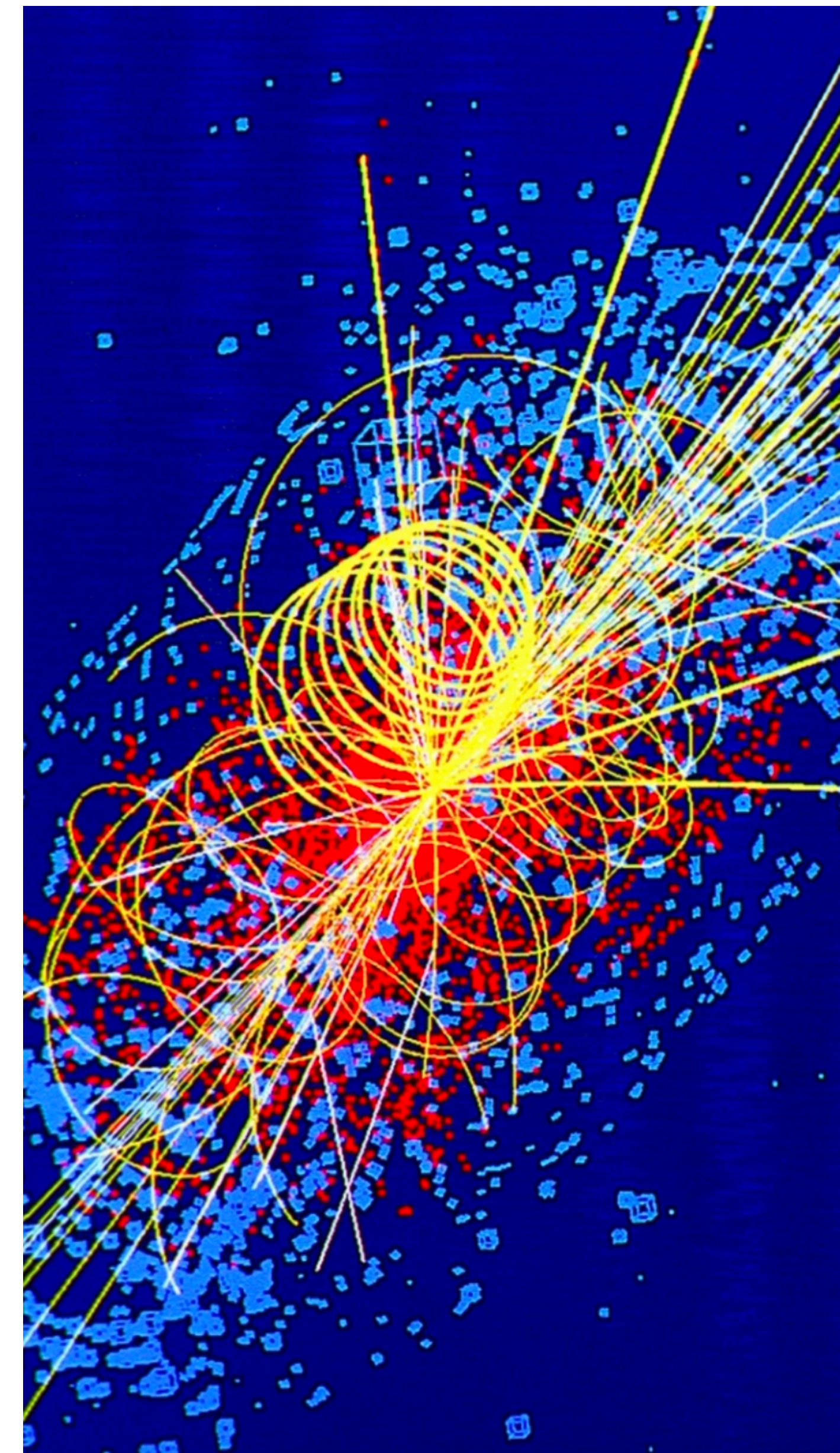
- ▶ Simulate $X \sim \mathbb{P}_{\text{noHiggs}}^{\otimes n}$ and $Y \sim \mathbb{P}_{\text{Higgs}}^{\otimes n}$



$$\bullet = Z_i$$

$S = \text{classifies}^* X \text{ vs } Y$

$$\text{Output} = \#\{\bullet \in S\} \leq \gamma$$



Simulation of the birth of a Higgs boson (CERN, Lucas Taylor)

*Boosted decision trees in the case of the Higgs boson discovery

Minimax setup

Likelihood-free hypothesis testing (LFHT)

Questions we will address:

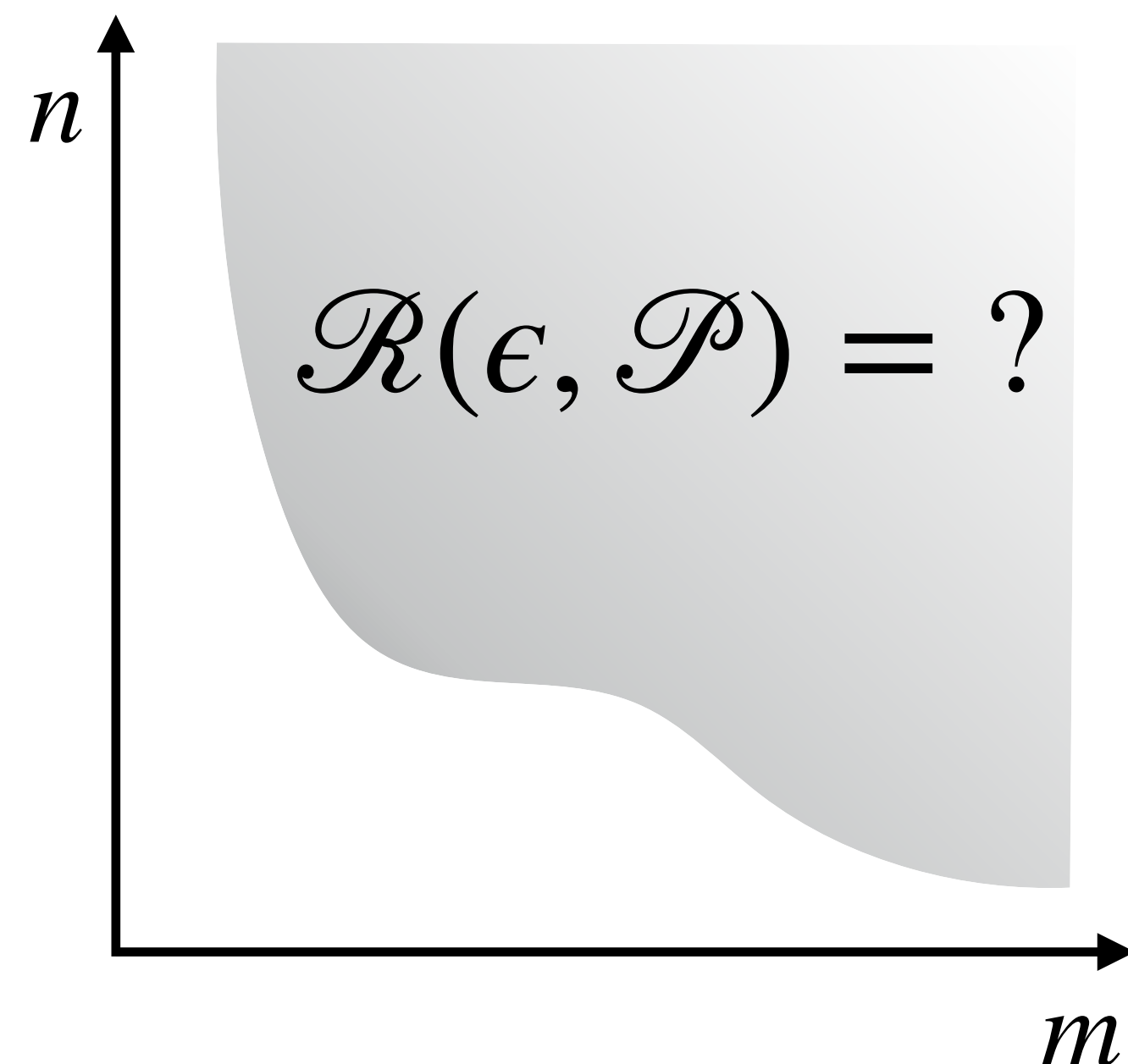
- 1) Fix
- 2) Let
 - Can we avoid learning distributions P_X, P_Y ?
 - Is there a tradeoff m vs n ?
- 3) Sim
- 4) Depending on H_0 or H_1 nature generates $Z \sim \mathbb{P}_X^{\otimes m}$ or $\mathbb{P}_Y^{\otimes m}$ respectively
- 5) Statistician observes $(X, Y, Z, \mathcal{P}, \epsilon)$ and decides H_0 or H_1

Likelihood-free hypothesis testing (LFHT)

$\mathcal{R}(\epsilon, \mathcal{P}) \subseteq \mathbb{N}^2$ is set of (m, n) s.t. exists test that given X, Y, Z performs

$$H_0 : \mathbb{P}_X = \mathbb{P}_Z \quad \text{versus} \quad H_1 : \mathbb{P}_Y = \mathbb{P}_Z$$

with (Type-I + Type II error) $< 1\%$.



278

IEEE TRANSACTIONS ON INFORMATION THEORY, VOL. 34, NO. 2, MARCH 1988

On Classification with Empirically Observed Statistics and Universal Data Compression

JACOB ZIV, FELLOW, IEEE

Michael Gutman
Aaron Wagner
Sean Meyn
Ashish Khisti
Vincent Tan

....

- ▶ ... and other prior work. But only for **discrete distributions** and fixed

$$\text{TV}(\mathbb{P}_X, \mathbb{P}_Y) \asymp 1, \quad m, n \rightarrow \infty$$

Statistical Problems

$$\begin{array}{l} X \sim \mathbb{P}_X^{\otimes n} \\ Y \sim \mathbb{P}_Y^{\otimes n} \\ Z \sim \mathbb{P}_Z^{\otimes m} \end{array} \begin{array}{l} \swarrow \\ \leftarrow \\ \swarrow \end{array} \text{unknown, } \mathbb{P}_0 \text{ known, all in } \mathcal{P}$$

The classes \mathcal{P}

Choices of \mathcal{P} we considered:

- ▶ $\mathcal{P}_H(\beta, d) = \{\beta\text{-H\"older densities over } [0,1]^d \text{ with } \|\cdot\|_{\mathcal{C}^\beta} \leq C_H\}$

aka β times differentiable densities.

