

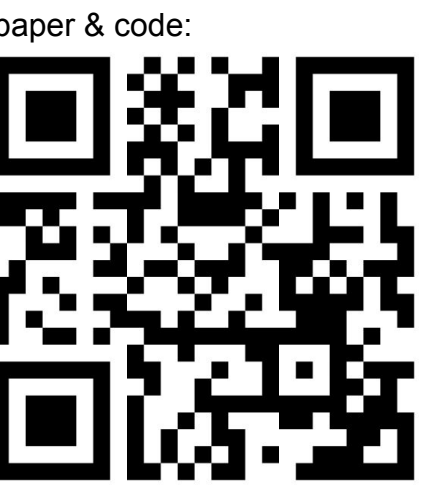
# Estimating the Rate-Distortion Function by Wasserstein Gradient Descent

Yibo Yang, Stephan Eckstein, Marcel Nutz, and Stephan Mandt

UCI University of California, Irvine

ETH Zürich

COLUMBIA UNIVERSITY



## Overview:

- We apply ideas and techniques from optimal transport to make advances on a basic problem in information theory — estimating the rate-distortion (R-D) function from data.
- Our R-D estimator is based on minimizing an appropriate functional in the space of probability measures, approximated by moving particles.
- We draw close connections between R-D estimation, entropic optimal transport, and deconvolution, and leverage the connections to:
  - introduce a new class of sources with known solutions to the R-D problem as test cases for algorithms, and
  - derive sample complexity for our R-D / deconvolution estimator.
- We obtain comparable or improved R-D estimates compared to SOTA methods based on neural networks [Yang & Mandt 2022; Lei et al. 2023], while requiring significantly less computation and tuning.

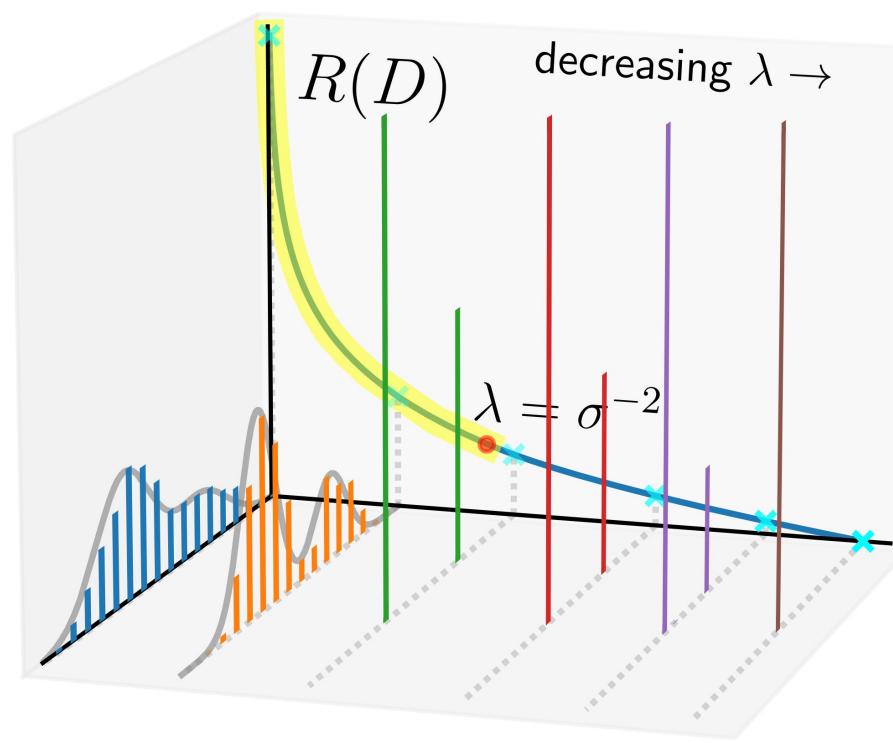
## Background: lossy compression and $R(D)$

In lossy compression, we are given

- The spaces (“alphabets”) of data and reproductions,  $(\mathcal{X}, \mathcal{Y})$ .
- The source (data) distribution  $\mu$  (a prob. measure) on  $\mathcal{X}$ .
- A distortion function  $\rho: \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty)$

A lossy compression algorithm maps the source measure  $\mu$  on  $\mathcal{X}$  to a reproduction measure  $\nu$  on  $\mathcal{Y}$ , incurring

- a **distortion**/transportation cost (“reconstruction error”) and
- a **rate** cost (“avg. file size”).



