

Are generative models the new sparsity*?

Bruno Loureiro
IdePHICS, EPFL

*from <https://solevillar.github.io/2018/03/28/SUNLayer.html>

Joint work with



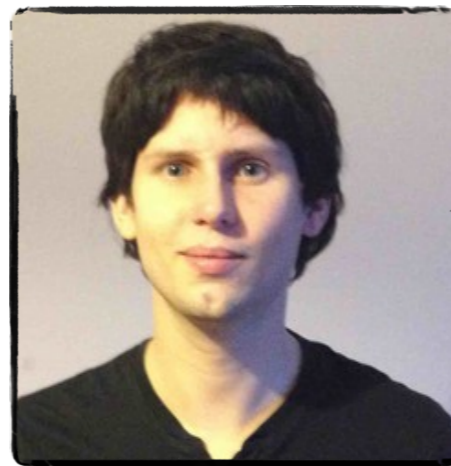
Benjamin Aubin
(ClipDrop)



Lenka Zdeborová
(EPFL)



Antoine Maillard
(ENS Paris)

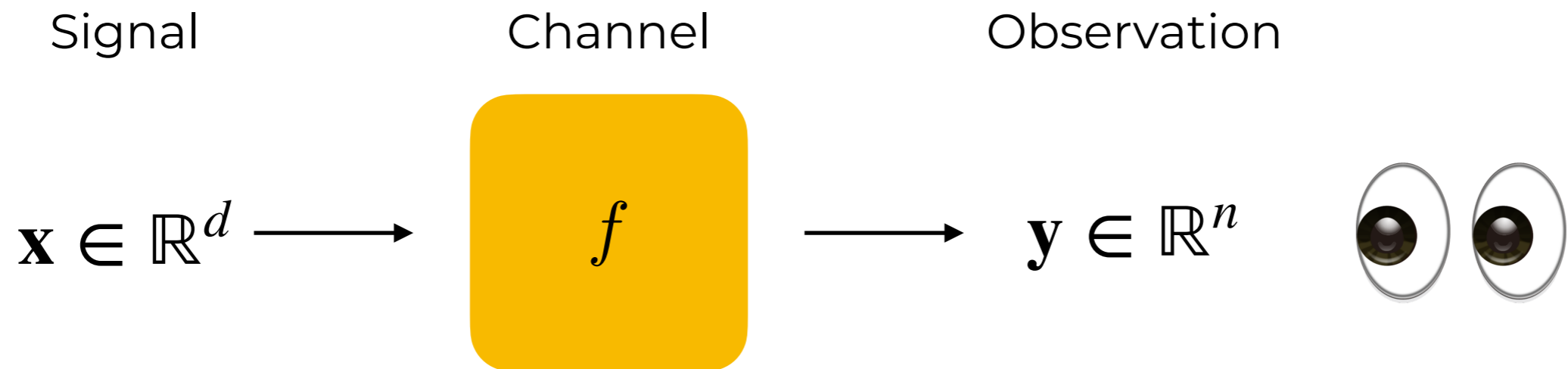


Antoine Baker
(ENS Paris)

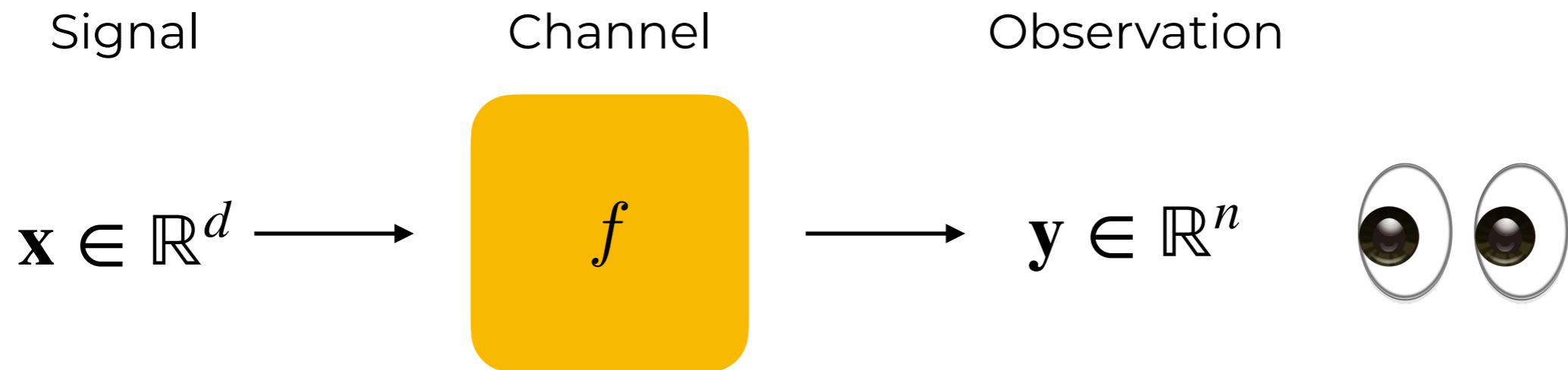


Florent Krzakala
(EPFL)

Inverse problems



Inverse problems



- Questions:**
- How many samples?
 - How to do it efficiently?
 - What is the role of structure?

Examples

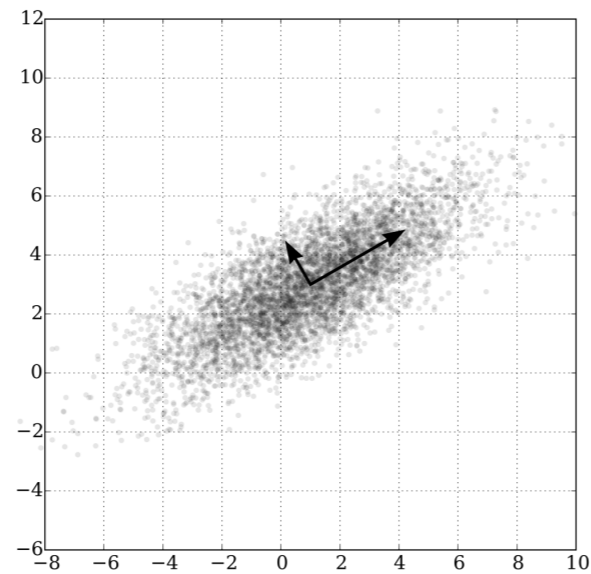
Denoising

$$\mathbf{y} = \mathbf{x} + \xi$$



Matrix factorisation

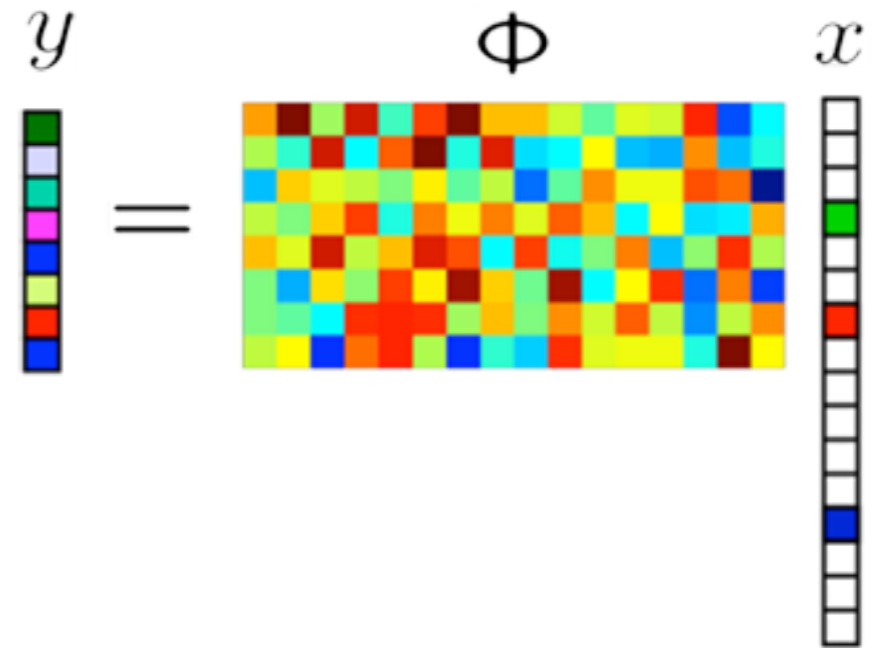
$$\mathbf{Y} = \mathbf{x}\mathbf{x}^T + \mathbf{Z}$$



Compressed sensing
Phase retrieval

$$\mathbf{y} = \Phi\mathbf{x} + \xi$$

$$\mathbf{y} = |\mathbf{A}\mathbf{x}| + \xi$$

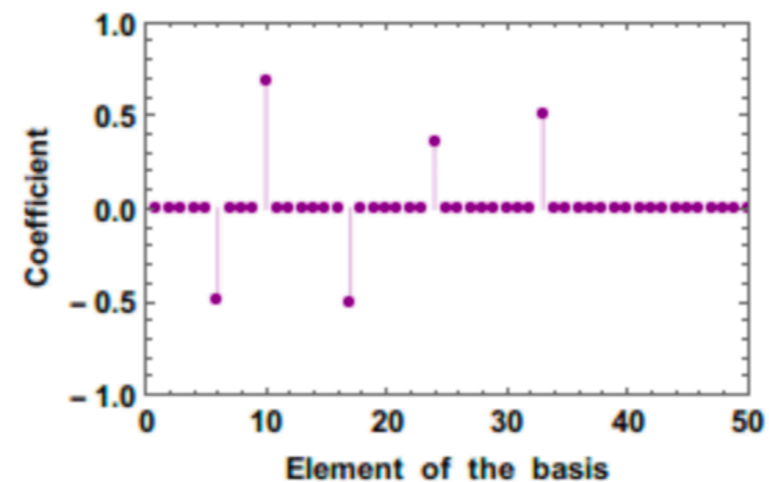
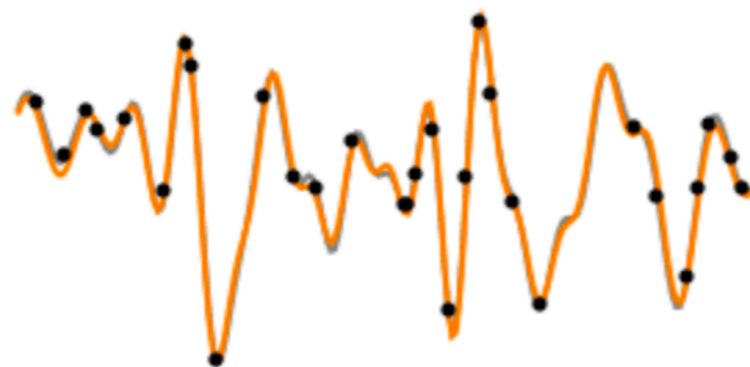


Sparsity

Structure helps!

If there is a basis (possibly learned) for which signal is sparse:

$$\mathbf{x} = (-\mathbf{z} - \underbrace{0 \dots 0}_{d-k})$$

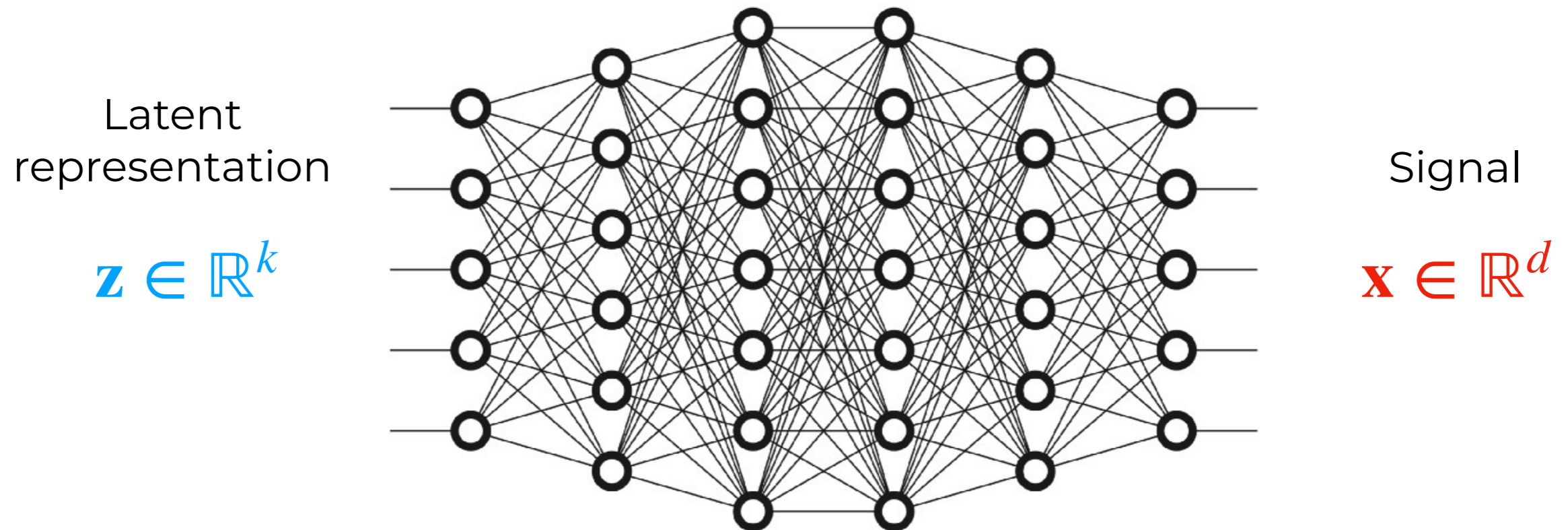


Efficient reconstruction might be possible both **statistically** (NP-hard?) and **algorithmically**, e.g. L1-minimisation in CS.