- ▶ $\mathcal{P}_G(s) = \{\otimes_{j=1}^d \mathcal{N}(\theta_j, 1) : \theta_j \text{ in Sobolev ellipsoid } \sum_j \theta_j^2 j^{2s} \leq C_G\}$

- ▶ $\mathcal{P}_{Db}(k) = \{\text{discrete distributions on } [k] \text{ with } \|\cdot\|_{TV} \leq C_{Db}/k\}$

This talk: focus on \mathcal{P}_H (smooth densities)

- ▶ $\mathcal{P}(k) = \{\text{all discrete distributions on } [k]\}$

- ▶ arbitrary densities on $[0,1]^d$ (with MMD separation instead of TV)

Rates vs sample complexity

Famous results for $\mathcal{P}_H(\beta, d)$

	Rate	Sample complexity
Goodness-of-fit	$n^{-\frac{\beta}{2\beta + d/2}}$	$\epsilon^{-\frac{2\beta + d/2}{\beta}} = n_{\text{GoF}}$
Estimation	$n^{-\frac{\beta}{2\beta + d}}$	$\epsilon^{-\frac{2\beta + d}{\beta}} = n_{\text{Est}}$

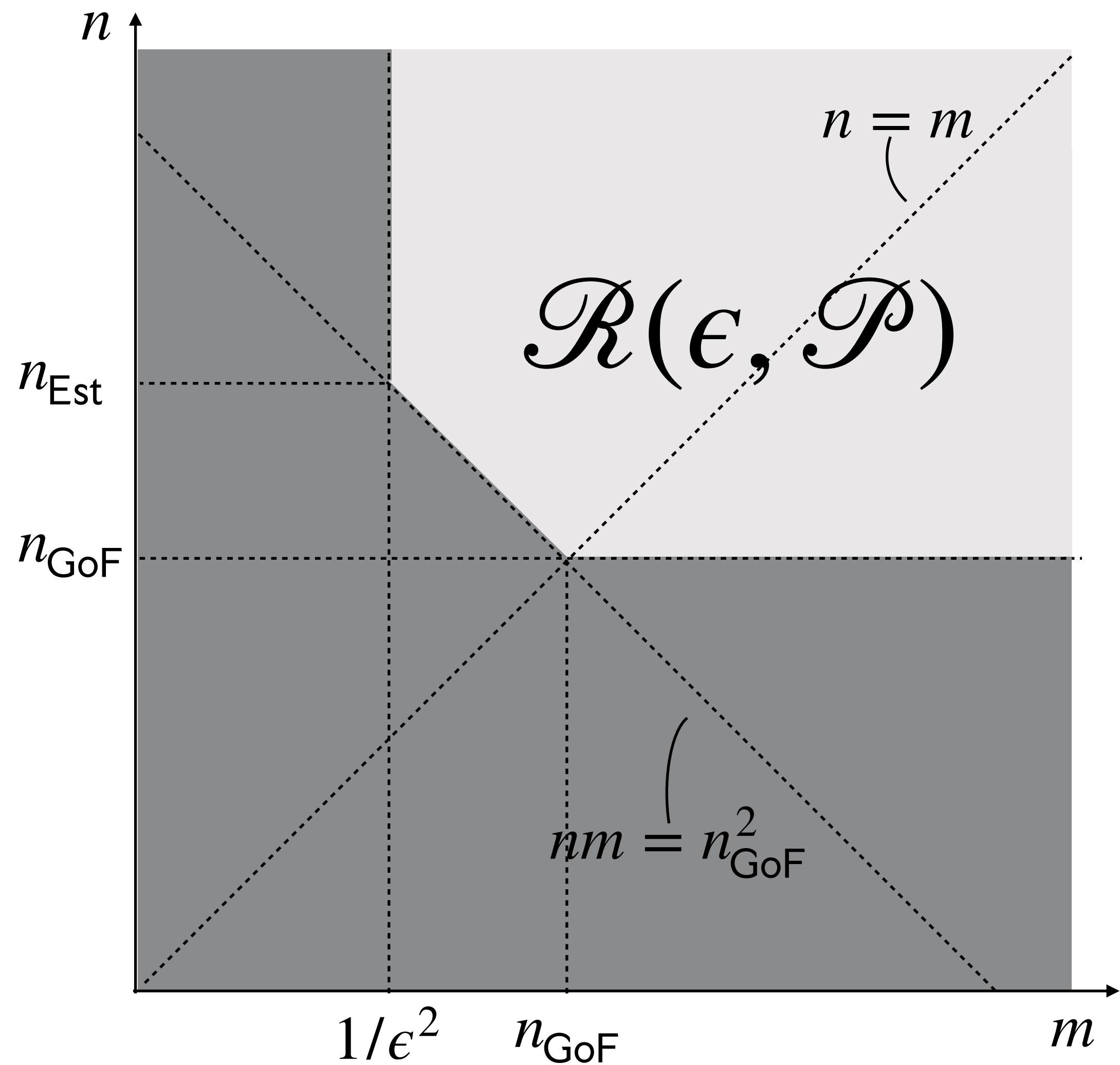
Results for \mathcal{P}_H , \mathcal{P}_G and \mathcal{P}_{Db}

Theorem (Gerber-P.'2022)

Up to constant factors:

$$\mathcal{R}(\epsilon, \mathcal{P}) \asymp \left\{ \begin{array}{l} m \geq 1/\epsilon^2 \text{ and } n \geq n_{\text{GoF}} \\ \text{and } n \cdot m \geq n_{\text{GoF}}^2 \end{array} \right\}$$

Interpreting the results



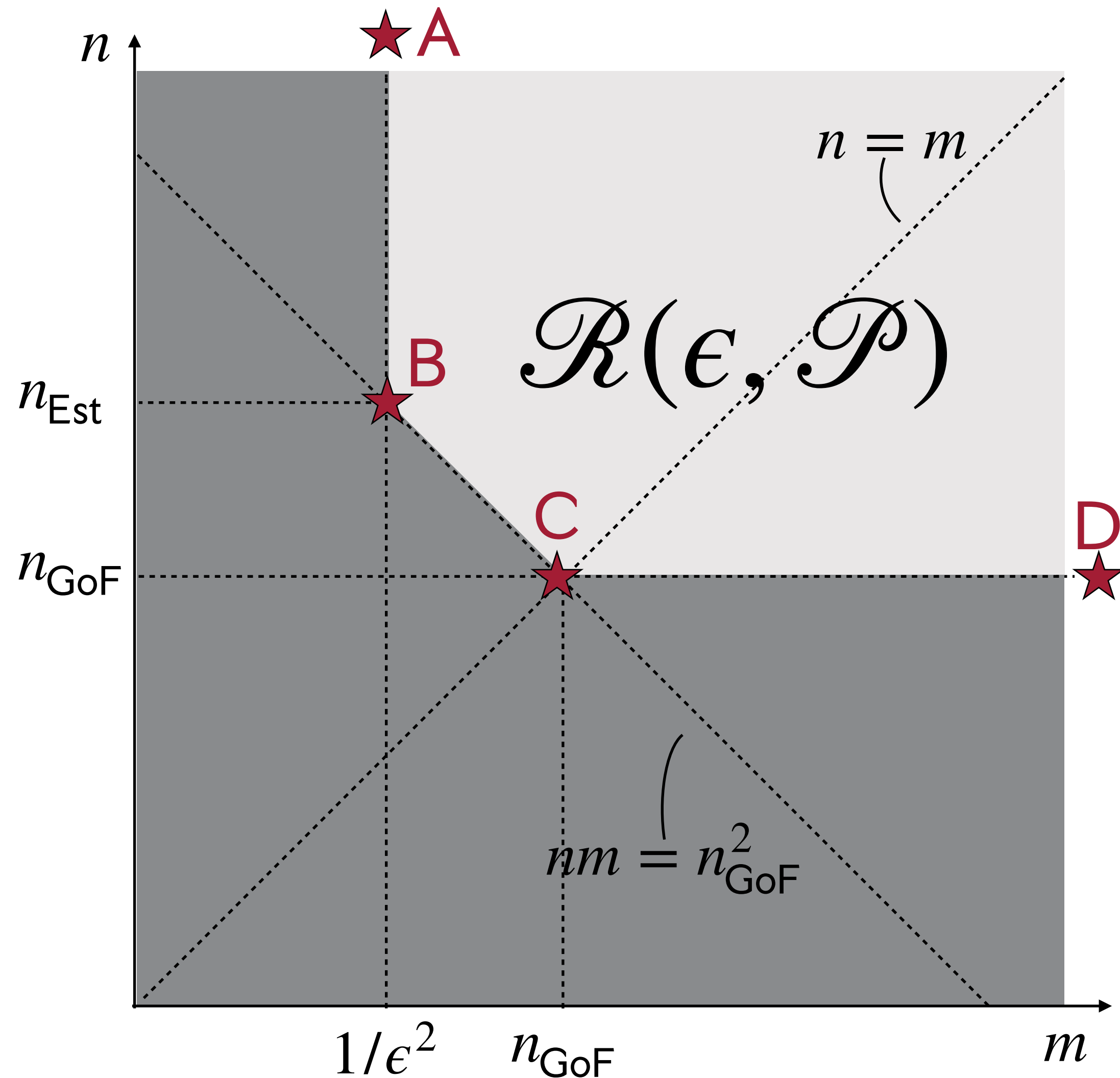
$$\mathcal{R}(\epsilon, \mathcal{P}) \asymp \left\{ \begin{array}{l} m \geq 1/\epsilon^2 \text{ and } n \geq n_{\text{GoF}} \\ \text{and } n \cdot m \geq n_{\text{GoF}}^2 \end{array} \right\}$$

	n_{GoF}	n_{Est}
$\mathcal{P}_H(\beta, d)$	$\epsilon^{-\frac{2\beta + d/2}{\beta}}$	$\epsilon^{-\frac{2\beta + d}{\beta}}$

Target: minimal m (as in Higgs)

$$n_{\text{Est}} = n_{\text{GoF}}^2 \epsilon^2$$

Interpreting the results



$$\mathcal{R}(\epsilon, \mathcal{P}) \asymp \left\{ \begin{array}{l} m \geq 1/\epsilon^2 \text{ and } n \geq n_{\text{GoF}} \\ \text{and } n \cdot m \geq n_{\text{GoF}}^2 \end{array} \right\}$$

Point	Algorithm	Lower bd
$A \leftrightarrow (1/\epsilon^2, \infty)$	Binary HT	Trivial
$B \leftrightarrow (1/\epsilon^2, n_{\text{Est}})$	Est + robust HT	New
$C \leftrightarrow (n_{\text{TS}}, n_{\text{TS}})$	Two-sample*	Reduction to TS
$D \leftrightarrow (\infty, n_{\text{GoF}})$	Goodness-of-fit	New but easy

Can estimate \mathbb{P}_X and \mathbb{P}_Y

15 * $n_{\text{GoF}} = n_{\text{TS}}$ for each of these classes

The test statistic

- ▶ Based on Ingster's L^2 -comparison idea
- ▶ Discretize $[0,1]^d$ cube into $k = \epsilon^{-\frac{d}{\beta}}$ bins
- ▶ Empirical pmfs $\hat{p}_X, \hat{p}_Y, \hat{p}_Z$ based on (n, n, m) observations
- ▶ **Theorem:** All points on the optimal tradeoff are achieved by

$$T = \|\hat{p}_X - \hat{p}_Z\|_2^2 - \|\hat{p}_Y - \hat{p}_Z\|_2^2$$

$$\Psi = \mathbb{I}\{T \geq 0\}$$

Note: no training, distance estimates are both wrong.