**Q:** what is the best possible rate-distortion trade-off?  
**A:** the rate-distortion function  $R(D)$ .

$$R(D) := \inf_{\pi \in \Pi(\mu, \cdot): \int \rho d\pi \leq D} H(\pi | \pi_1 \otimes \pi_2)$$

Following [Blahut 1972, Arimoto 1972], we work with an equivalent variational “Lagrangian” representation of  $R(D)$ :

$$F(\lambda) := \inf_{\nu \in \mathcal{P}(\mathcal{Y})} \underbrace{\inf_{\pi \in \Pi(\mu, \cdot)} \lambda \int \rho d\pi + H(\pi | \mu \otimes \nu)}_{\mathcal{L}_{BA}^\lambda(\mu, \nu)}$$

## Connections to EOT and denoising

We show that the Lagrangian R-D problem (1) is equivalent to:

- (2) projecting the source measure under an entropic optimal transport (EOT) cost;
- (3) denoising/deconvolving the source by maximum-likelihood.

$$\begin{aligned} (1) \min_{\nu \in \mathcal{P}(\mathcal{Y})} \mathcal{L}_{BA}^\lambda(\mu, \nu) &\longleftrightarrow (2) \min_{\nu \in \mathcal{P}(\mathcal{Y})} \mathcal{L}_{EOT}^{1/\lambda}(\mu, \nu) \\ &\quad \mathcal{L}_{EOT}^\epsilon(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int \rho d\pi + \epsilon H(\pi | \mu \otimes \nu) \\ &\quad \text{[Rigollet and Weed, 2018]} \\ (3) \max_{\nu \in \mathcal{P}(\mathcal{Y})} \mathbb{E}_{x \sim \mu} [\log \left( \int e^{-\lambda \rho(x, y)} \nu(dy) \right)] &\quad \begin{array}{c} \text{Y} \\ \downarrow \\ \text{X} \end{array} \quad \begin{array}{c} Y \sim \nu^* \\ X|Y=y \sim \mathcal{N}(y, \frac{1}{\lambda}) \end{array} \end{aligned}$$

Thus our algorithm/results transfer across all three problems.

In particular:

- (1)  $\leftrightarrow$  (2): we leverage sample complexity results for EOT [Mena and Niles-Weed, 2019] to obtain finite-sample bounds for R-D estimation / projection under EOT / maximum-likelihood deconvolution;
- (1)  $\leftrightarrow$  (3): we leverage the solution to the deconvolution problem to introduce a new class of sources with closed-form  $R(D)$  segments.

## Proposed: Wasserstein Gradient Descent (WGD)

Let  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ ,  $\rho$  continuously differentiable. We aim to solve

$$\min_{\nu \in \mathcal{P}(\mathbb{R}^d)} \mathcal{L}(\nu), \quad \mathcal{L}(\cdot) \in \{\mathcal{L}_{BA}(\mu, \cdot), \mathcal{L}_{EOT}(\mu, \cdot)\}$$

Idea: simulate the gradient flow of  $\mathcal{L}$  in the 2-Wasserstein space of probability measures:

$$\nu^{(t)} = \left( \text{id} - \underbrace{\gamma \nabla \frac{\delta \mathcal{L}}{\delta \nu}}_{\text{Wasserstein gradient: } \mathbb{R}^d \rightarrow \mathbb{R}^d}(\nu^{(t-1)}) \right) \# \nu^{(t-1)}$$

Particle scheme in practice:

$$\nu^{(t)} = \frac{1}{n} \sum_{i=1}^n \delta_{y_i^{(t)}} \quad y_i^{(t)} = y_i^{(t-1)} - \gamma \nabla \frac{\delta \mathcal{L}}{\delta \nu}(\nu^{(t-1)})[y_i^{(t-1)}], \quad \forall i$$

The Wasserstein gradient can be tractably computed by

- Sinkhorn’s algorithm, for  $\mathcal{L} = \mathcal{L}_{EOT}$ , or
- A **single** Sinkhorn iteration, for  $\mathcal{L} = \mathcal{L}_{BA}$  (orders of magnitude faster!)