Learning a basis

[Ulyanov et al. 17,
Bora et al. 17',
Hecklel, Hand 18']

Deep generative networks: VAEs, GANs, etc.



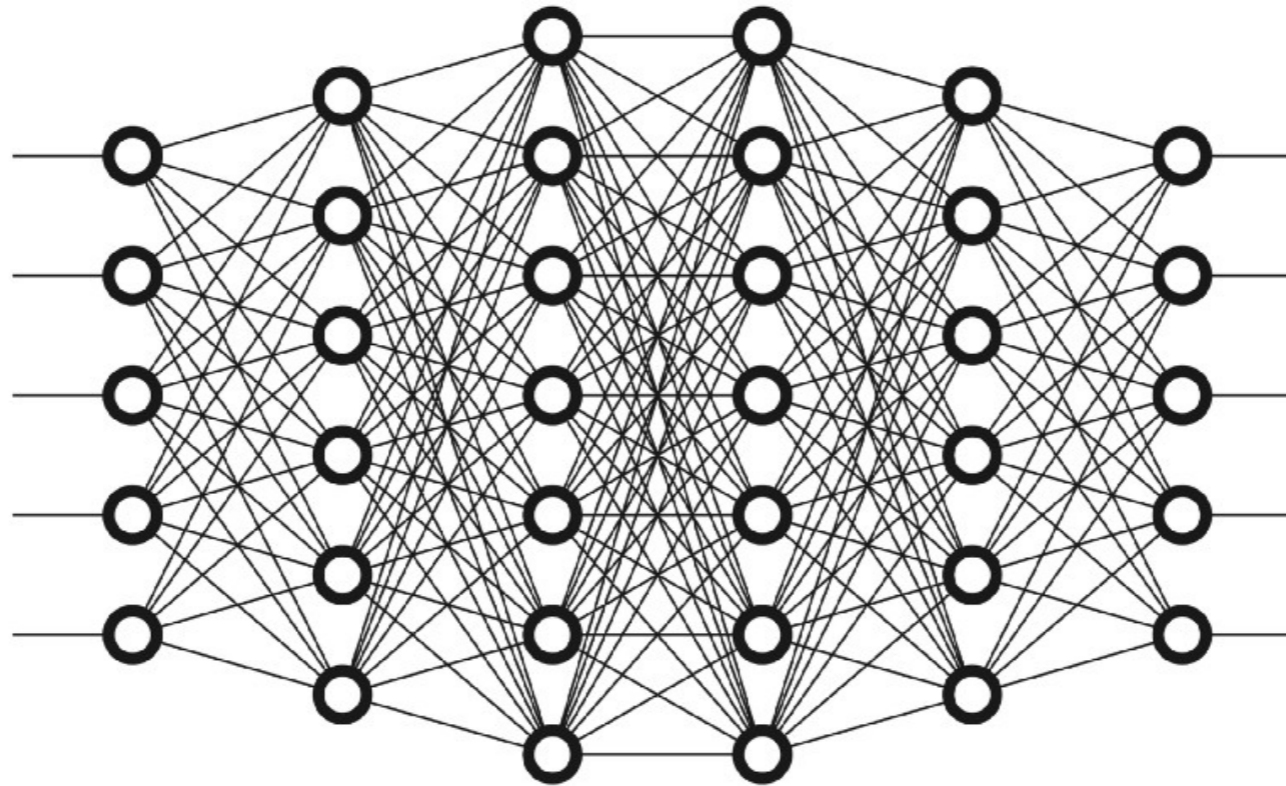
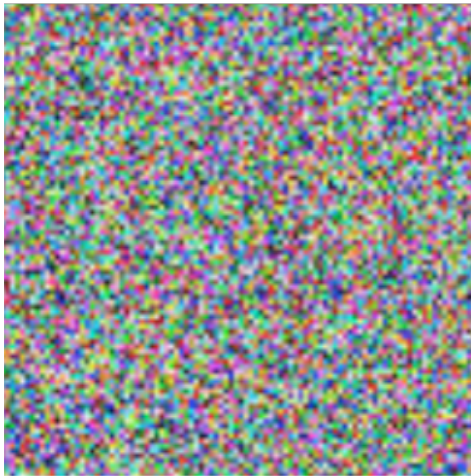
$$\mathbf{x} = \sigma^{(L)} \left(\mathbf{W}^{(L)} \sigma^{(L-1)} \left(\mathbf{W}^{(L-1)} \dots \sigma^{(1)} \left(\mathbf{W}^{(1)} \mathbf{z} \right) \dots \right) \right) \in \mathbb{R}^d$$

Learning a basis

[Ulyanov et al. 17,
Bora et al. 17',
Hecklel, Hand 18']

Deep generative networks: VAEs, **GANs**, etc.

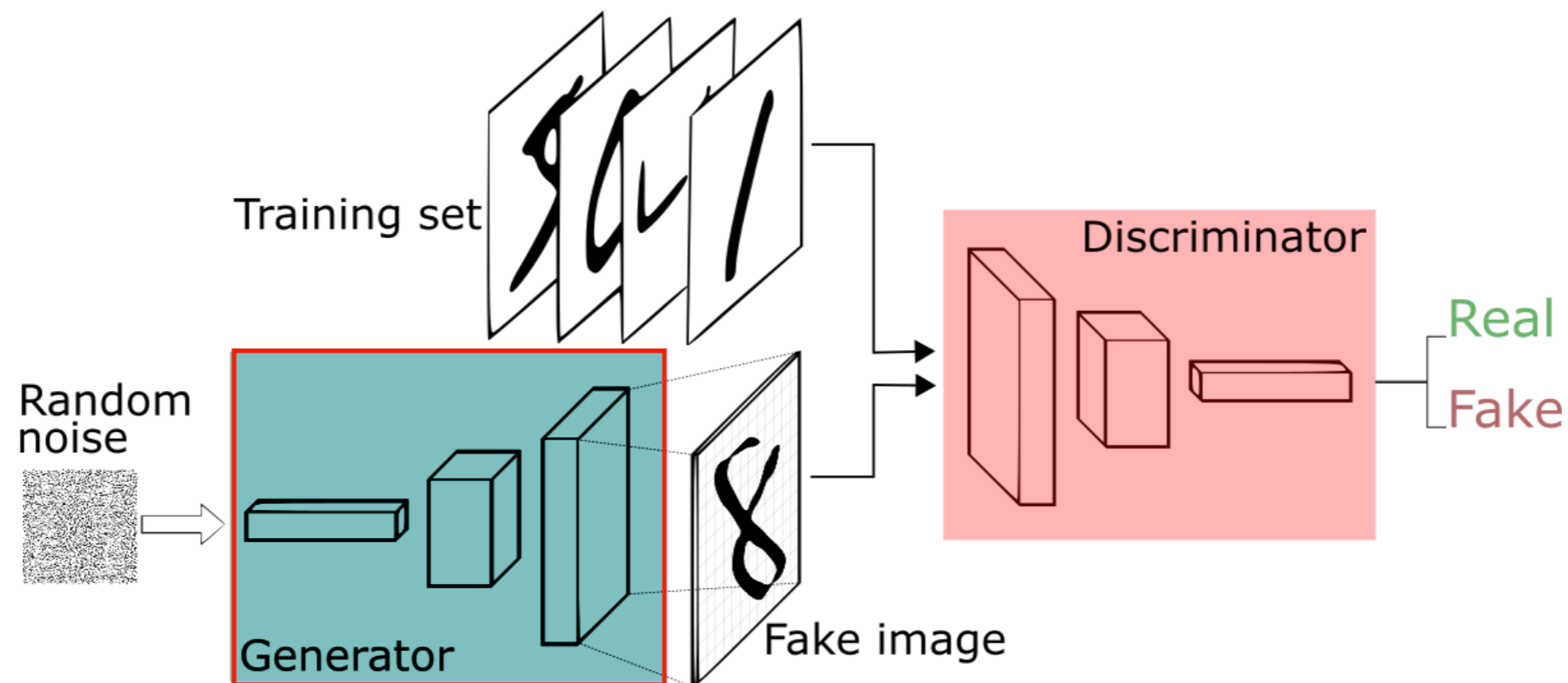
$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$



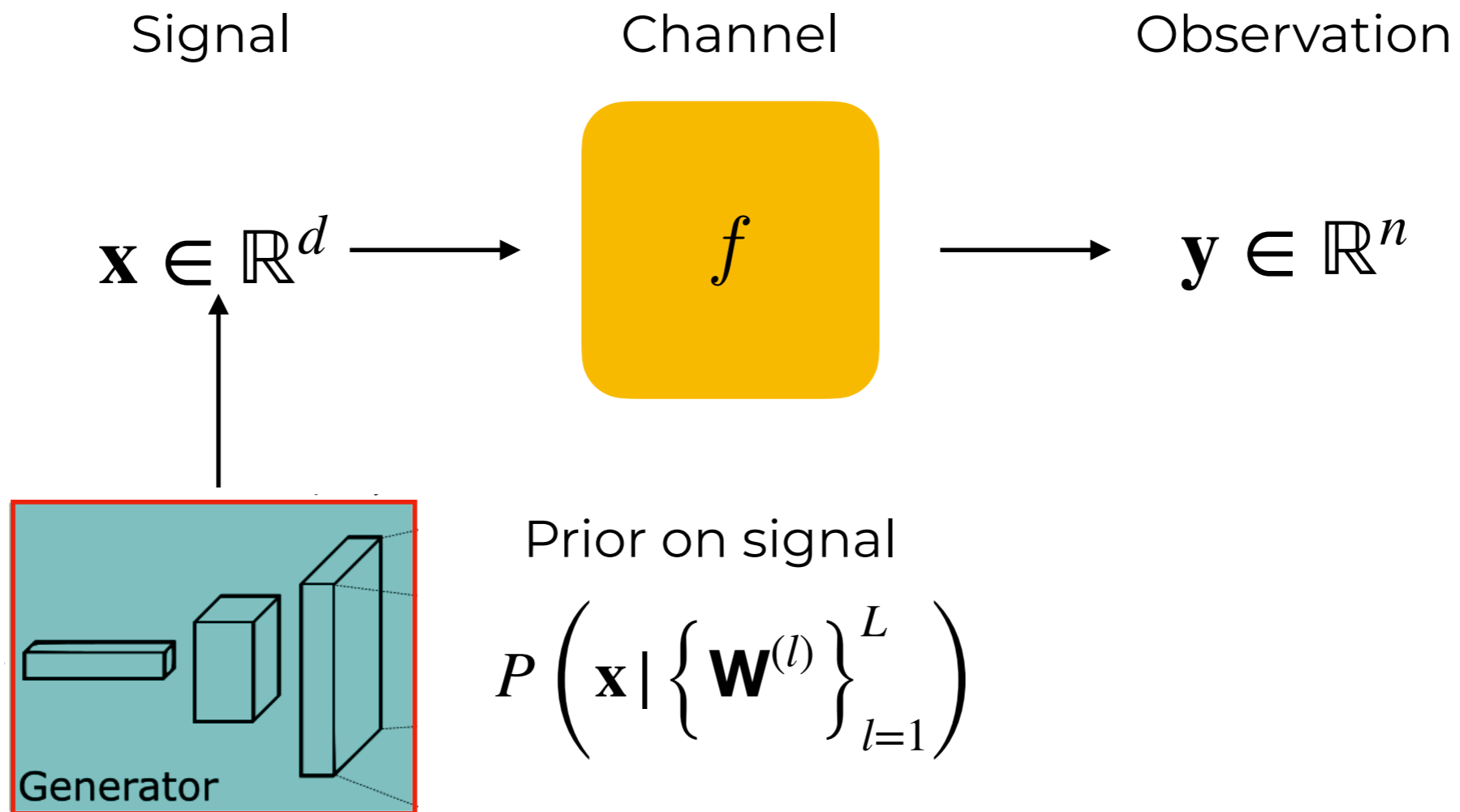
<https://www.thispersondoesnotexist.com/>