Enter Machine Learning: Practical tests

Kernel-based L_2 test (MMD)

- ▶ Real-world distributions are high-dimensional \implies discretization impractical.

- ▶ Given $\hat{\mathbb{P}}_X, \hat{\mathbb{P}}_Z$ measure distance after applying **feature map** ϕ :

$$\text{MMD}^2(\hat{\mathbb{P}}_X, \hat{\mathbb{P}}_Z) = \|\hat{\mathbb{E}}\phi(X) - \hat{\mathbb{E}}\phi(Z)\|_2^2$$

(proposed for two-sample testing [[Sutherland et al, ICLR'17](#)])

- ▶ We adopt this to LFHT via the same mechanism:

$$T(X, Y, Z) = \|\hat{\mathbb{E}}\phi(X) - \hat{\mathbb{E}}\phi(Z)\|_2^2 - \|\hat{\mathbb{E}}\phi(Y) - \hat{\mathbb{E}}\phi(Z)\|_2^2 \quad \geq 0$$

- ▶ Has the same LFHT region wrt $\text{MMD}(P_X, P_Y) \geq \epsilon$ [[Gerber, Jiang, Sun, P., NeurIPS'23](#)]

- ▶ **Train feature map** to maximize $\frac{\mathbb{E}[T | H_0]}{\sqrt{\text{Var}[T | H_0]}}$ ratio (gradient descent in kernel space)

LFHT for CIFAR

[NeurIPS'23]

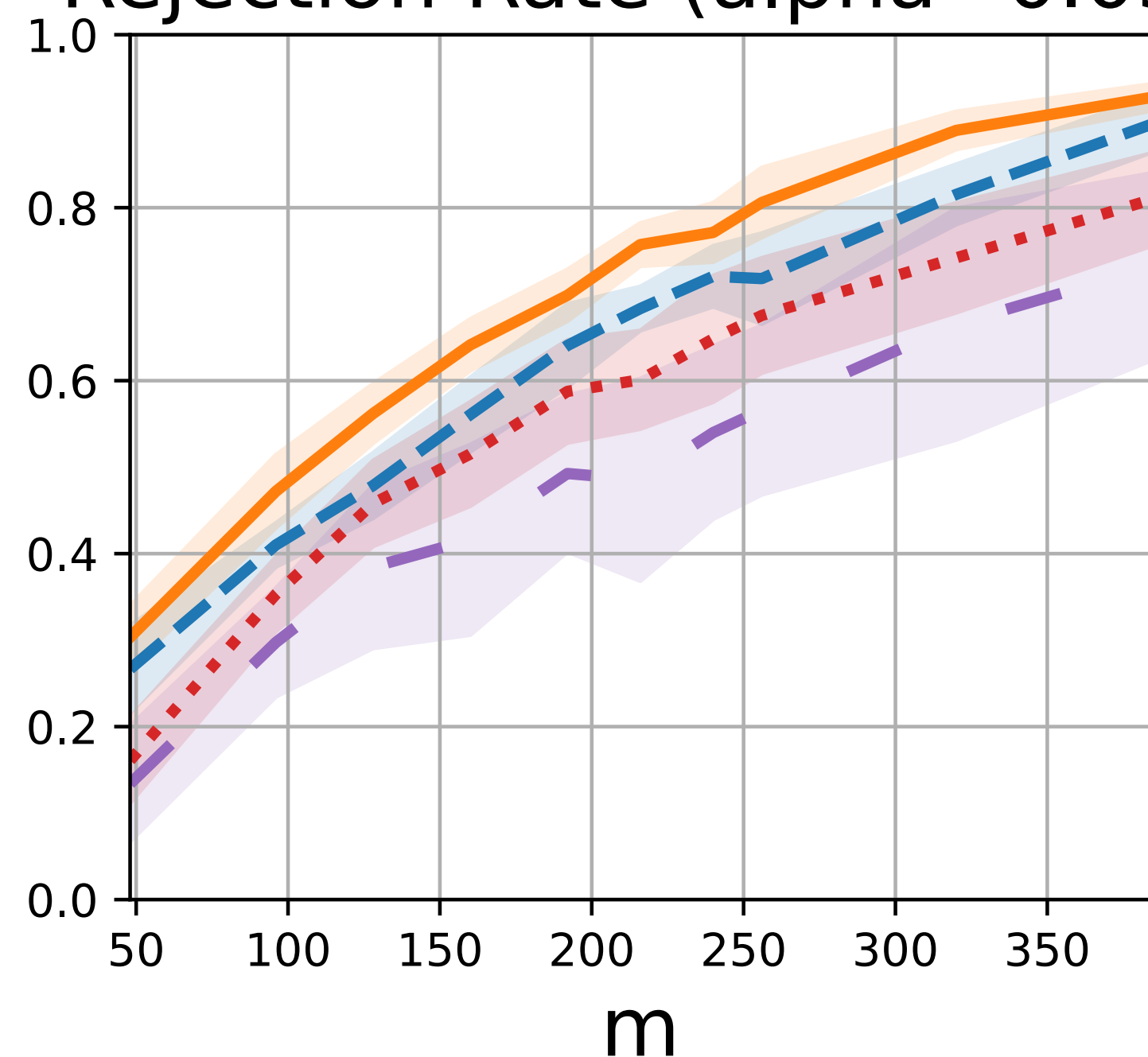
- ▶ So our test:

$$T(X, Y, Z) = \|\hat{E}\phi(X) - \hat{E}\phi(Z)\|_2^2 - \|\hat{E}\phi(Y) - \hat{E}\phi(Z)\|_2^2 \geq 0$$

- ▶ Here is an example: **X=CIFAR10 vs Y=1/3 CIFAR + 2/3 Diffusion Model (DDPN)**

- ▶ ($n \approx 10^5, m \approx 10^1$)

Rejection Rate (alpha=0.05)



Can detect fakes with accuracy 90% from about 300 examples.

Best of 4 types of tests = MMD

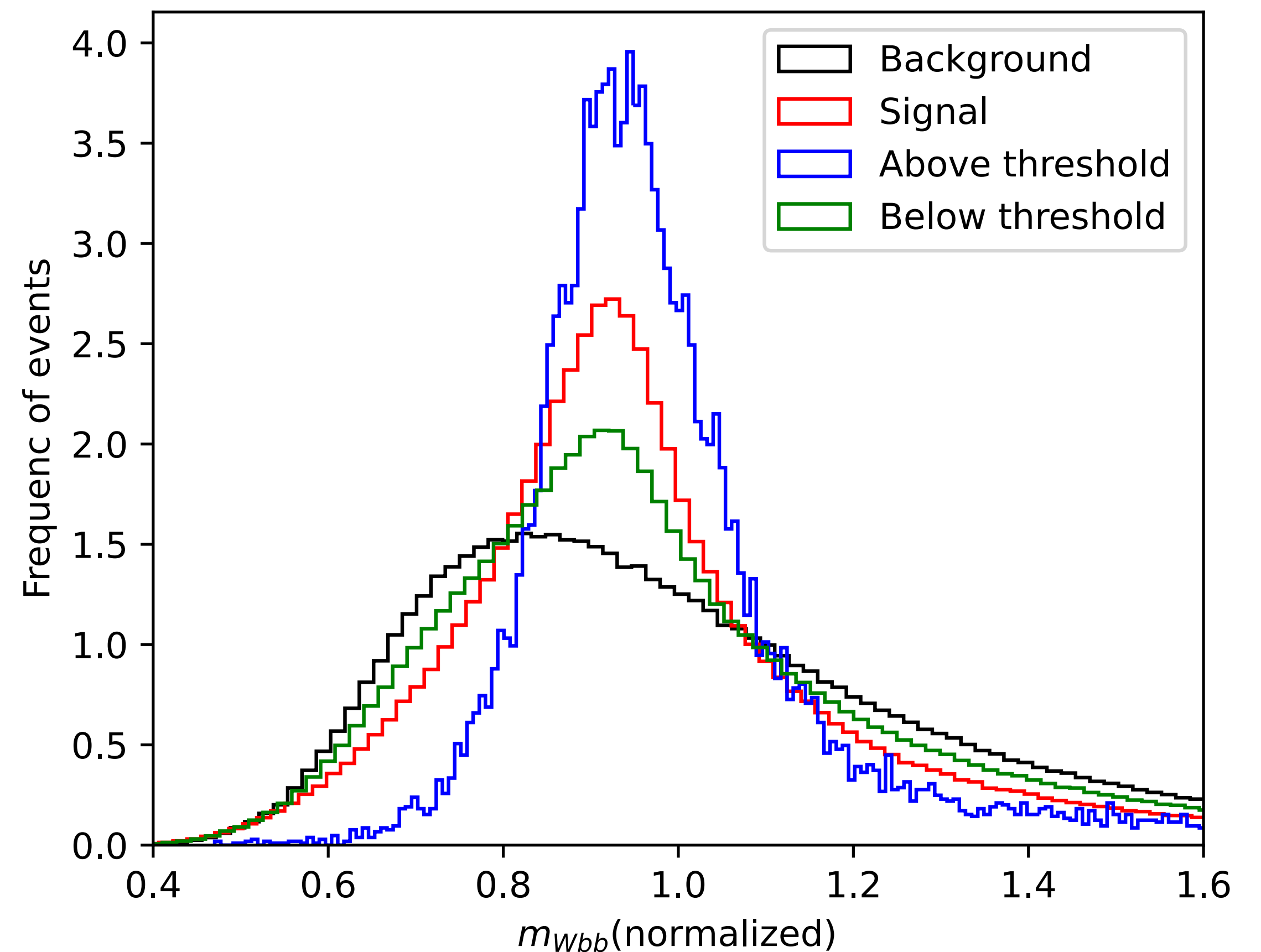
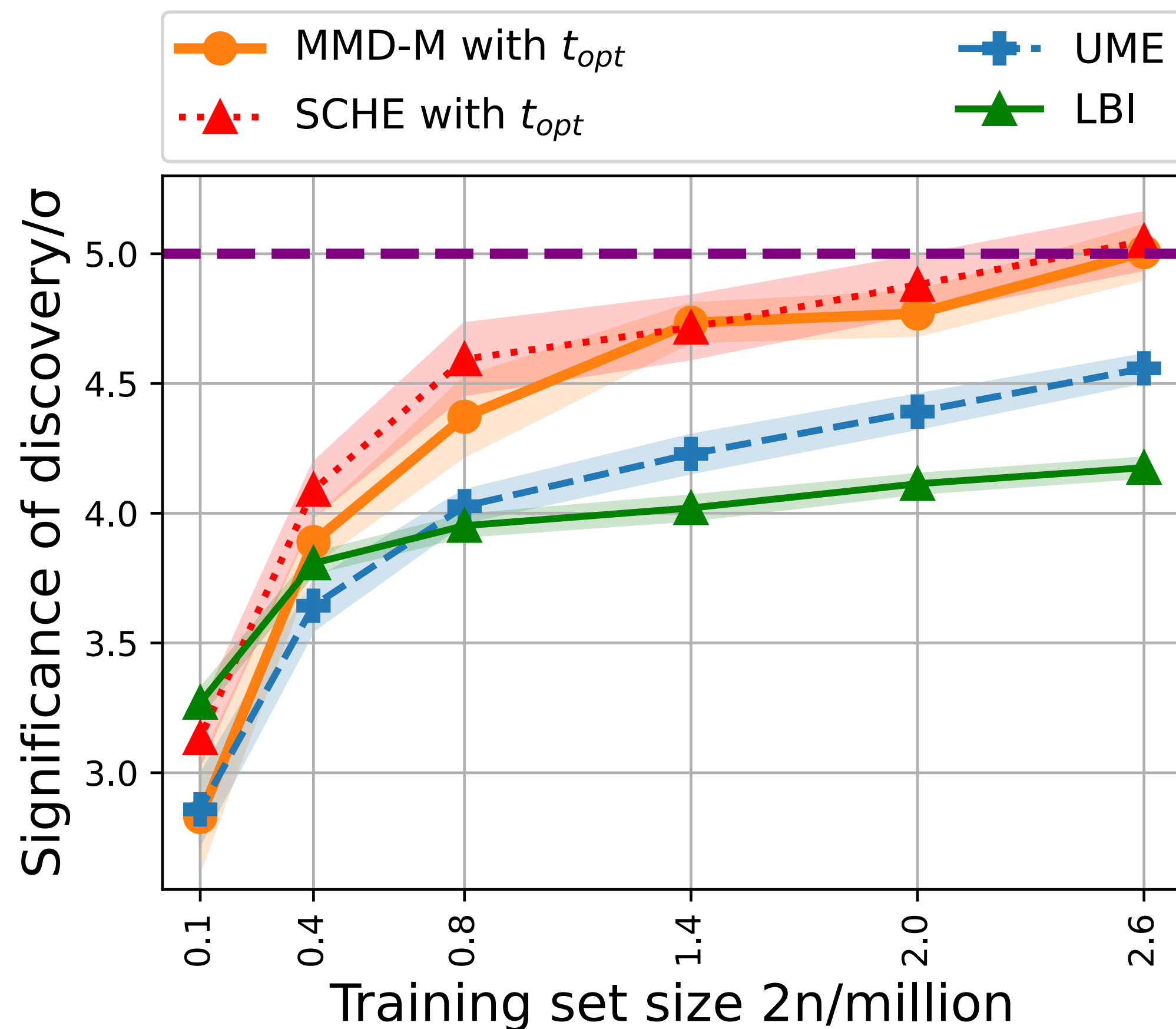
— MMD-M - - - UME ···· SCH - - - MMD-G