## Finite-sample bounds for R-D estimation

Given  $m$  samples from a  $\sigma^2$ -sub-Gaussian source, the best empirical loss achievable with  $n$  particles converges to the true optimum as follows ( $\epsilon = 1/\lambda$ ):

$$\mathbb{E} \left[ \left| \min_{\nu \in \mathcal{P}(\mathbb{R}^d)} \mathcal{L}(\mu, \nu) - \min_{\nu_n \in \mathcal{P}_n(\mathbb{R}^d)} \mathcal{L}(\mu^m, \nu_n) \right| \right] \leq C_d \epsilon \left( 1 + \frac{\sigma^{[5d/2]+6}}{\epsilon^{[5d/4]+3}} \right) \left( \frac{1}{\sqrt{m}} + \frac{1}{\sqrt{n}} \right)$$

## Empirical results

- We compare against other R-D upper bound algorithms: Blahut-Arimoto [Blahut 1972; Arimoto 1972] and SOTA deep learning methods RD-VAE [Yang & Mandt 2022] and NERD [Lei et al. 2023].
- For a given per-iteration compute budget, we obtain much faster convergence and better approximation quality (deconv example):

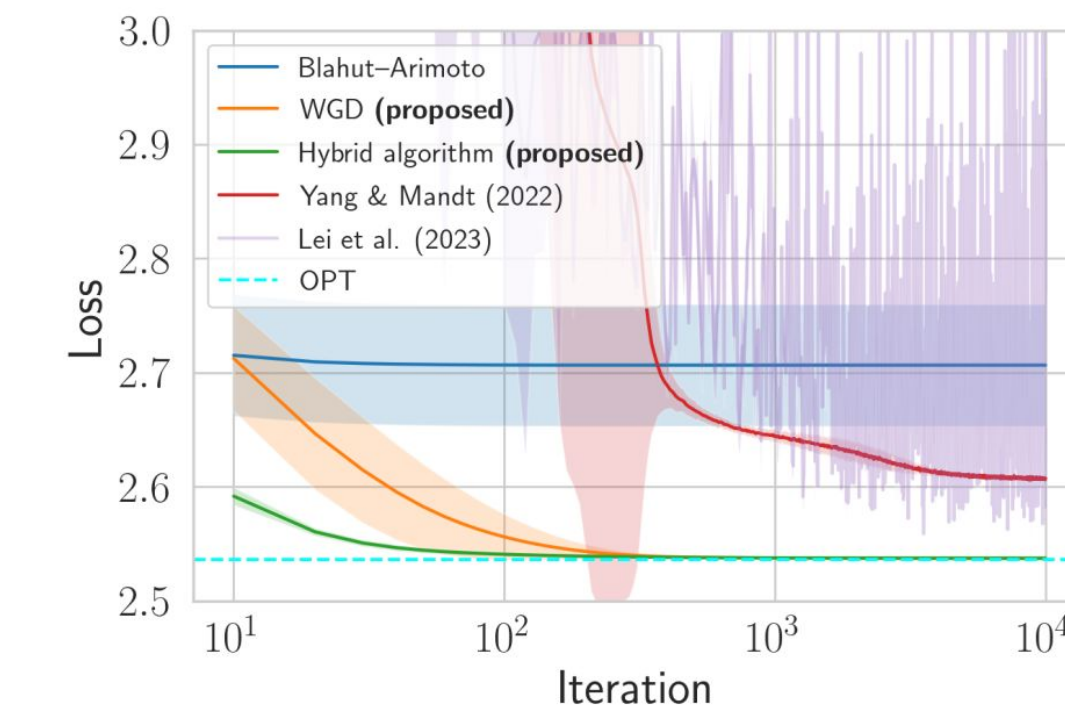


Figure 2: Losses over iterations. Shading corresponds to one standard deviation over random initializations.