Setting

Someone, somewhere



Setting



Matrix factorisation

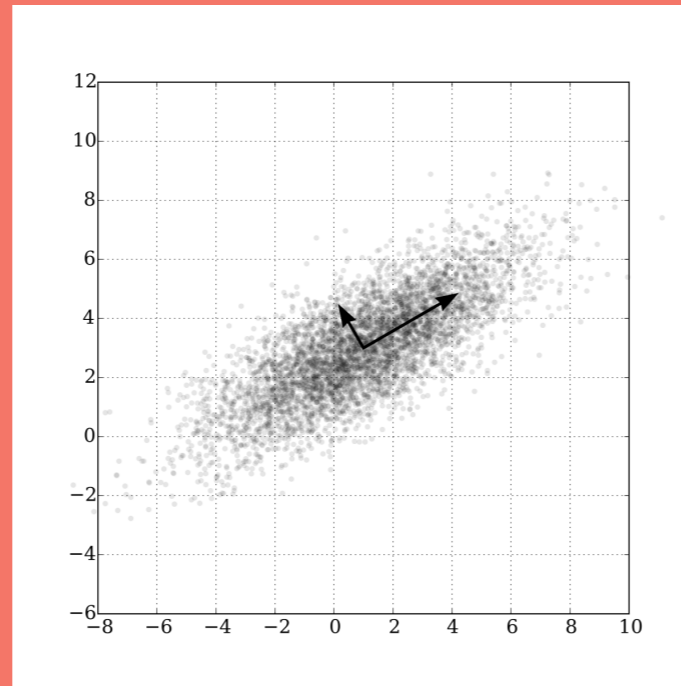
Denoising

$$\mathbf{y} = \mathbf{x} + \xi$$



Matrix factorisation

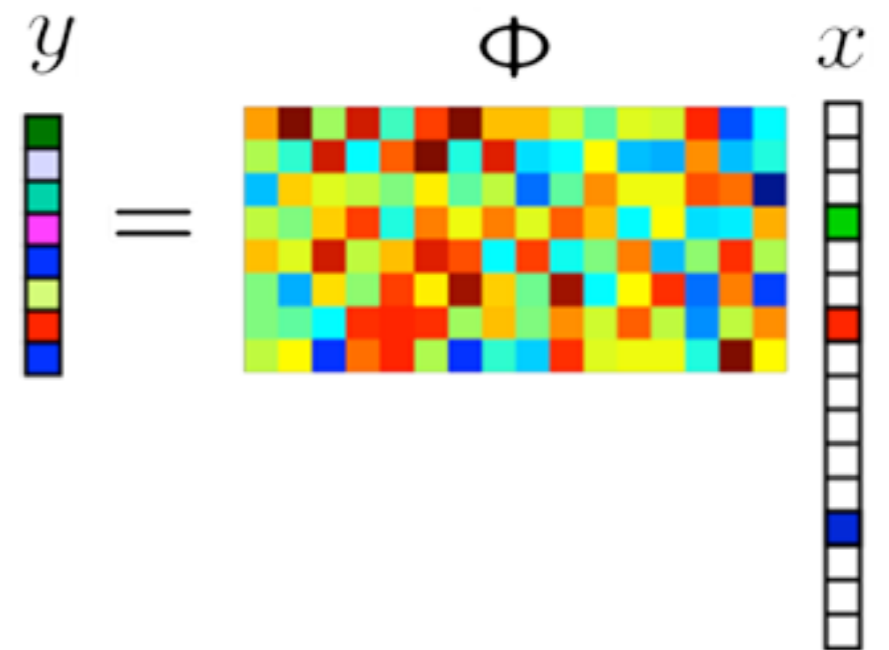
$$\mathbf{Y} = \mathbf{x}\mathbf{x}^T + \mathbf{Z}$$



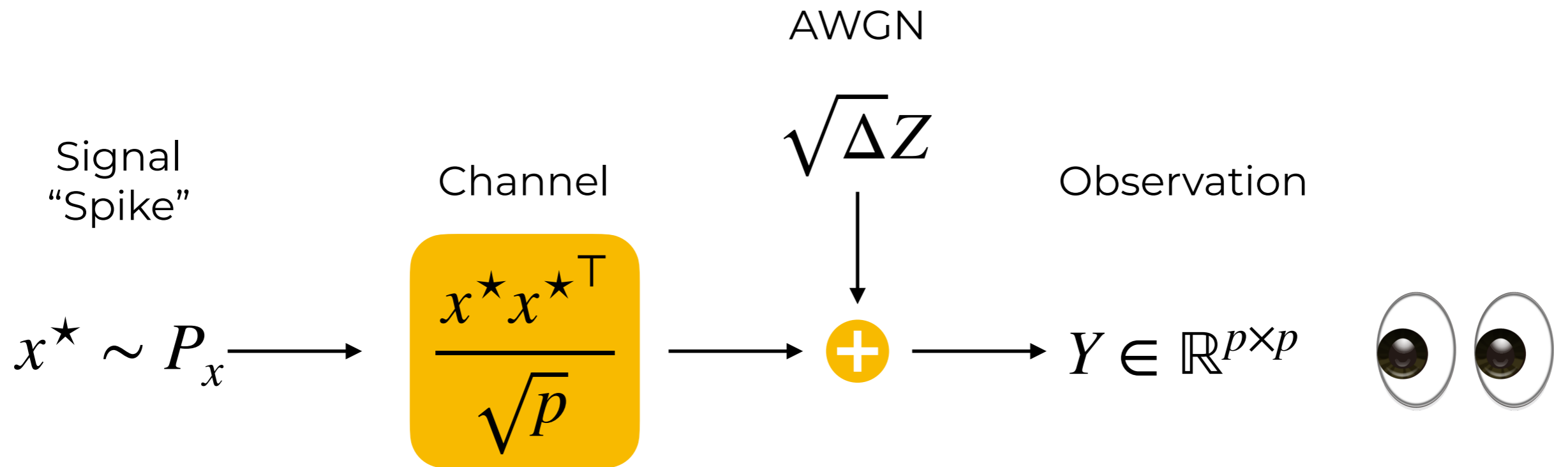
Compressed sensing
Phase retrieval

$$\mathbf{y} = \Phi\mathbf{x} + \xi$$

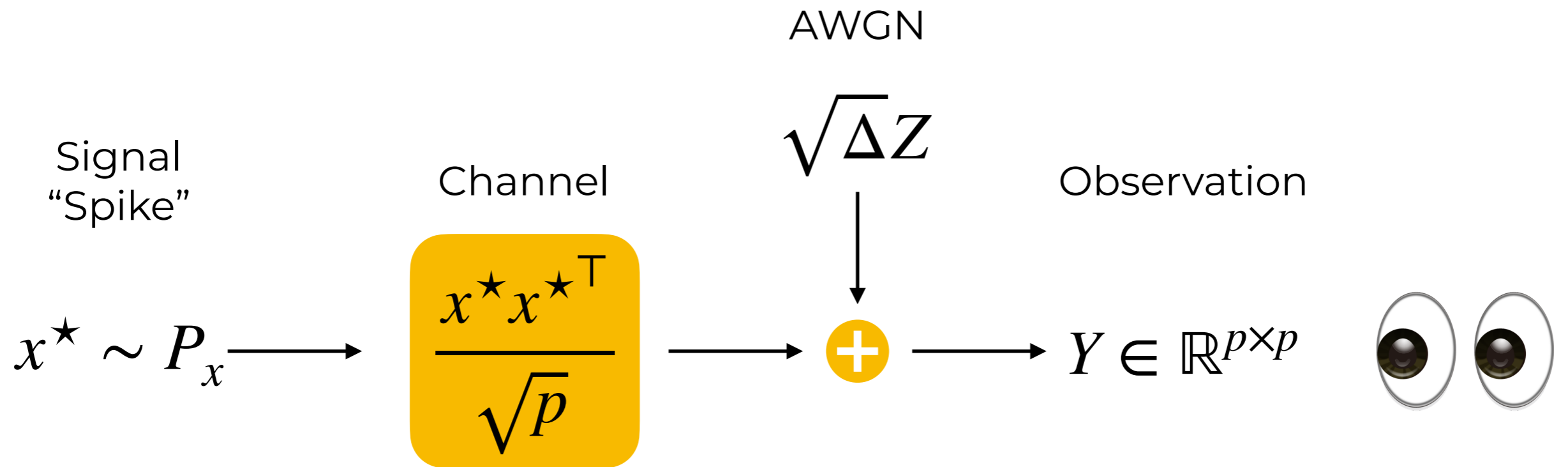
$$\mathbf{y} = |\mathbf{A}\mathbf{x}| + \xi$$



Matrix factorisation



Matrix factorisation

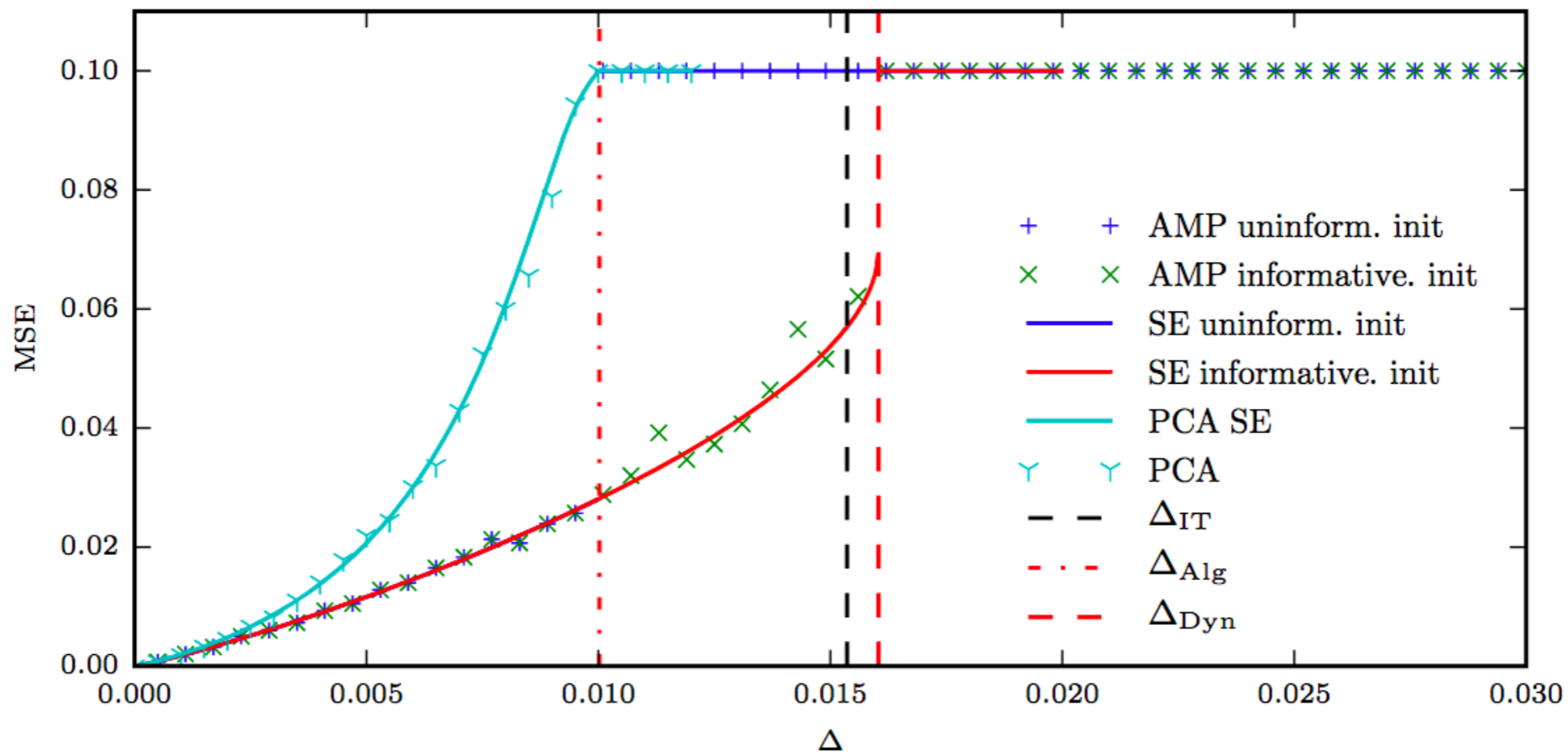


Goal: reconstruct x^\star from knowledge of P_x, Δ
in the high-dimensional regime
 $p \rightarrow \infty$ and $\Delta = O(1)$