Some examples of bad diffusion images



Back to Higgs [NeurIPS'23]

- ▶ Instead of fixed two-sided error physicists use significance of discovery
- ▶ Expressed in σ 's. For the *new particle* need 5σ . Our road to 5σ ...

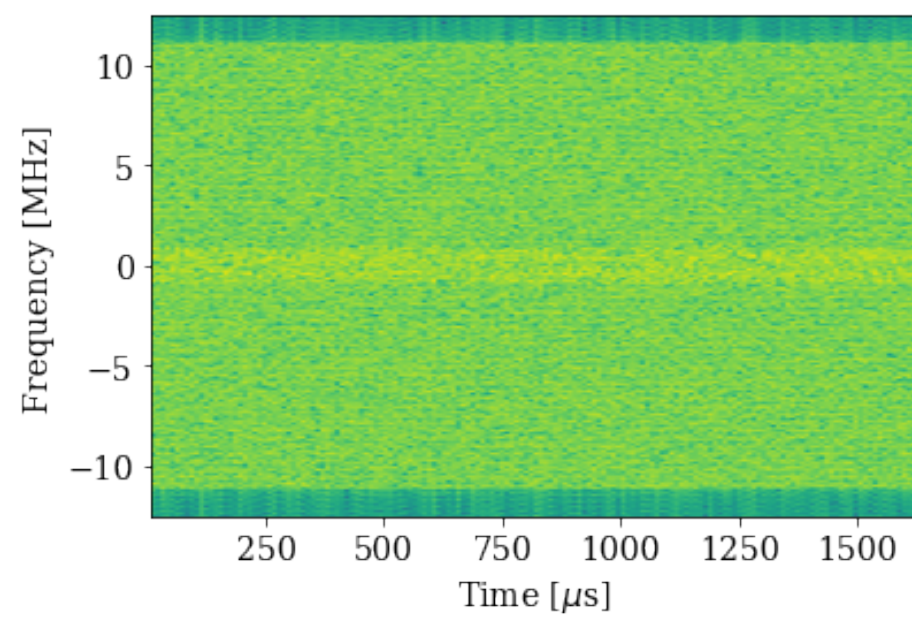
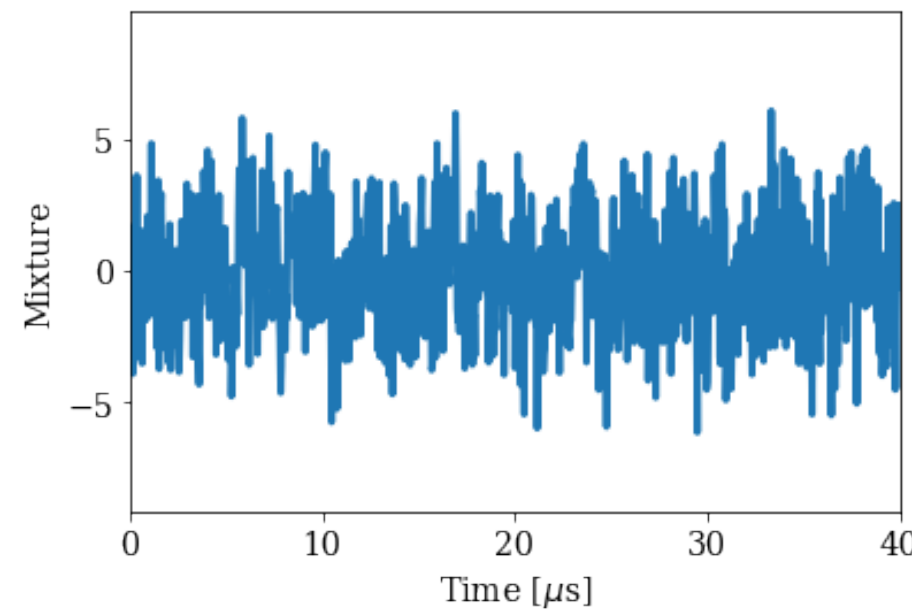


Interference rejection

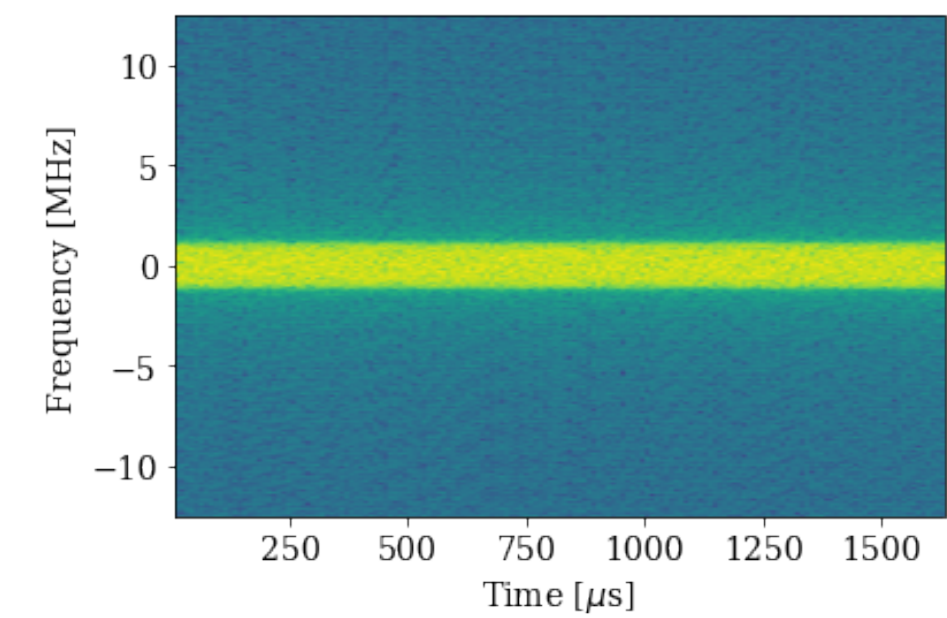
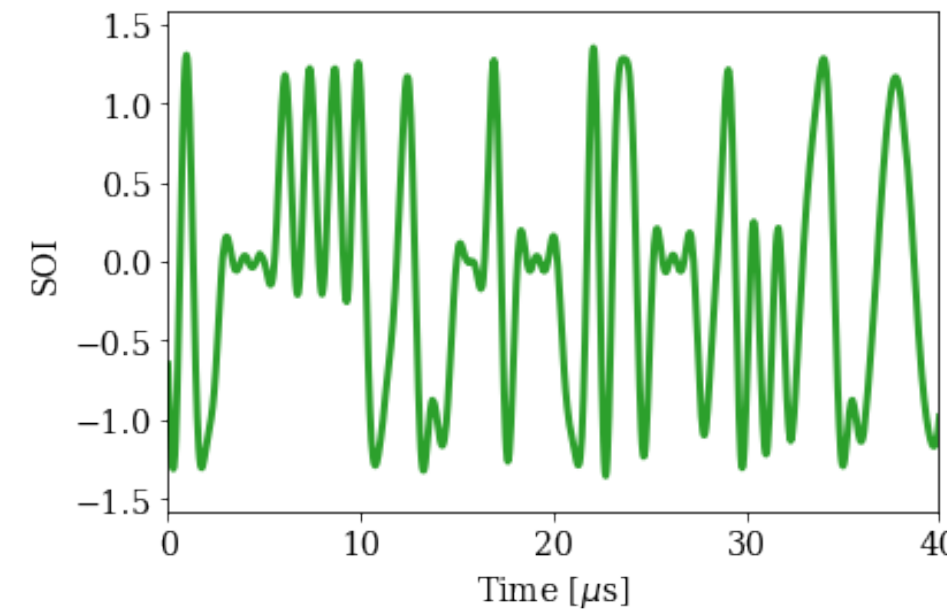
Demodulation task in communication