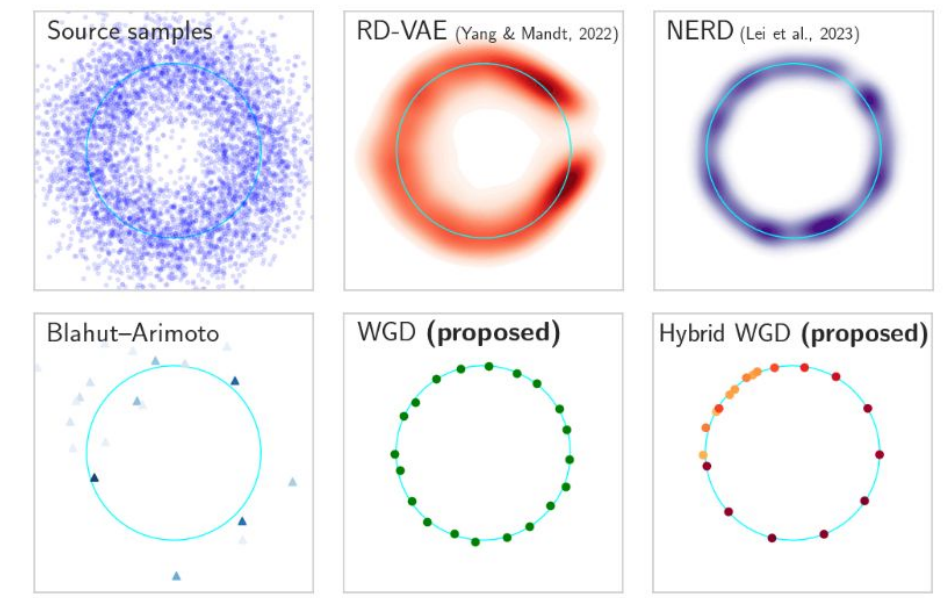
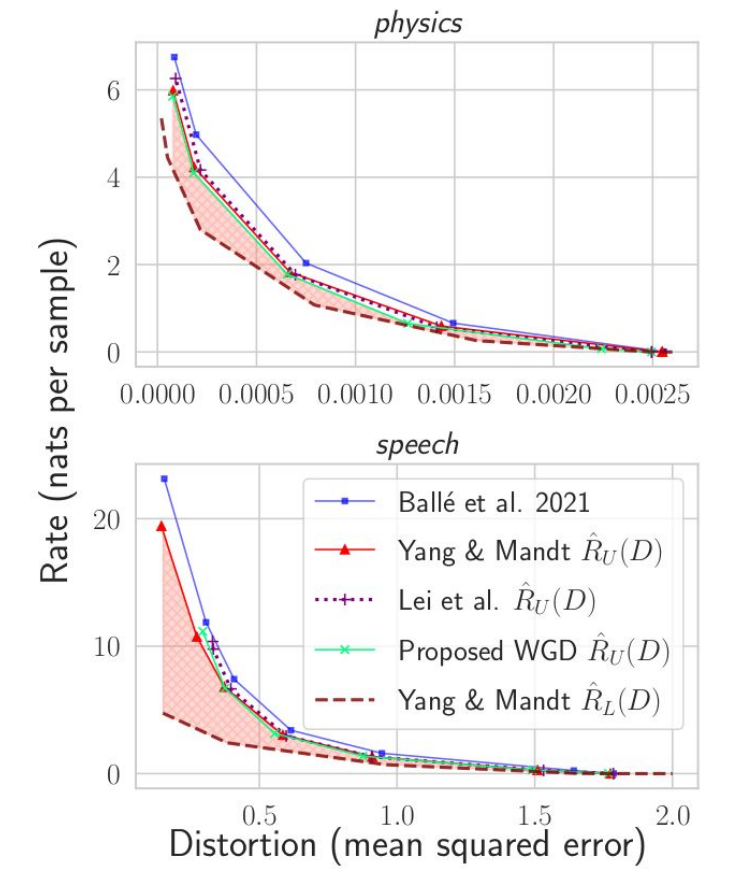
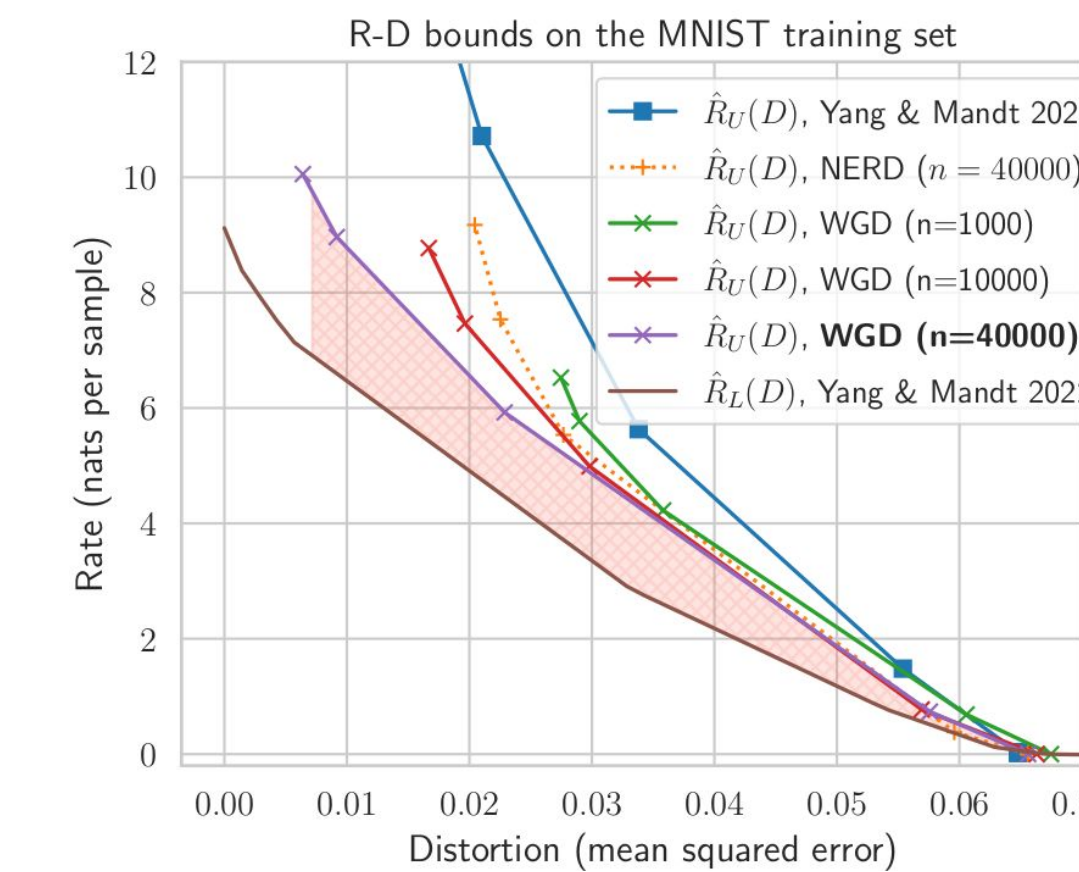