Sparse case

Consider a **sparse** spike

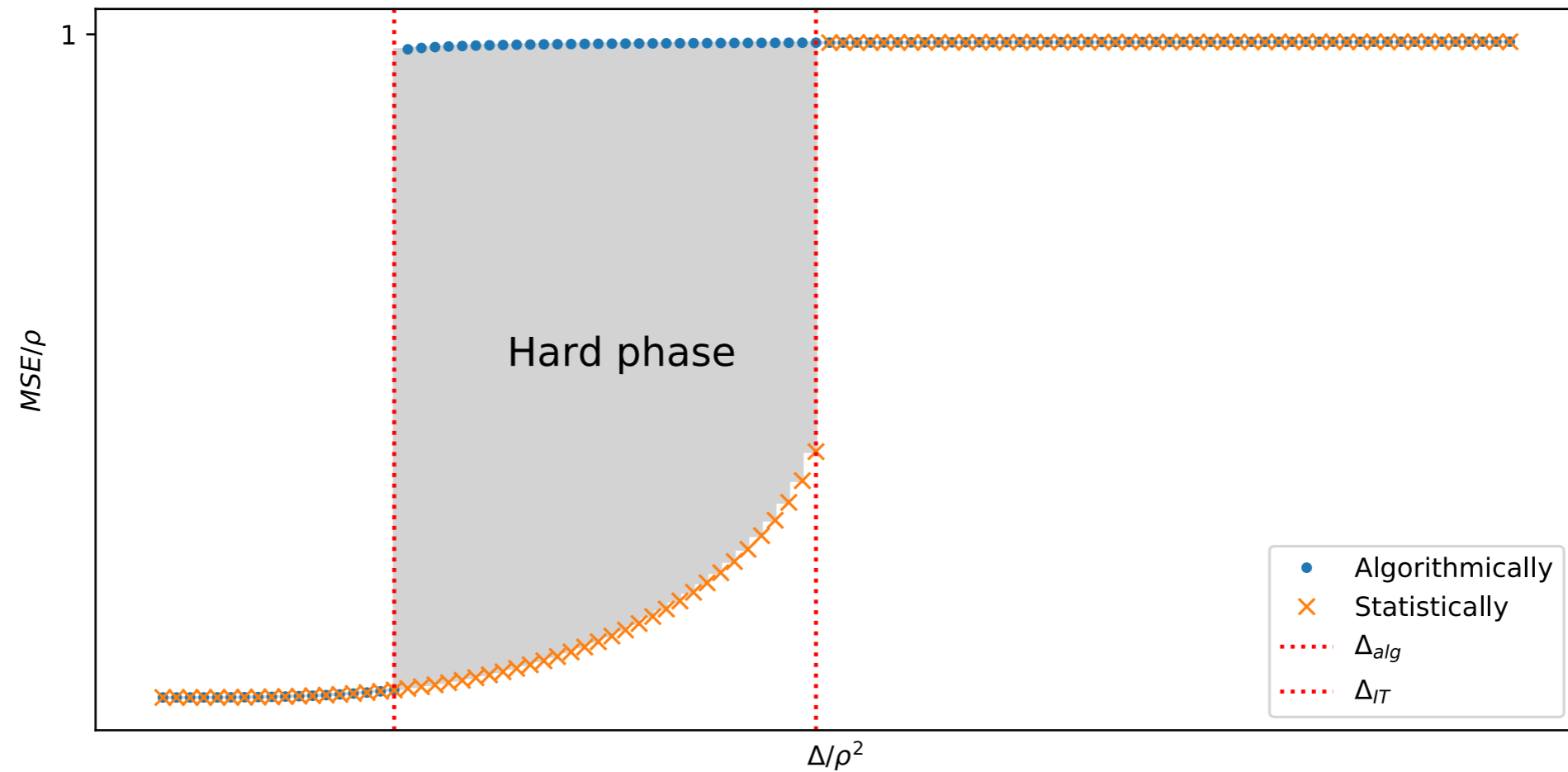
$$x^* \sim \prod_{i=1}^p [(1 - \rho)\delta_0 + \rho\mathcal{N}(0, 1)] \quad \rho = 0.1$$



[Lesieur et al. 17']

Sparse case

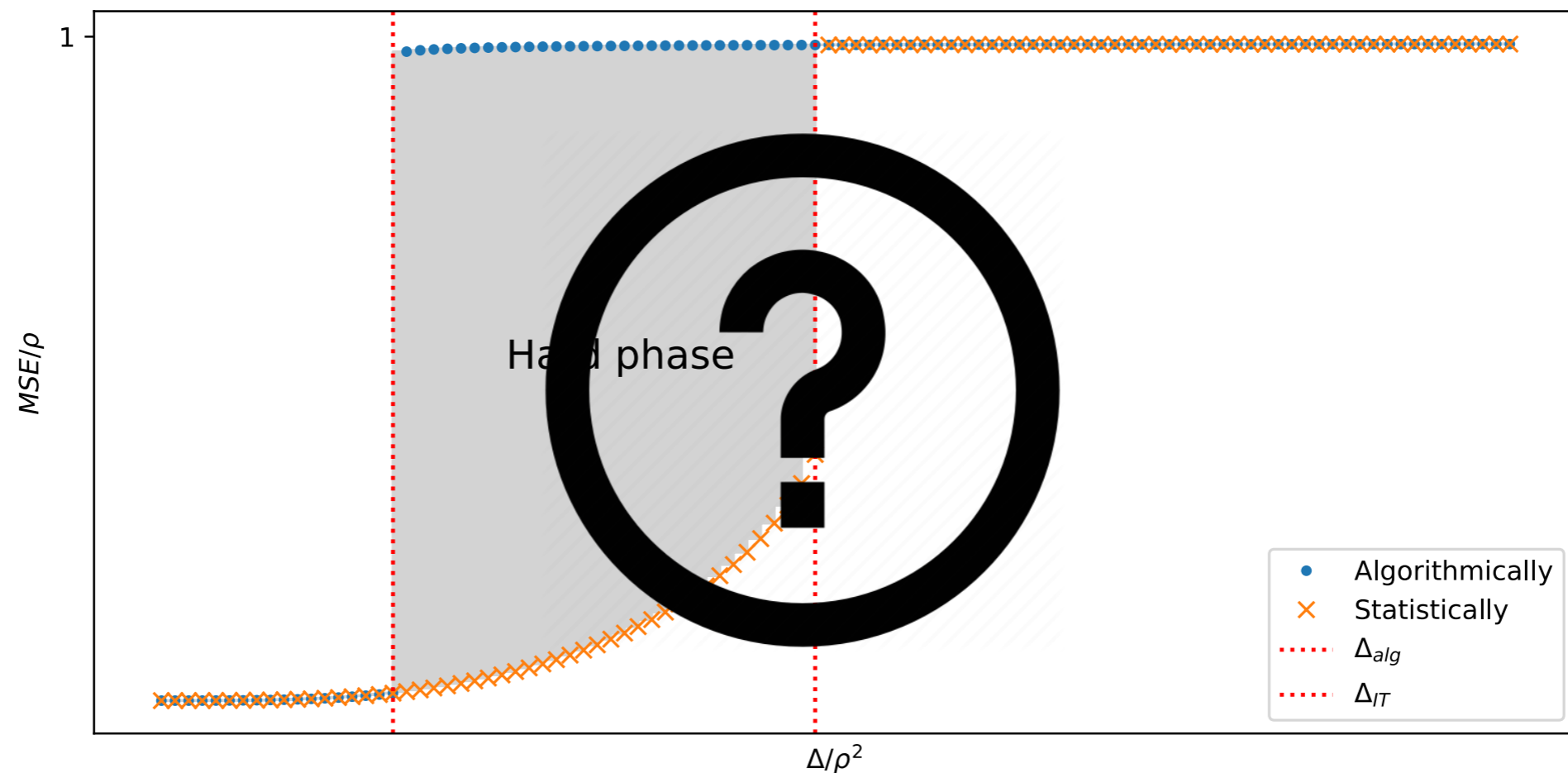
Large sparsity leads to
statistical-to-algorithmic gap



Generative networks

$$\mathbf{x} = \sigma^{(L)} \left(\mathbf{W}^{(L)} \sigma^{(L-1)} \left(\mathbf{W}^{(L-1)} \dots \sigma^{(1)} \left(\mathbf{W}^{(1)} \mathbf{z} \right) \dots \right) \right)$$

$$\mathbf{W}^{(l)} \in \mathbb{R}^{k_{l+1} \times k_l} \text{ i.i.d. Gaussian}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$$



Statistical guarantee

Theorem
(informal)

In the high-dimensional limit,

$$\text{MMSE} = \rho_x - q_x^*$$

where: $q_x^* = \operatorname{arginf}_{0 \leq q_x \leq \rho} i_{\text{RS}}(\Delta, q_x)$

$$i_{\text{RS}}(\Delta, q_x) = \frac{(\rho_x - q_x)^2}{4\Delta} + \lim_{p \rightarrow \infty} \frac{I\left(x; x + \sqrt{\frac{\Delta}{q_x}} \xi\right)}{p}$$

$$\rho_x = \lim_{p \rightarrow \infty} \mathbb{E}_{P_x} \left[\frac{x \cdot x}{p} \right], \quad \xi \sim \mathcal{N}(0, 1)$$

Statistical guarantee

Remark:

$$\frac{1}{p} I \left(x; x + \sqrt{\frac{\Delta}{q_x}} \xi \right)$$

Mutual information density for a denoising problem.

Has been **rigorously** computed for

$$P_x(x) = \prod_{i=1}^p P_x(x_i)$$

Uncorrelated
signal

[Lesieur et al. 15-17']

$$x = \varphi^{(2)} \left(W^2 \varphi^{(1)} \left(W^{(1)} z \right) \right)$$

2-layer NN with
random weights

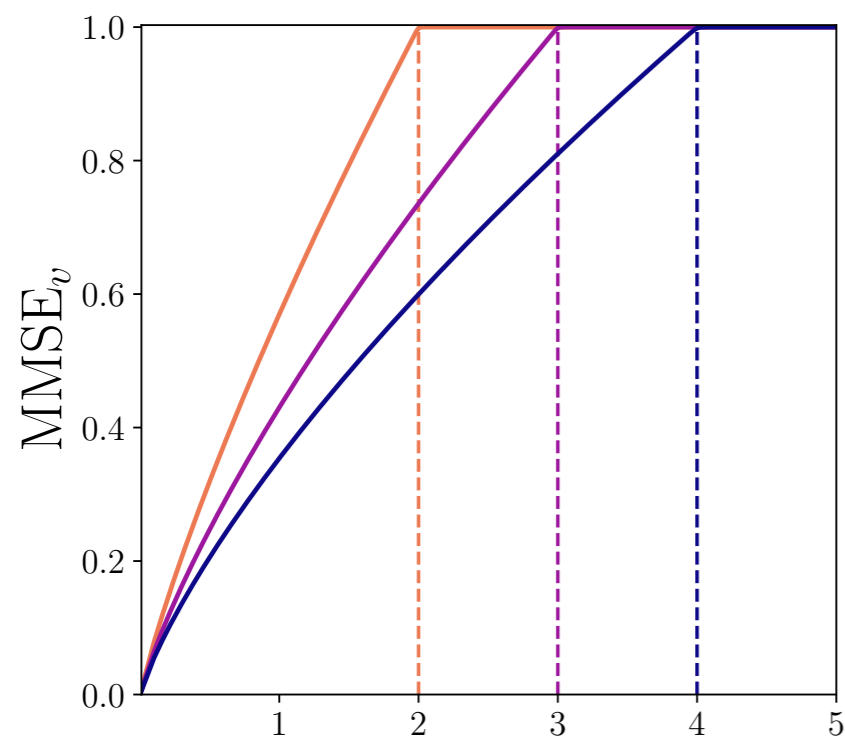
[Krzakala et al. 16']

[Barbier et al. 16']

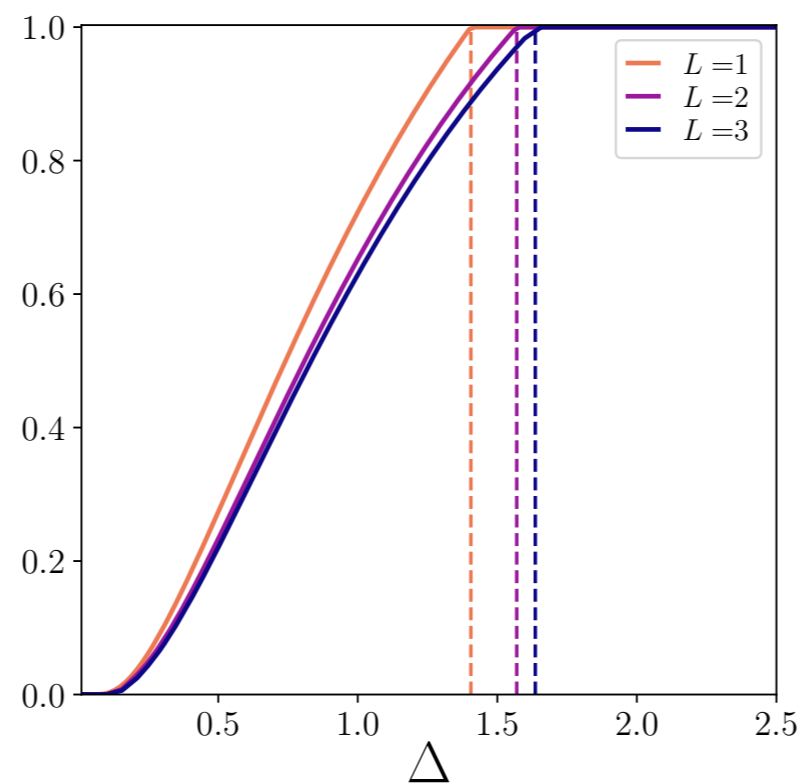
Role of depth

$$\mathbf{x}^* = \sigma^{(L)} \left(\mathbf{W}^{(L)} \sigma^{(L-1)} \left(\mathbf{W}^{(L-1)} \dots \sigma^{(1)} \left(\mathbf{W}^{(1)} \mathbf{z} \right) \dots \right) \right)$$