$$y = s + b$$

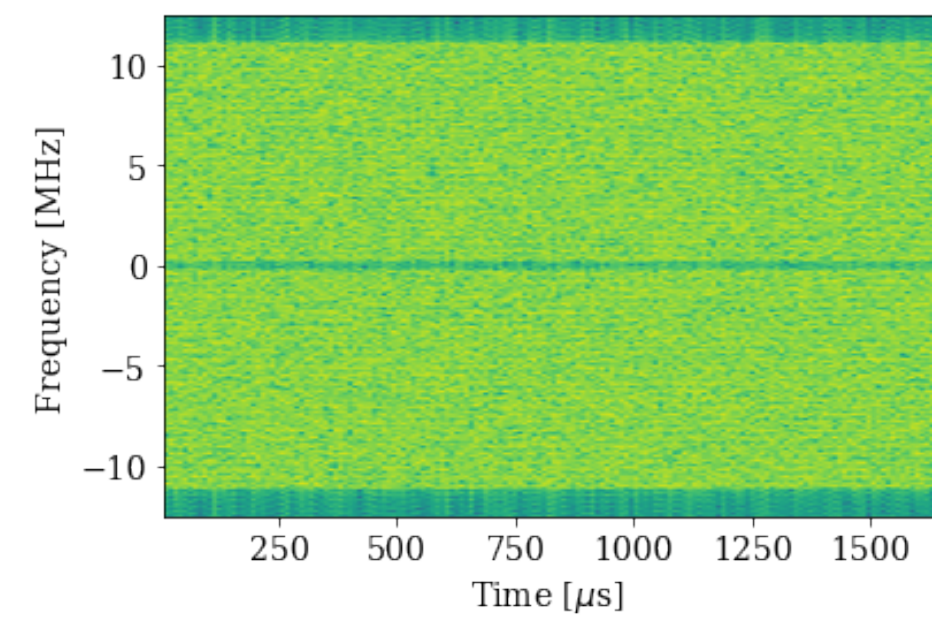
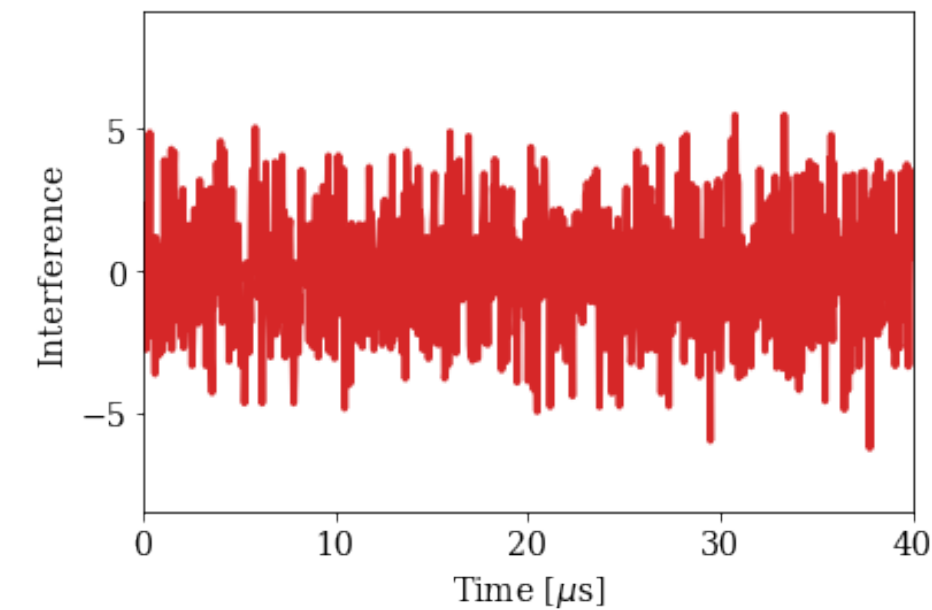
Received signal



*Signal of Interest (SOI)
e.g. BPSK/QPSK*



Noise and interference



Example at -9 dB Signal-to-Interference Ratio (SIR)

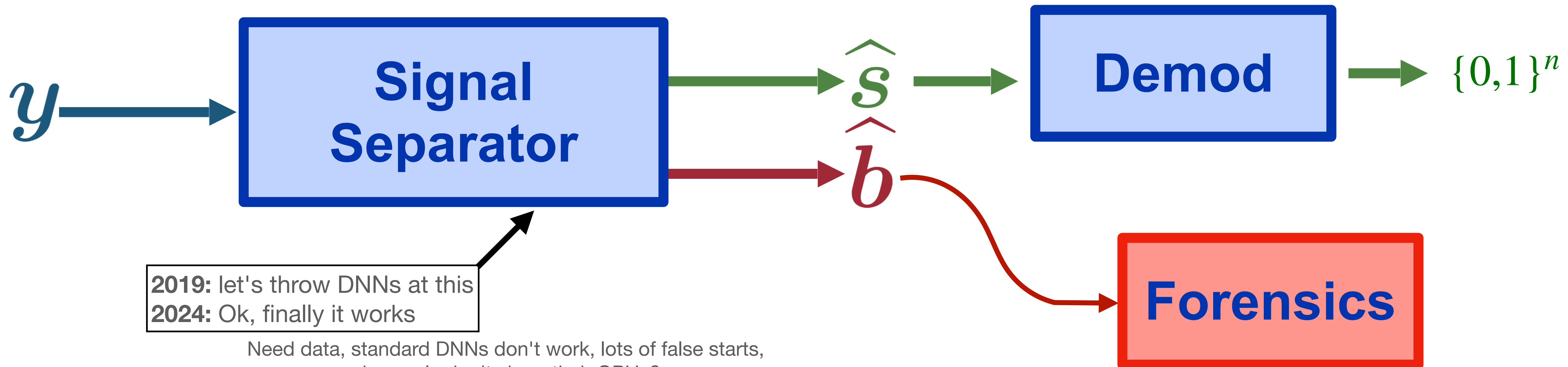
OFDM interference is marginally Gaussian => need to exploit time-frequency structure of the interference. How?

Idea: Use signal (source) separation



$$y = s + b$$

Observed *SOI* *Interference*



2019: let's throw DNNs at this
2024: Ok, finally it works

Need data, standard DNNs don't work, lots of false starts,
oh people don't share their GPUs?

Two types of architectures

$$\underset{\text{Observed}}{y} = \underset{\text{SOI}}{s} + \underset{\text{Interference}}{b}$$



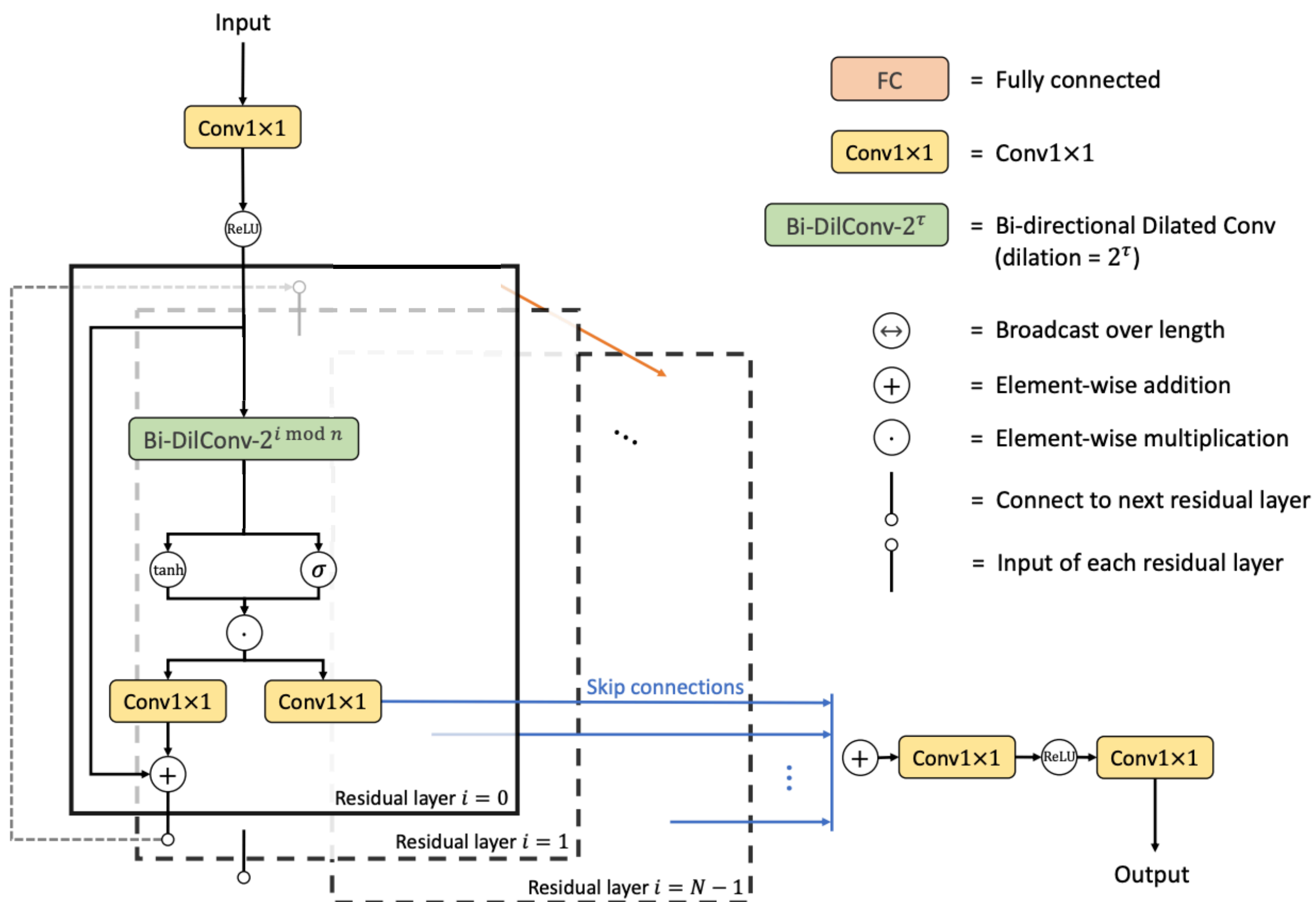
Supervised (end-to-end)

- Create many synthetic mixtures $s + b$
- Feed pairs (y, s) to DNN
- Force it to learn to recover s from y
- **Pros:** best performance
- **Cons:** need to retrain DNN for each signal-interference pair

Bayesian MAP

- Collect many samples of b
- Train a diffusion model to learn P_b
- Use MAP to recover s from y
- **Pros:** one model works for all SOI
- **Cons:** slow inference, performance