Figure 3: Visualizing  $\mu$  samples (top left), as well as the  $\nu$  returned by various algorithms compared to the ground truth  $\nu^*$  (cyan).

- as well as tighter R-D upper bounds:



## Limitations and future work

- Like NERD [Lei et al. 2023], our method can only target an  $R(D)$  point with a rate  $\leq \log(n)$  nats/sample, where  $n$  = number of particles.
- It remains to be studied how best to convert our R-D estimator into a practical data communication/compression algorithm.

## References

- [Blahut 1972] R. Blahut. “Computation of channel capacity and rate-distortion functions”. *IEEE Transactions on Information Theory*, 18(4):460–473, 1972.
- [Arimoto 1972] S. Arimoto. “An algorithm for computing the capacity of arbitrary discrete memoryless channels”. *IEEE Transactions on Information Theory*, 18(1):14–20, 1972.
- [Rigollet and Weed, 2018] Philippe Rigollet and Jonathan Weed. Entropic optimal transport is maximum-likelihood deconvolution. *Comptes Rendus Mathématique*, 356(11–12):1228–1235, 2018.
- [Mena and Niles-Weed, 2019] Gonzalo Mena and Jonathan Niles-Weed. Statistical bounds for entropic optimal transport: sample complexity and the central limit theorem. *Advances in Neural Information Processing Systems*, 32, 2019.
- [Yang & Mandt 2022] Y. Yang, S. Mandt. Towards empirical sandwich bounds on the rate-distortion function”. *International Conference on Learning Representations*, 2022.
- [Lei et al. 2023] E. Lei, H. Hassani, and S. Saeedi Bidokhti. “Neural estimation of the rate-distortion function with applications to operational source coding”. *IEEE Journal on Selected Areas in Information Theory*, 2023.