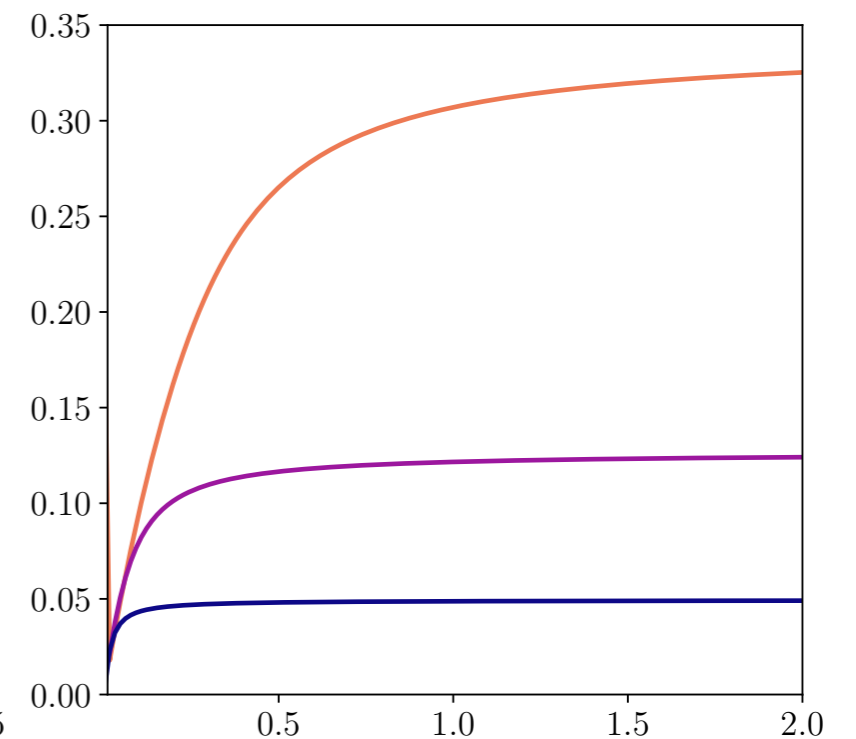
$\sigma(x) = x$



$\sigma(x) = \text{sign}(x)$

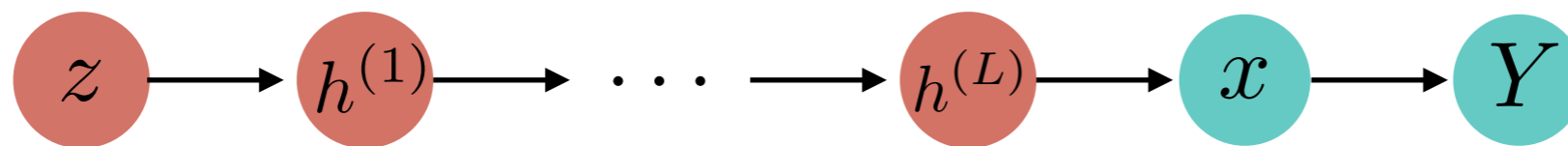


$\sigma(x) = \text{relu}(x)$



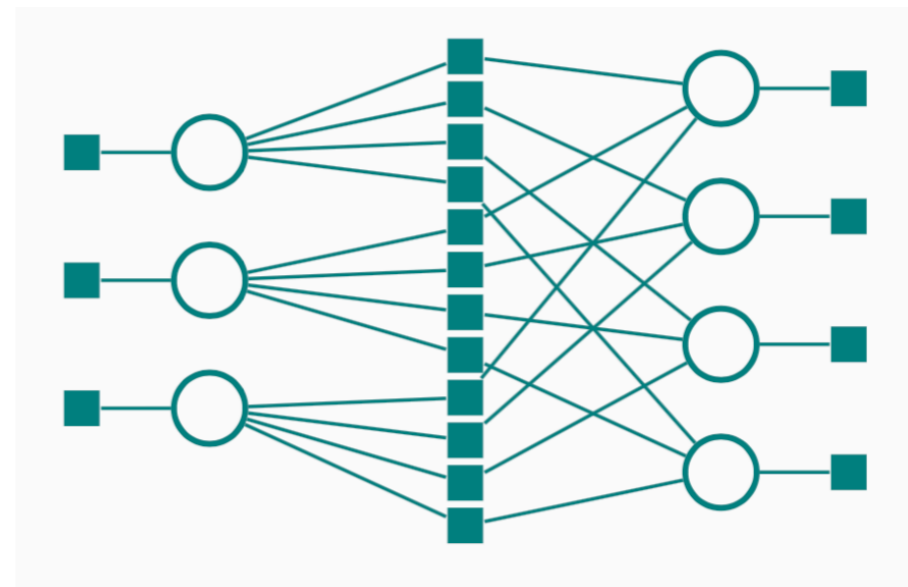
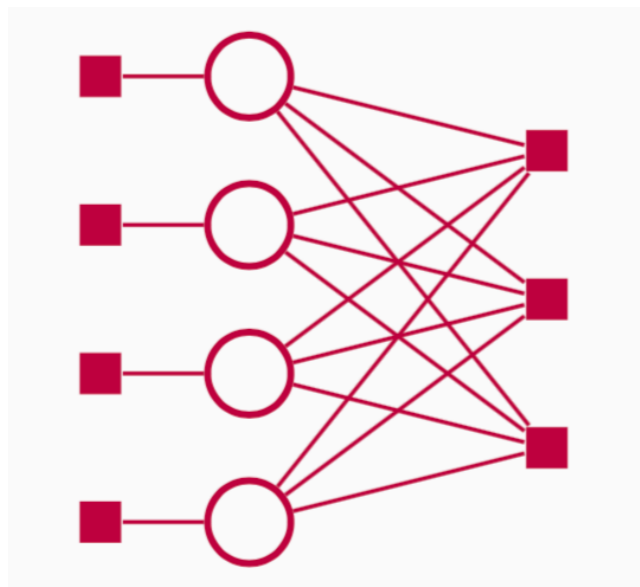
Compositional Principle

A compositional principle for AMP



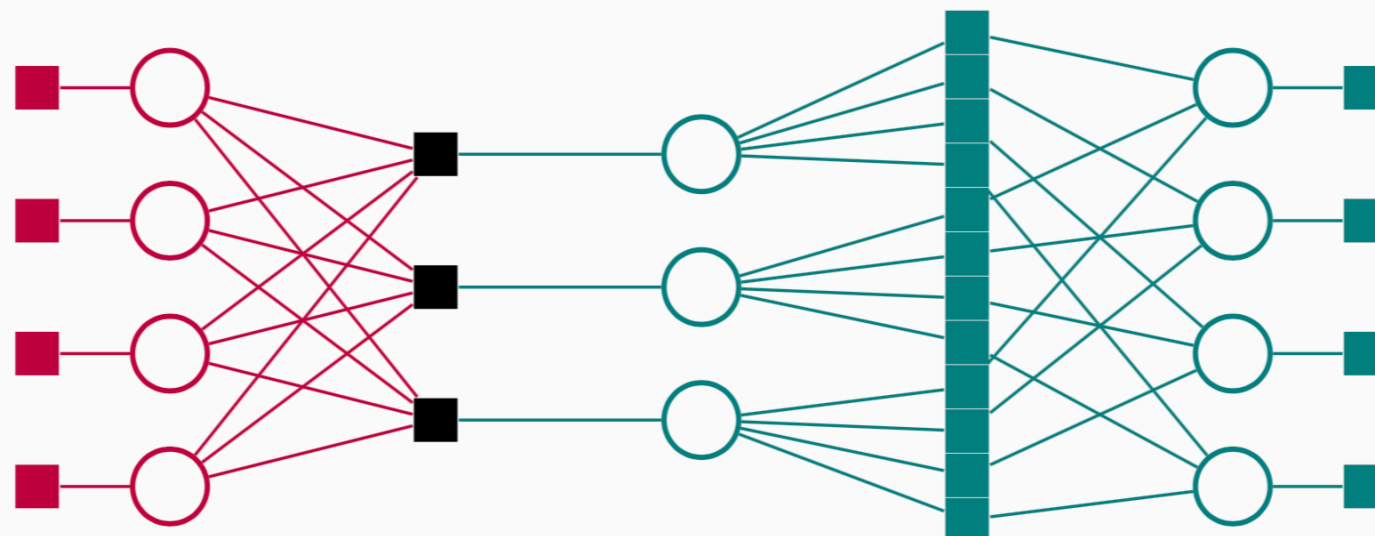
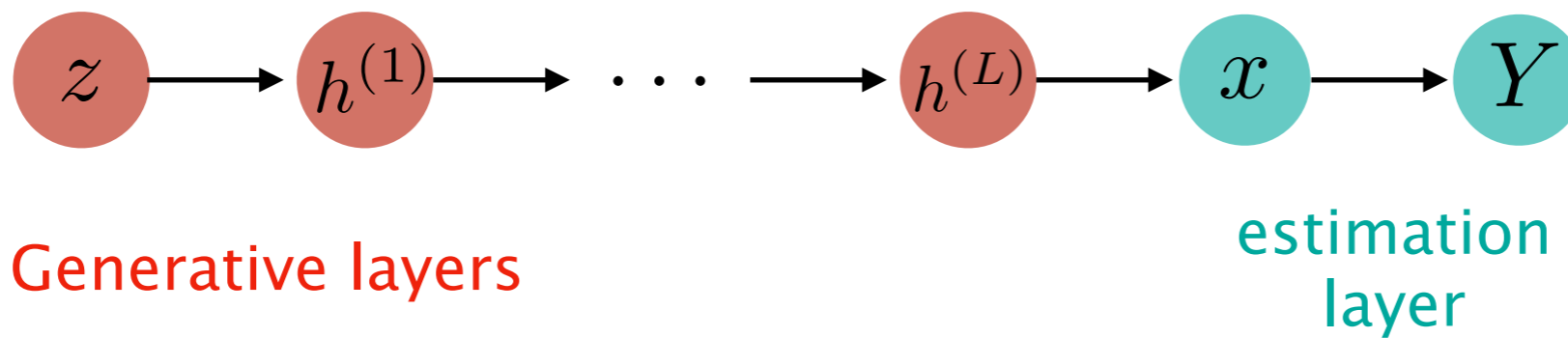
Generative layers

estimation layer



Compositional Principle

A compositional principle for AMP



AMP for one-layer prior

- 1: **Input:** $Y \in \mathbb{R}^{p \times p}$ and $W \in \mathbb{R}^{p \times k}$.
- 2: *Initialize with:* $\hat{\mathbf{v}}^{t=1} = \mathcal{N}(\mathbf{0}, \sigma^2 I_p)$, $\hat{\mathbf{z}}^{t=1} = \mathcal{N}(\mathbf{0}, \sigma^2 I_k)$, and $\hat{\mathbf{c}}_v^{t=1} = I_p$, $\hat{\mathbf{c}}_z^{t=1} = I_k$, $t = 1$.
- 3: **repeat**
- 4: **Spiked layer denoising:**
- 5: $\mathbf{B}_v^t = \frac{1}{\Delta} \frac{Y}{\sqrt{p}} \hat{\mathbf{v}}^t - \frac{1}{\Delta} \frac{(I_p^\top \hat{\mathbf{c}}_v^t)}{p} \hat{\mathbf{v}}^{t-1}$ and $A_v^t = \frac{1}{\Delta p} (\|\hat{\mathbf{v}}^t\|_2)^2 I_p$.
- 6: **Generative layer denoising:**
- 7: $V^t = \frac{1}{k} (I_k^\top \hat{\mathbf{c}}_z^t) I_p$, $\boldsymbol{\omega}^t = \frac{1}{\sqrt{k}} W \hat{\mathbf{z}}^t - V^t \mathbf{g}^{t-1}$ and $\mathbf{g}^t = f_{\text{out}}(\mathbf{B}_v^t, A_v^t, \boldsymbol{\omega}^t, V^t)$,
- 8: $\Lambda^t = \frac{1}{k} \|\mathbf{g}^t\|_2^2 I_k$ and $\gamma^t = \frac{1}{\sqrt{k}} W^\top \mathbf{g}^t + \Lambda^t \hat{\mathbf{z}}^t$.
- 9: **Marginals estimation:**
- 10: $\hat{\mathbf{v}}^{t+1} = f_v(\mathbf{B}_v^t, A_v^t, \boldsymbol{\omega}^t, V^t)$ and $\hat{\mathbf{c}}_v^{t+1} = \partial_{\mathbf{B}} f_v(\mathbf{B}_v^t, A_v^t, \boldsymbol{\omega}^t, V^t)$,
- 11: $\hat{\mathbf{z}}^{t+1} = f_z(\gamma^t, \Lambda^t)$ and $\hat{\mathbf{c}}_z^{t+1} = \partial_\gamma f_z(\gamma^t, \Lambda^t)$,
- 12: $t = t + 1$.
- 13: **until** Convergence.
- 14: **Output:** $\hat{\mathbf{v}}, \hat{\mathbf{z}}$.

TRAMP

Compositional Inference with Tree Approximate Message Passing

Antoine Baker

Florent Krzakala

Laboratoire de Physique

CNRS, École Normale Supérieure, PSL University

Paris, France

Benjamin Aubin

Lenka Zdeborová

Institut de Physique Théorique

CNRS, CEA, Université Paris-Saclay

Saclay, France

ANTOINE.BAKER@ENS.FR

FLORENT.KRZAKALA@ENS.FR

BENJAMIN.AUBIN@CEA.FR

LENKA.ZDEBOROVA@CEA.FR

Abstract

We introduce *tramp*, standing for *TRee Approximate Message Passing*, a python package for compositional inference that runs Expectation Propagation on high-dimensional tree-structured models. The package provides an unifying framework to study several approximate message passing algorithms previously derived for a variety of machine learning tasks such as generalized linear models, inference in multi-layer networks, matrix factorization, and reconstruction using non-separable penalties. For some models, the performance of the algorithm can be theoretically predicted by the State Evolution, and the measurements entropy estimated by the free entropy formalism. The implementation is modular by design: each module, which implements a factor, can be composed at will with other modules to solve complex inference tasks. The user only needs to declare the factor graph of the model: the Expectation Propagation, State Evolution and entropy estimation are fully automated. The source code is publicly available at <https://github.com/sphinxteam/tramp>.