NeurIPS'2023: WaveNet (dilated CNN)



	Description
Number of layers	30 residual layers with dilation cycle of {1, 2, ... 512} repeated three times
Total number of parameters	4M
GPU compute (training)	8 GPU days

Additional training tricks:

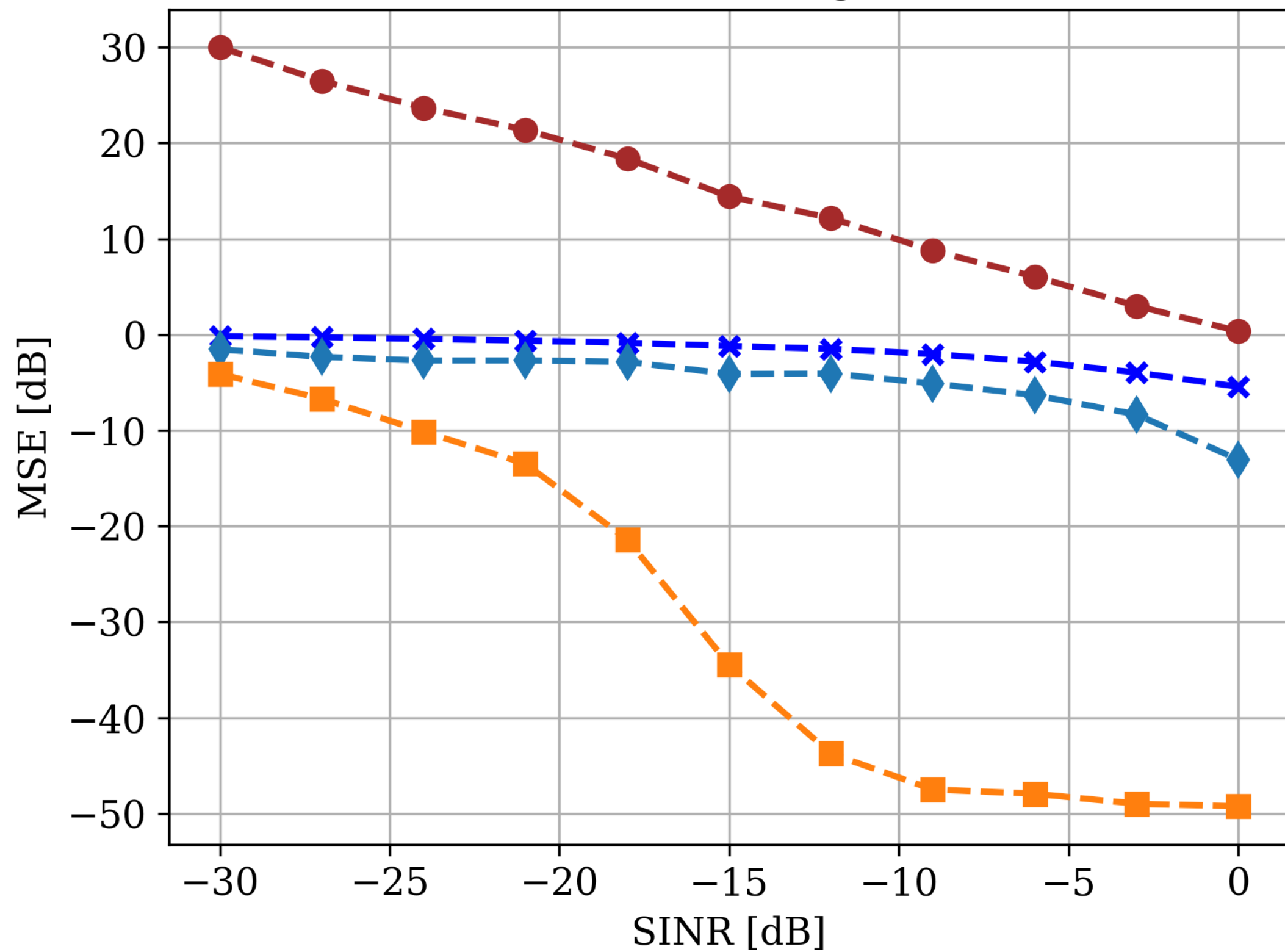
Adaptive learning rate scheduler based on validation loss

Mixed precision training with fp16 to speed up inference

So does it work?

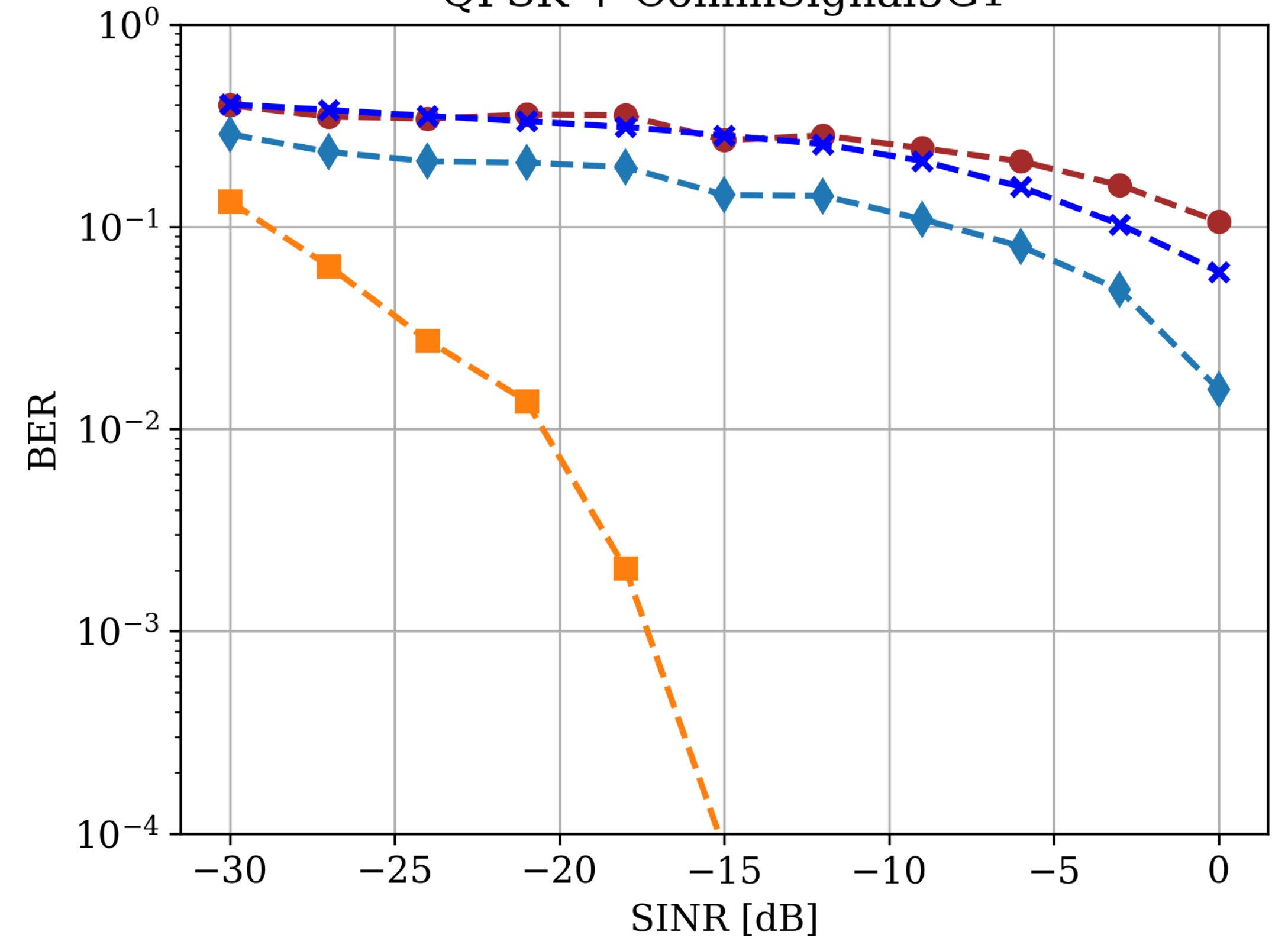
QPSK vs OFDM (5GNR) example

QPSK + CommSignal5G1



- Matched Filter Demod Only (No Mitigation)
- ×— LMMSE Separator + Matched Filter Demod
- ◆— UNet Separator + Matched Filter Demod
- WaveNet Separator + Matched Filter Demod