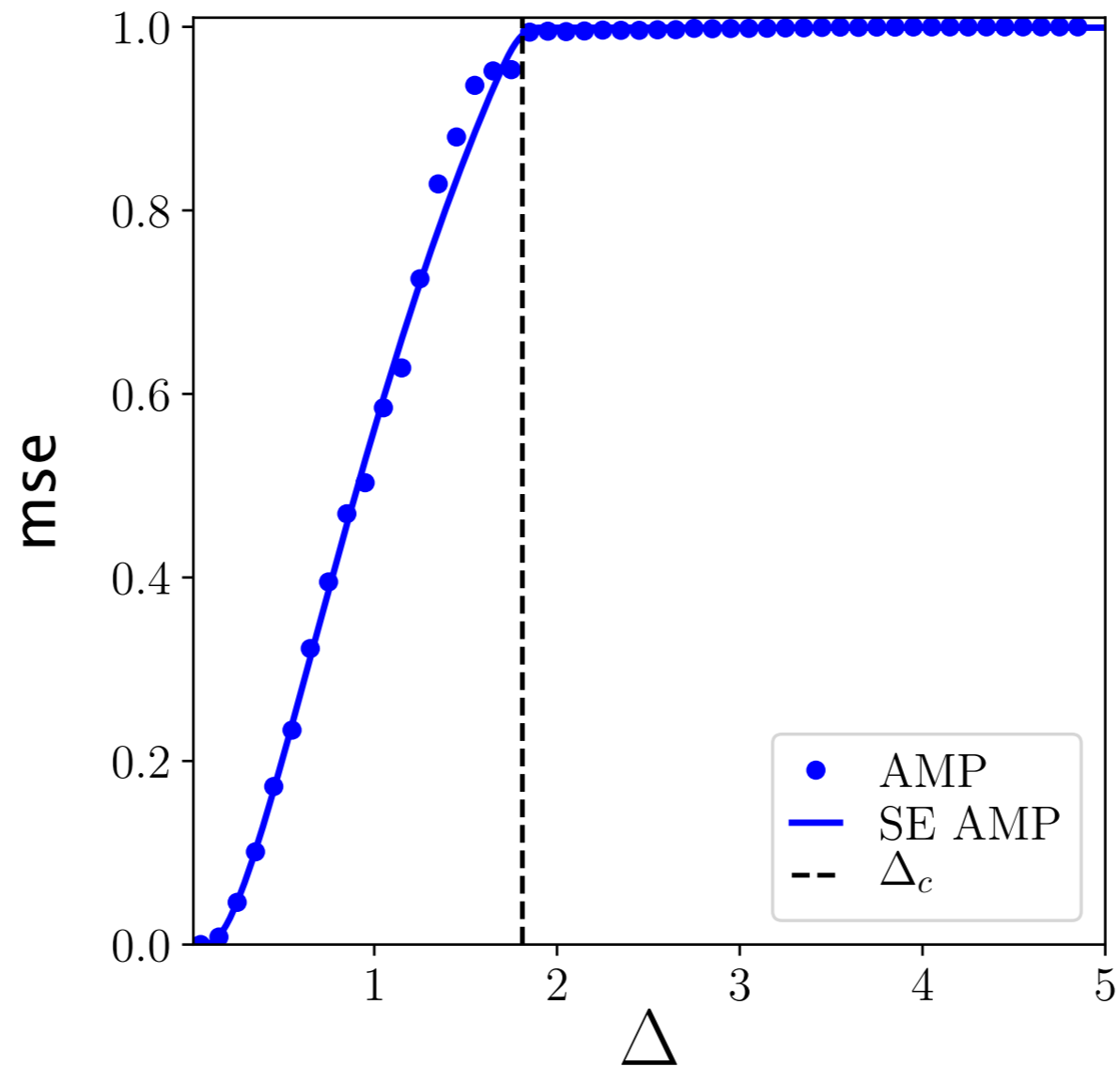
Keywords: Probabilistic programming, Bethe entropy, State Evolution, Expectation Propagation

<https://sphinxteam.github.io/tramp.docs>

[arXiv2004.01571](https://arxiv.org/abs/2004.01571)

Surprise

$$\mathbf{x} = \text{sign}(\mathbf{Wz})$$



No **statistical-to-algorithmic** gap!

Non-asymptotic: [Cocola, Hand, Voroninski 21']

Spectral method

PCA

Input: $Y \in \mathbb{R}^{p \times p}$

Output: $\hat{\mathbf{v}} = \sigma_{\max}(Y)$

Spectral method

PCA

Input: $Y \in \mathbb{R}^{p \times p}$

Output: $\hat{\mathbf{v}} = \sigma_{\max}(Y)$

L-AMP

Input: $Y \in \mathbb{R}^{p \times p}$, prior P_v

Output: $\hat{\mathbf{v}} = \sigma_{\max} \left(K_p \left(Y - I_p \right) \right)$

$$K_p = \mathbb{E}_{P_v} [\mathbf{v}\mathbf{v}^\top]$$

Example: for a linear layer

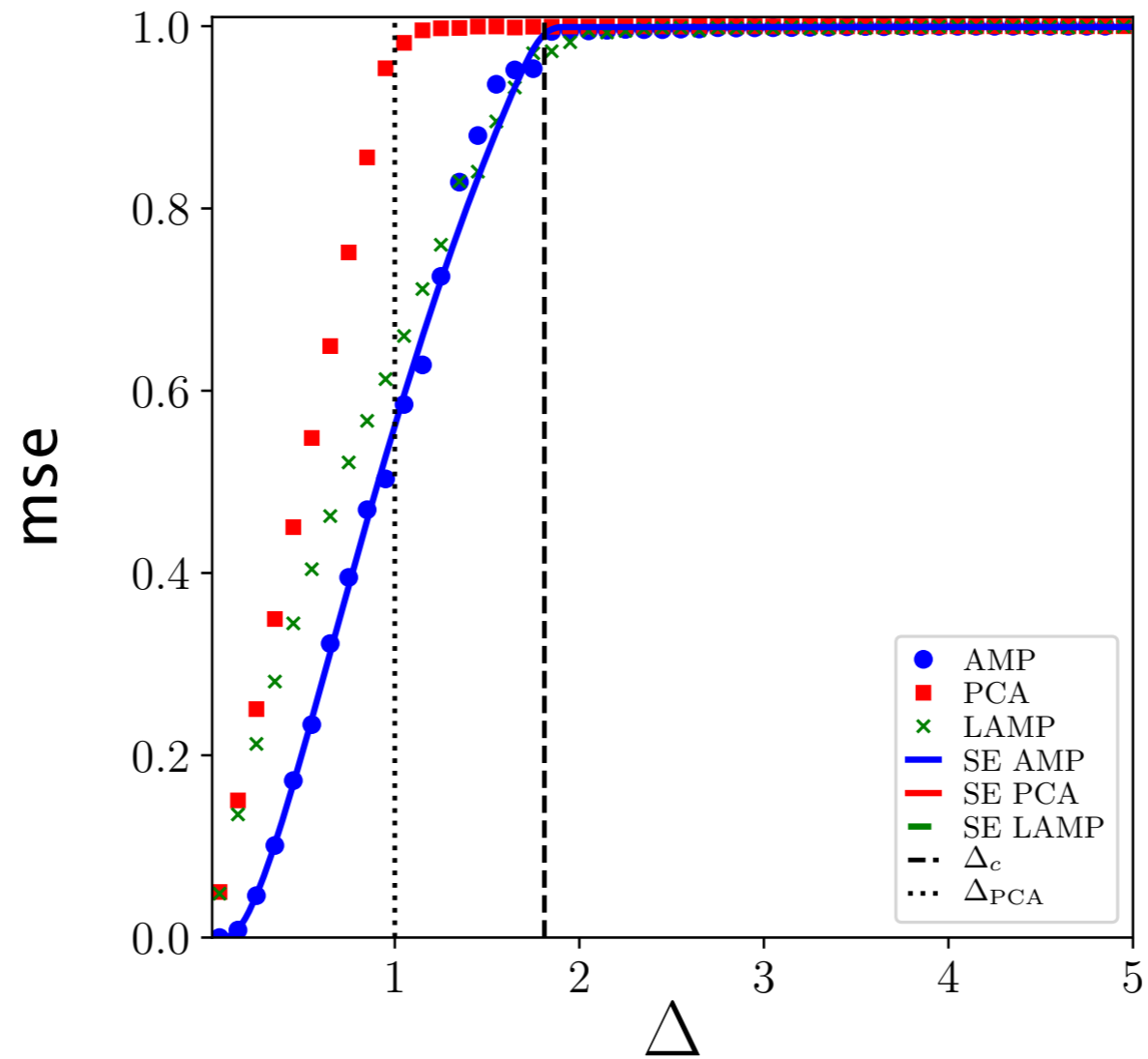
$$K_p = \frac{WW^\top}{p}$$

for a sign layer

$$K_p = \left(1 - \frac{2}{\pi} \right) I_p + \frac{2}{\pi} \frac{WW^\top}{p}$$

L-AMP

$$\mathbf{x} = \text{sign}(\mathbf{Wz})$$

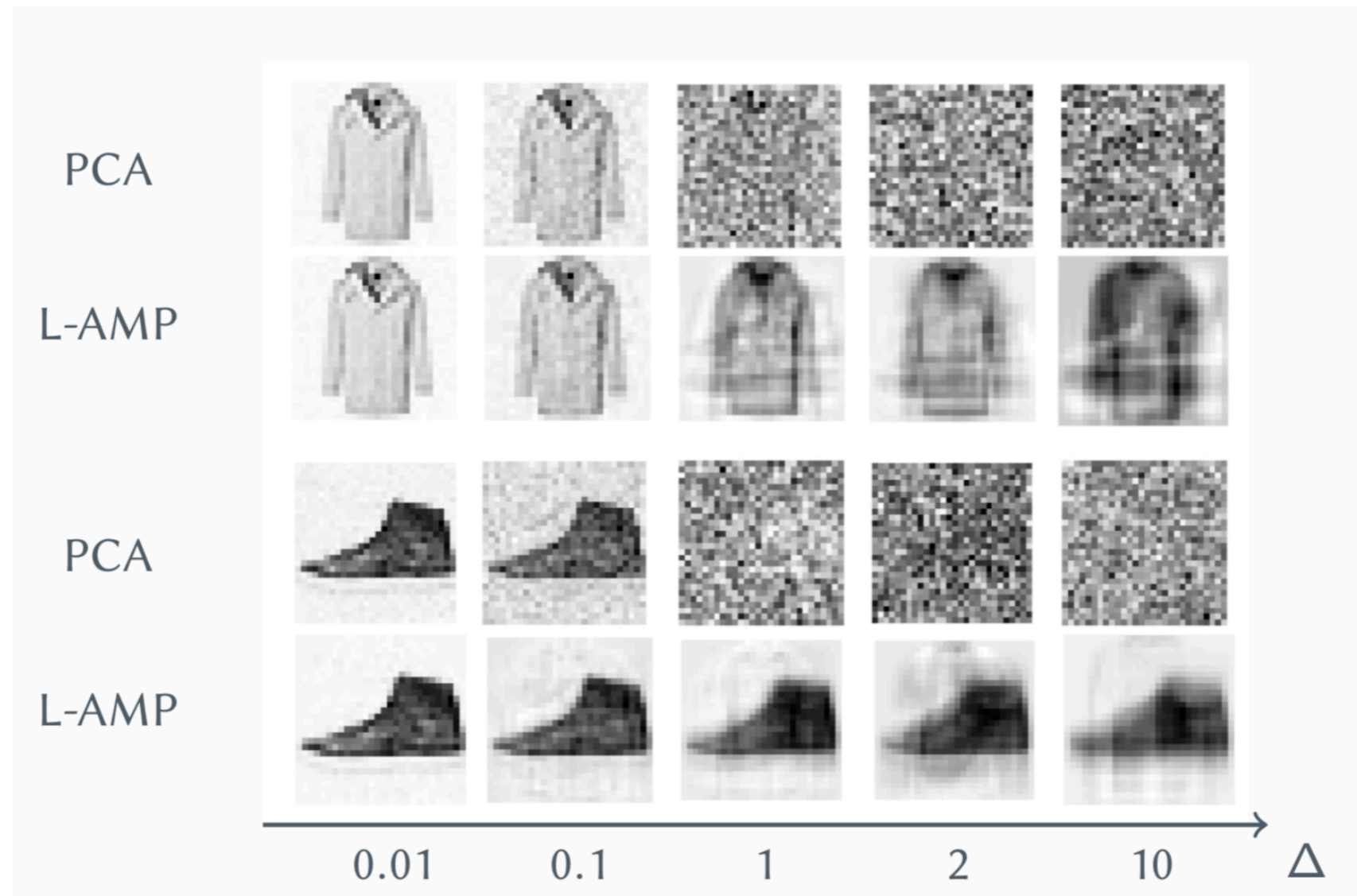


Same recovery threshold!