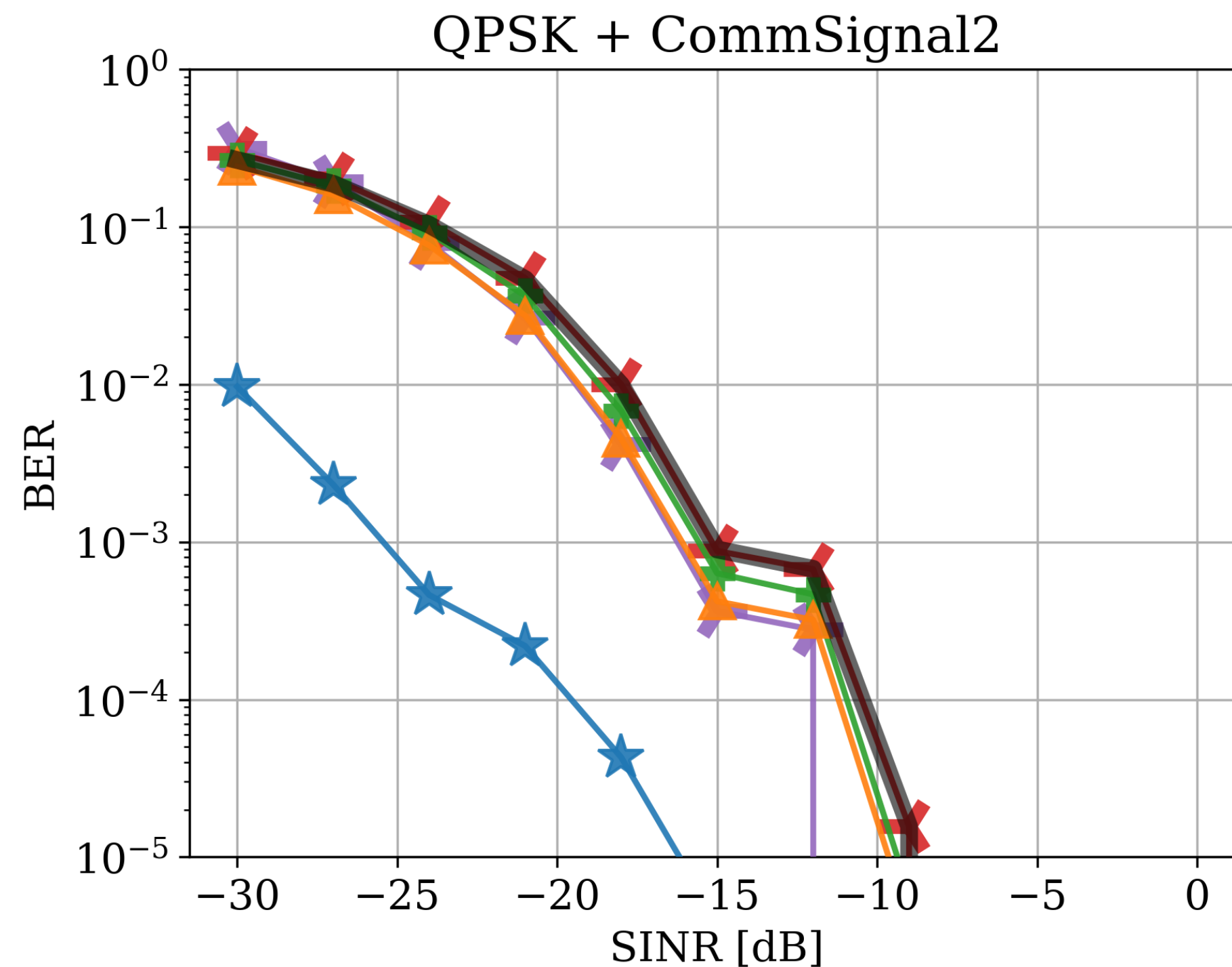
QPSK + CommSignal5G1



- Matched Filter Demod Only (No Mitigation)
- ×— LMMSE Separator + Matched Filter Demod
- ◆— UNet Separator + Matched Filter Demod
- WaveNet Separator + Matched Filter Demod

ICASSP'2024: Session on RF Challenge

Can we obtain further gains from other novel architectures?



- WaveNet Baseline
- ★ KU-TII
- ▲ OneInAMillion
- TUB
- ✦ LHen
- ✧ imec_IDLab_Wireless_UGent



Learnable dilations and new data augmentation schemes

Number of parameters: **16M**
 Number of GPUs: 4 x RTX 3090
 GPU Compute: **13 GPU days**



Attention-based UNet and fine-tuning of our WaveNet baseline

Number of parameters: **350M**
 Number of GPUs: 4 x A100
 GPU Compute: **8 GPU days**



New UNet architecture with bi-directional LSTM bottleneck layer

Number of parameters: **60M**
 Number of GPUs: 1 x RTX 6000
 GPU Compute: **4 GPU days**

Two types of architectures

$$\underset{\text{Observed}}{y} = \underset{\text{SOI}}{s} + \underset{\text{Interference}}{b}$$



Supervised (end-to-end)

- Create many synthetic mixtures $s + b$
- Feed pairs (y, s) to DNN
- Force it to learn to recover s from y
- **Pros:** best performance
- **Cons:** need to retrain DNN for each signal-interference pair

Bayesian MAP

- Collect many samples of b
- Train a **diffusion model** to learn P_b
- Use MAP to recover s from y
- **Pros:** one model works for all SOI
- **Cons:** slow inference, performance

Diffusion models

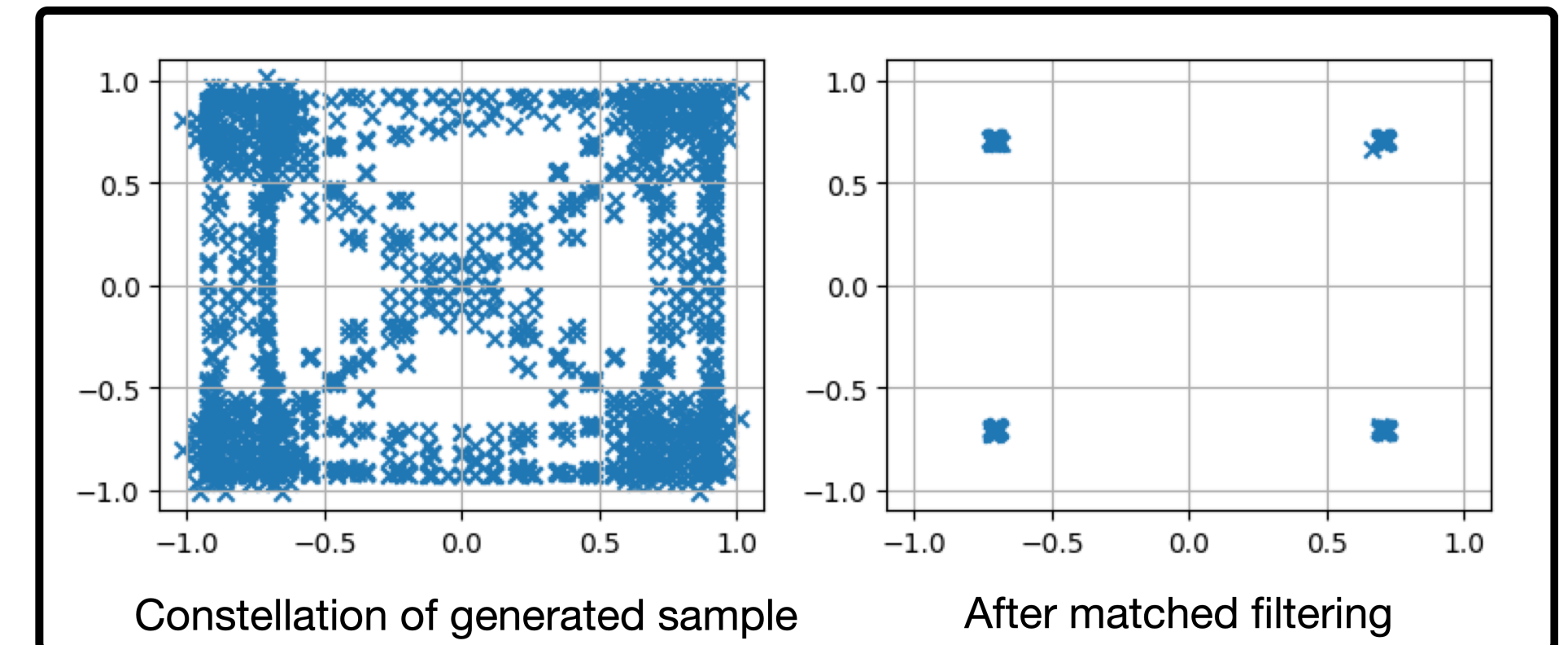
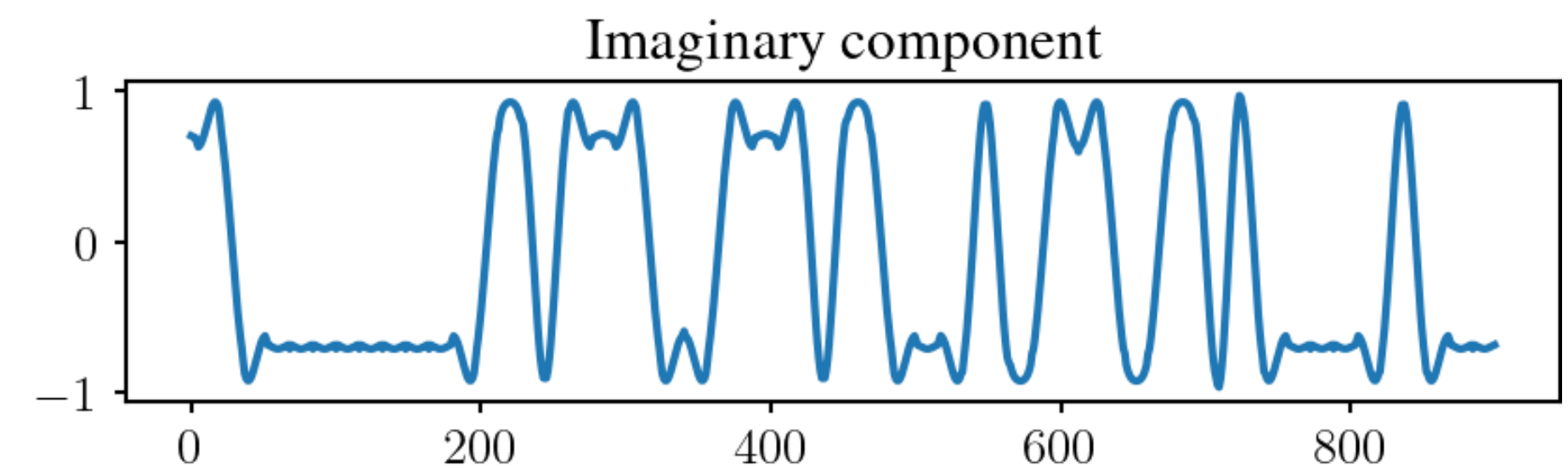
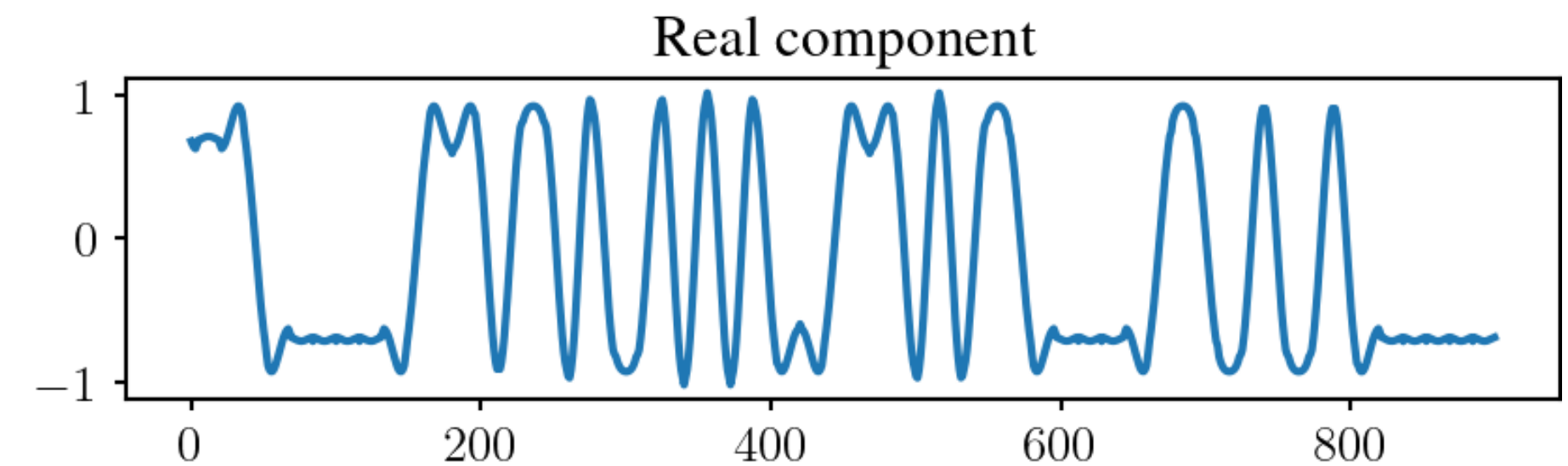
Images and RF

SOTA *generative model* that can learn complex structures from signal datasets



Can diffusion models capture the underlying discrete statistical structures of RF signals?