LAMP on Fashion MNIST

Spikes $\{x^\mu\}$ drawn from fashion MNIST

Use empirical covariance: $K_p = \mathbb{E}[xx^\top] \approx \frac{1}{m} \sum_{\mu=1}^m x^\mu x^{\mu\top}$



Other problems

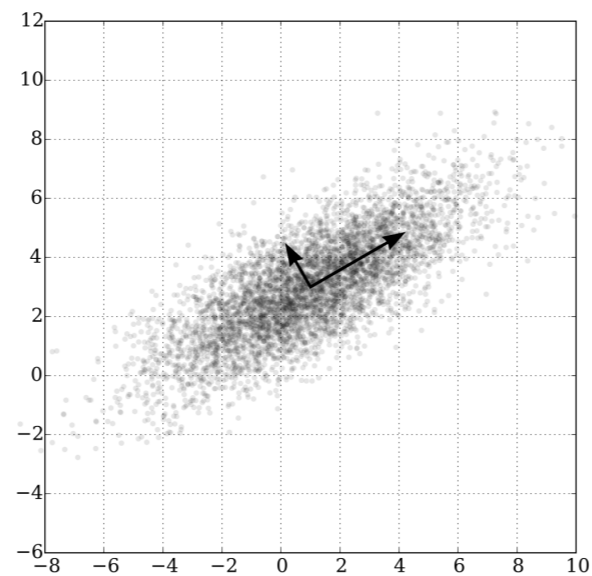
Denoising

$$\mathbf{y} = \mathbf{x} + \xi$$



Matrix factorisation

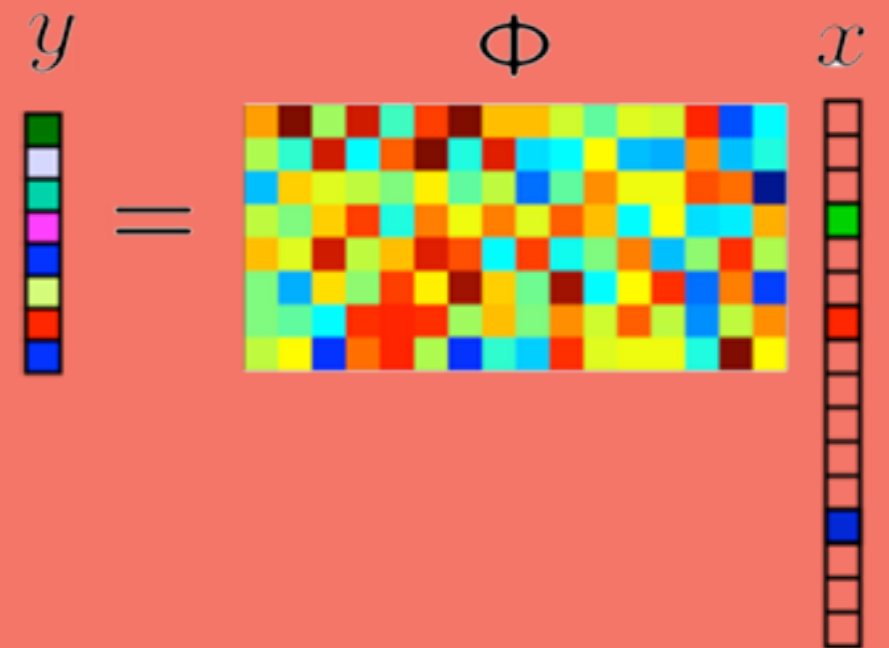
$$\mathbf{Y} = \mathbf{x}\mathbf{x}^T + \mathbf{Z}$$



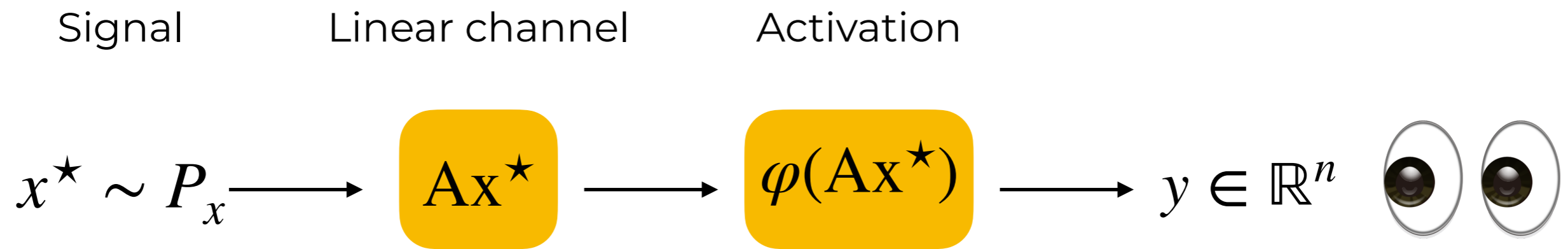
Compressed sensing
Phase retrieval

$$\mathbf{y} = \Phi\mathbf{x} + \xi$$

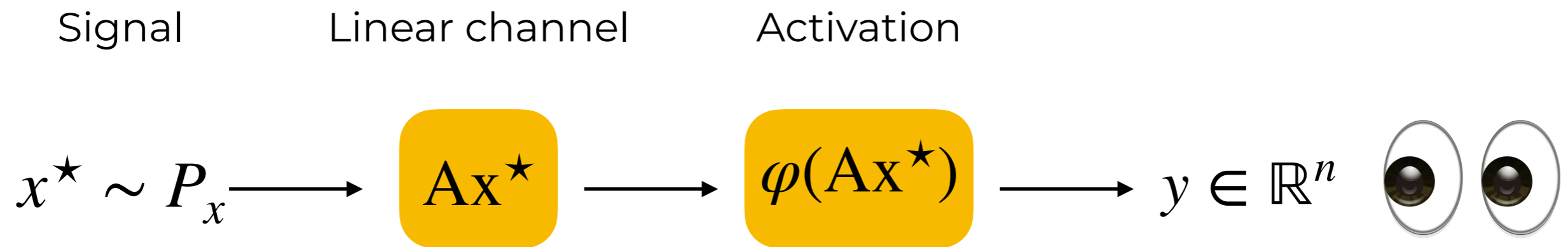
$$\mathbf{y} = |\mathbf{A}\mathbf{x}| + \xi$$



Generalised linear estimation



Generalised linear estimation

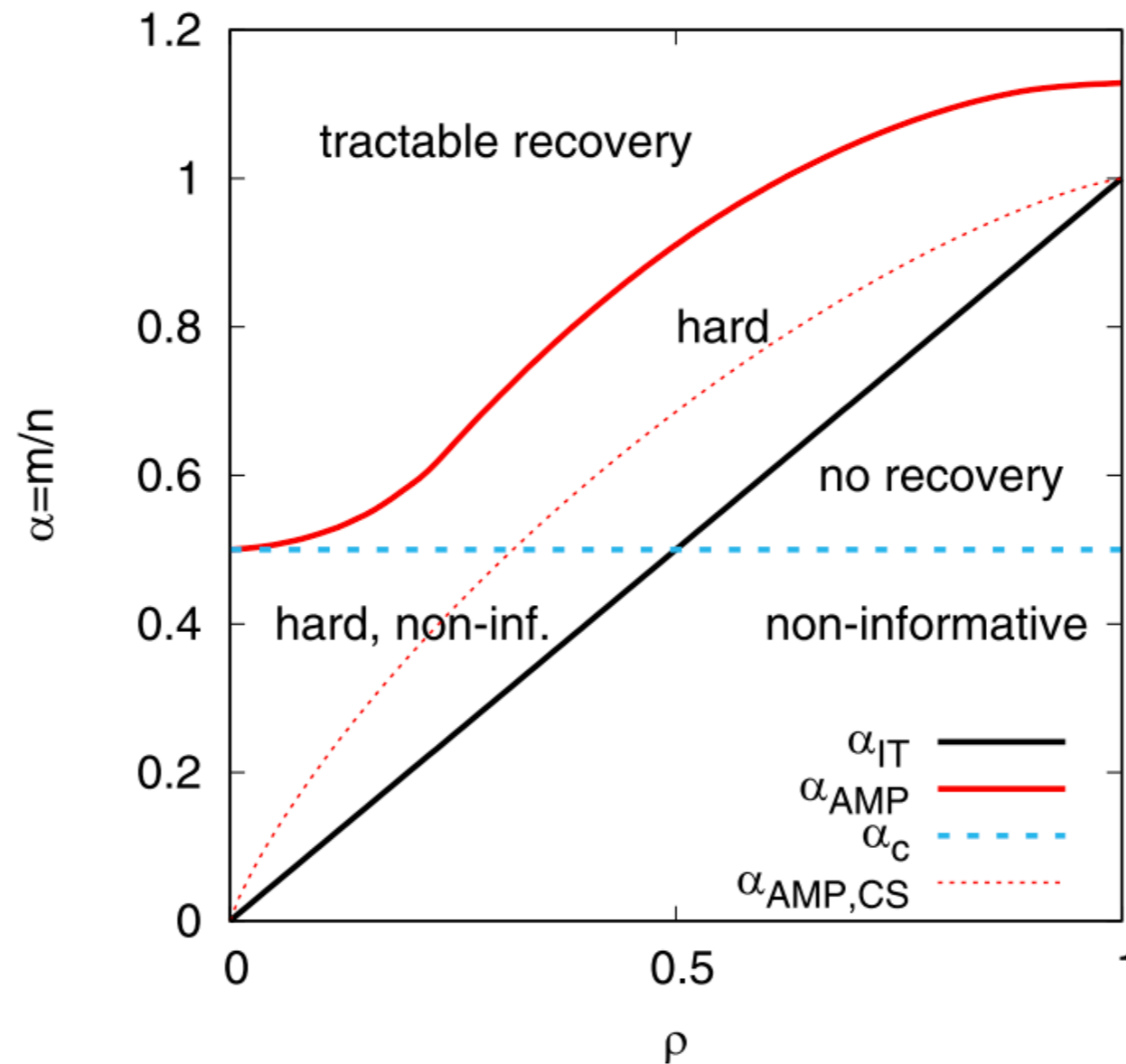


Goal: reconstruct x^\star from knowledge of P_x, A, σ
in the high-dimensional regime
 $n, p \rightarrow \infty$ and $n/p = O(1)$

Sparse case

$$x^* \sim \prod_{i=1}^p [(1 - \rho)\delta_0 + \rho\mathcal{N}(0, 1)]$$

$$y = |Ax^*|$$

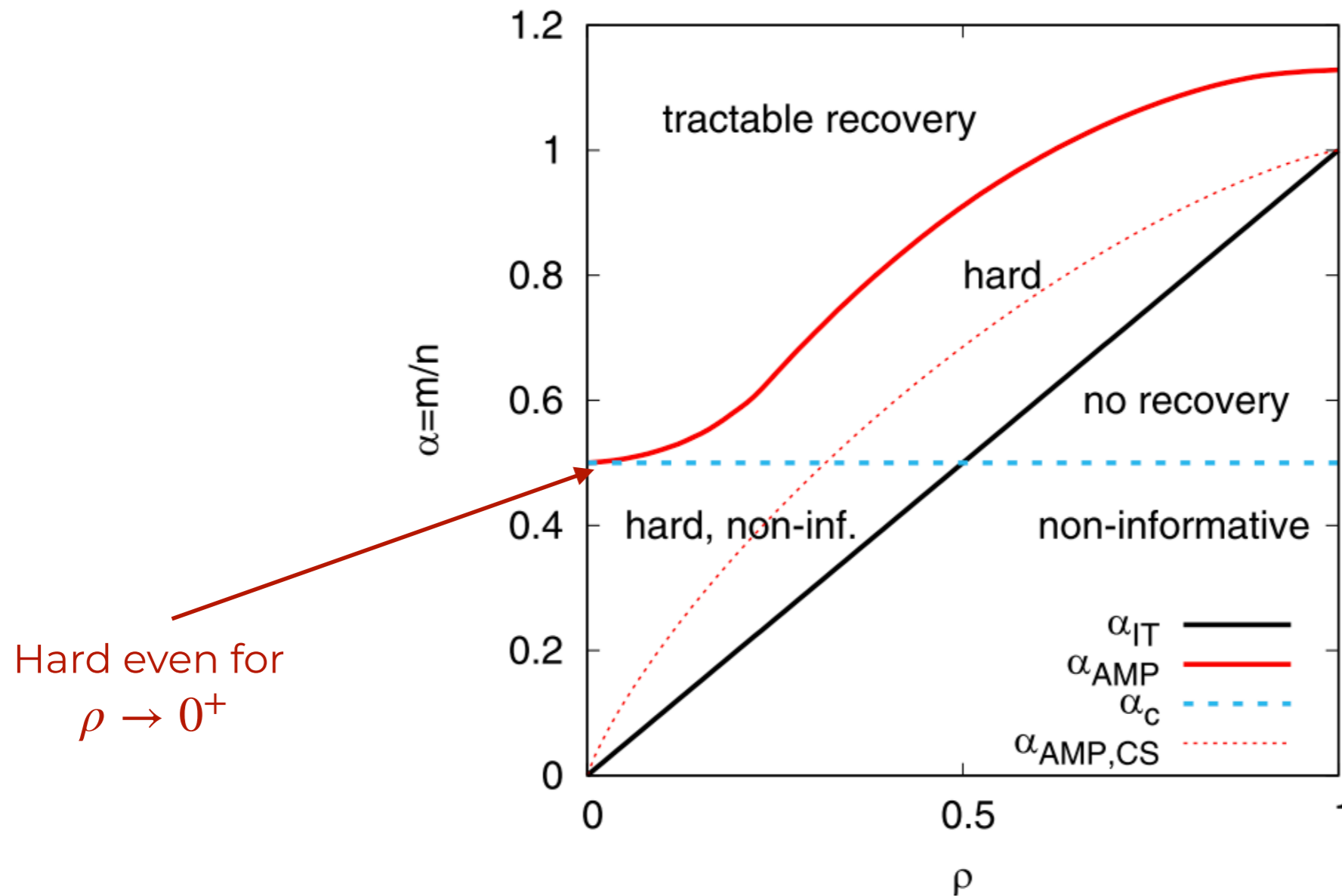


[Barbier et al. 16']