Sample from diffusion model trained on QPSK signals

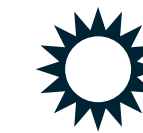


Score-based Source separation (α -RGS)

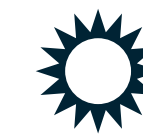
$\mathbf{s} \in \mathcal{S} \subset \mathbb{C}^D, \mathbf{b} \in \mathbb{C}^D$ statistically independent sources

MAP Estimation Given $\mathbf{y} = \mathbf{s} + \mathbf{b}$

$$\hat{\mathbf{s}} = \arg \max_{\mathbf{s} \in \mathcal{S}: \mathbf{y} = \mathbf{s} + \mathbf{b}} p_{\mathbf{s}|\mathbf{y}}(\mathbf{s}|\mathbf{y}) = \arg \min_{\mathbf{s} \in \mathcal{S}} -\log P_{\mathbf{s}}(\mathbf{s}) - \log p_{\mathbf{b}}(\mathbf{y} - \mathbf{s})$$



Combinatorially hard



Non-differentiable

Gradient Descent Estimate $\bar{\mathbf{s}} = \mathbf{s} + \epsilon, \epsilon \rightarrow 0$

$$\mathbf{s}_{i+1} \leftarrow \mathbf{s}_i + \underbrace{\nabla \log p_{\mathbf{s}}(\mathbf{s}_i)}_{\text{Score}} - \nabla \log p_{\mathbf{b}}(\mathbf{y} - \mathbf{s}_i)$$

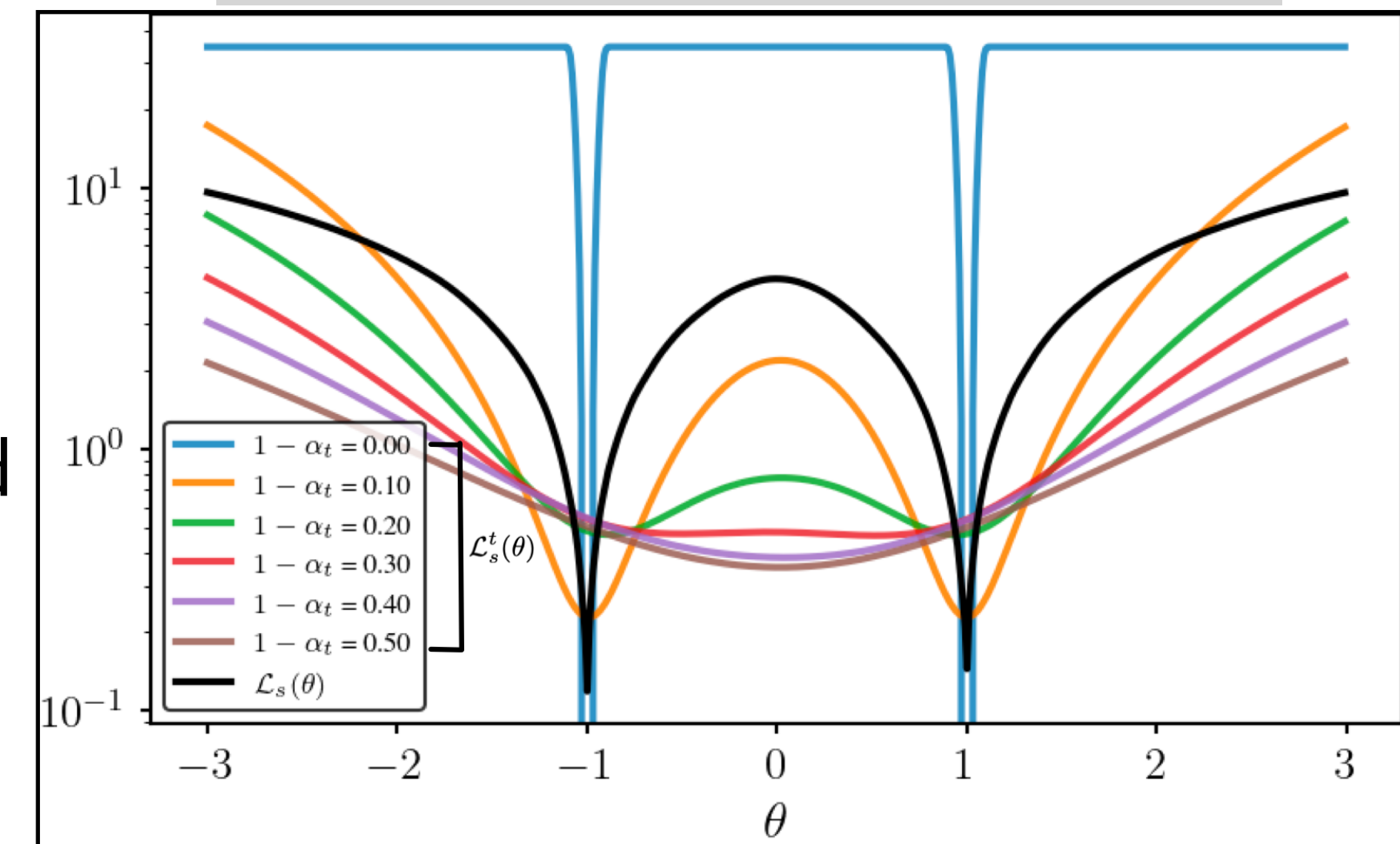
Randomized Gaussian Smoothing with an α -posterior (α -RGS)

Diffusion Models Model **unknown priors (score functions)** over \mathbf{s} and

Gaussian Smoothing Use noise variance levels α_t and α_u

α -posterior Reweight likelihood with weight $\alpha = \omega$

Smoothed optimization landscape

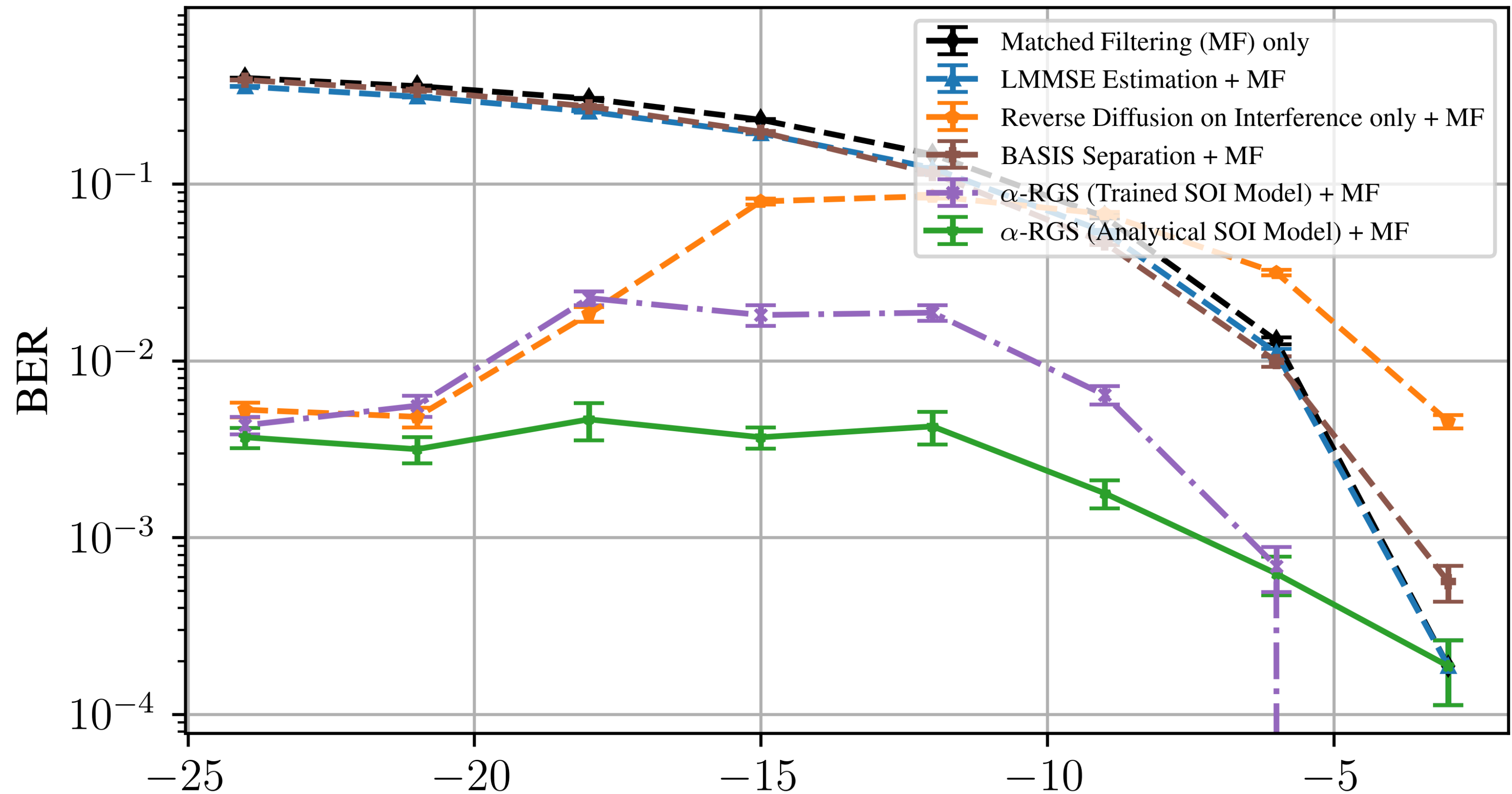


$$\mathcal{L}(\theta) \triangleq -\mathbb{E}_{t, \mathbf{z}_s} \left[\log p_{\tilde{\mathbf{s}}_t} \left(\tilde{\mathbf{s}}_t(\theta) \right) \right] - \omega \mathbb{E}_{u, \mathbf{z}_b} \left[\log p_{\tilde{\mathbf{b}}_u} \left(\tilde{\mathbf{b}}_u(\theta, \mathbf{y}) \right) \right]$$

Results: improving SOTA

Other algos based on approximating MAP via score-learning

RRC-QPSK SOI + OFDM (QPSK) Interference



Averaging over regularization + α -posterior give us an edge

Conclusion

- (i)* We studied signal detection (hypothesis testing) when hypotheses are only specified through examples.
- (ii)* We saw minimax optimal bounds and practical algorithms
- (iii)* **Next** : Study notion of regret or instance-optimality.
- (iv)* **More generally:** Study parameter estimation, confidence intervals, channel coding, constellation design,...

Thank you!