Sparse case

$$x^* \sim \prod_{i=1}^p [(1 - \rho)\delta_0 + \rho\mathcal{N}(0, 1)]$$

$$y = |Ax^*|$$

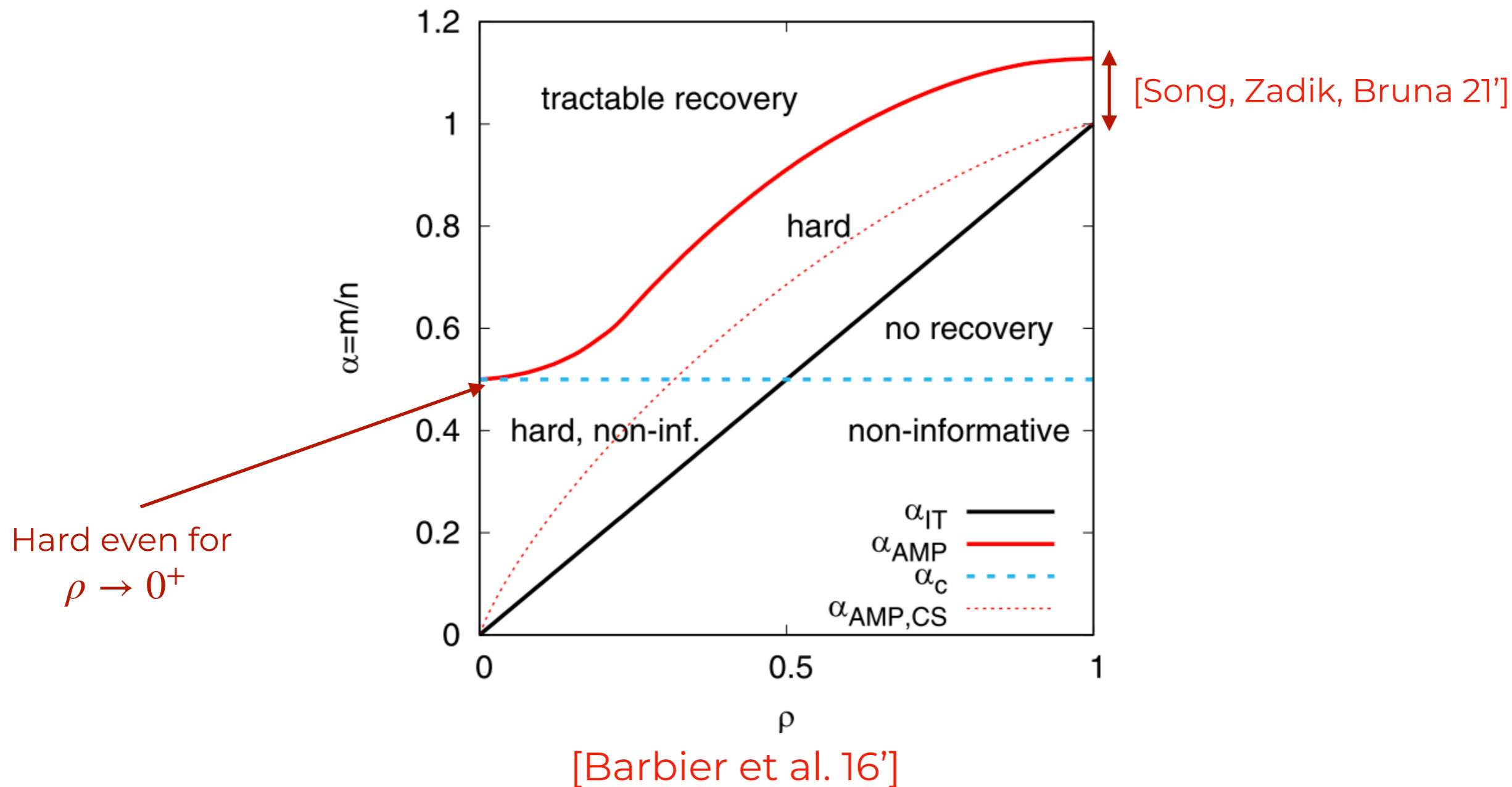


[Barbier et al. 16']

Sparse case

$$x^* \sim \prod_{i=1}^p [(1 - \rho)\delta_0 + \rho\mathcal{N}(0, 1)]$$

$$y = |Ax^*|$$



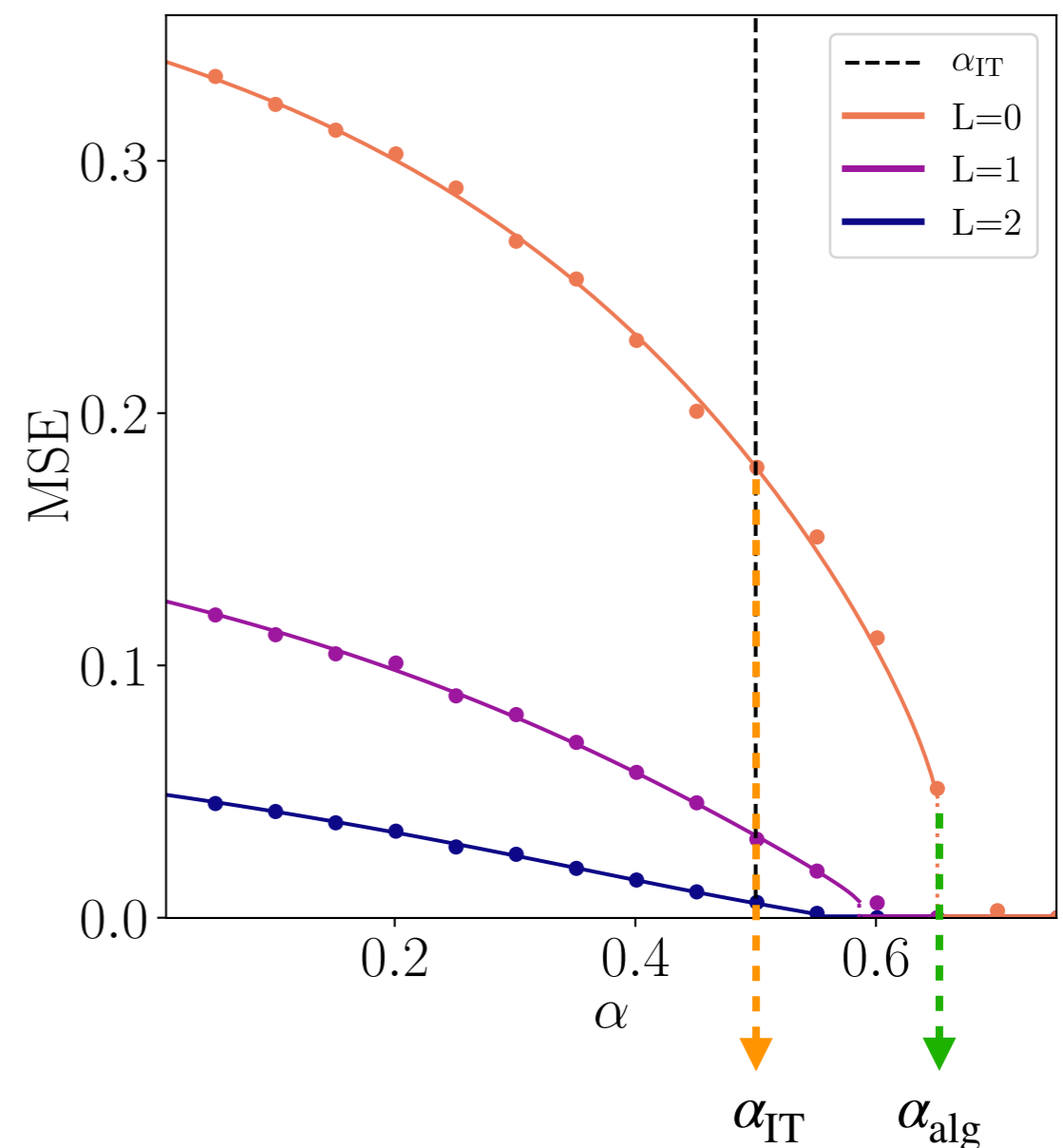
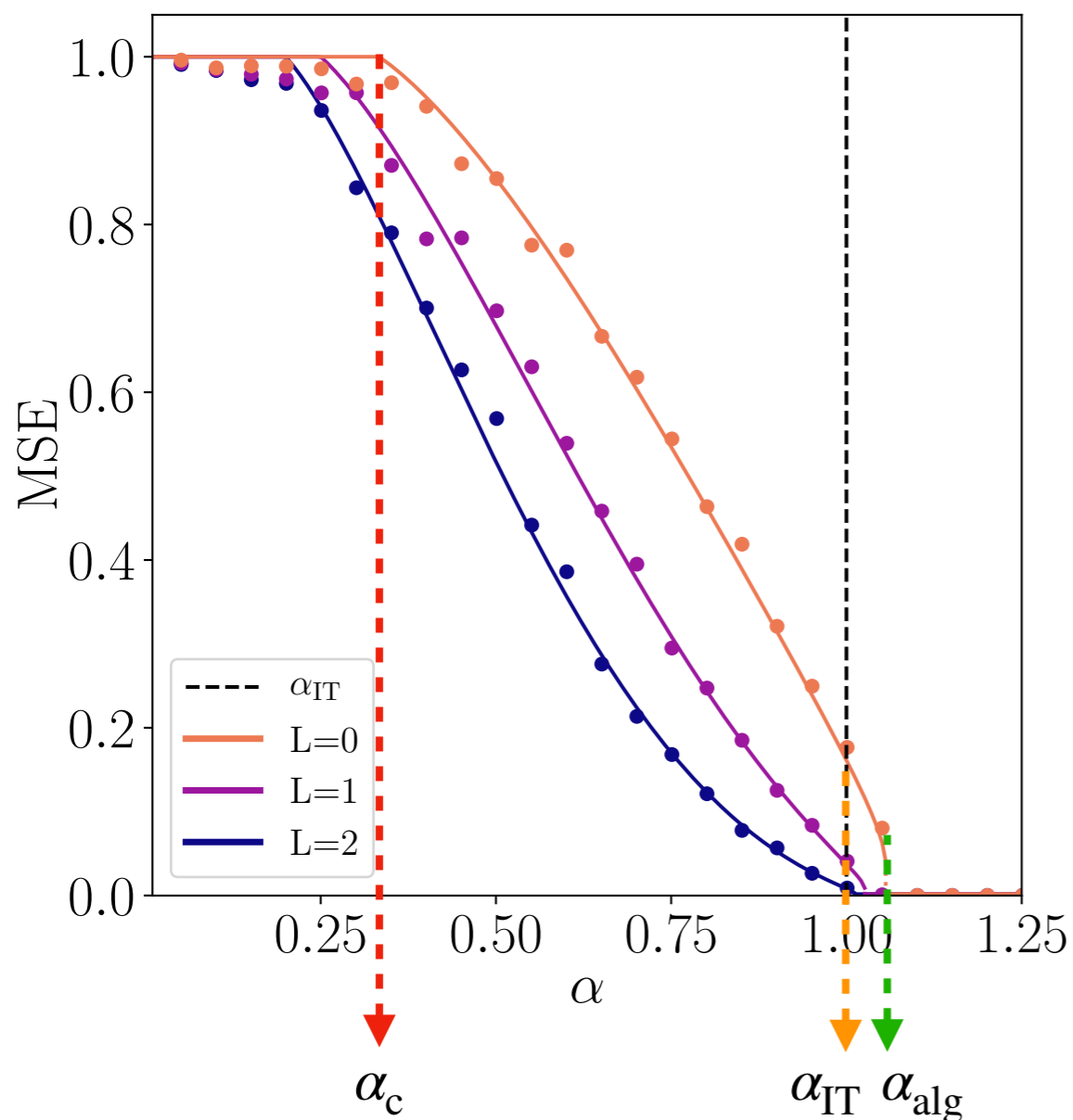
Generative network

$$\mathbf{x}^* = \sigma^{(L)} \left(\mathbf{W}^{(L)} \sigma^{(L-1)} \left(\mathbf{W}^{(L-1)} \dots \sigma^{(1)} \left(\mathbf{W}^{(1)} \mathbf{z} \right) \dots \right) \right)$$

$$y = |\mathbf{A} \mathbf{x}^*|$$

$$\sigma(x) = x$$

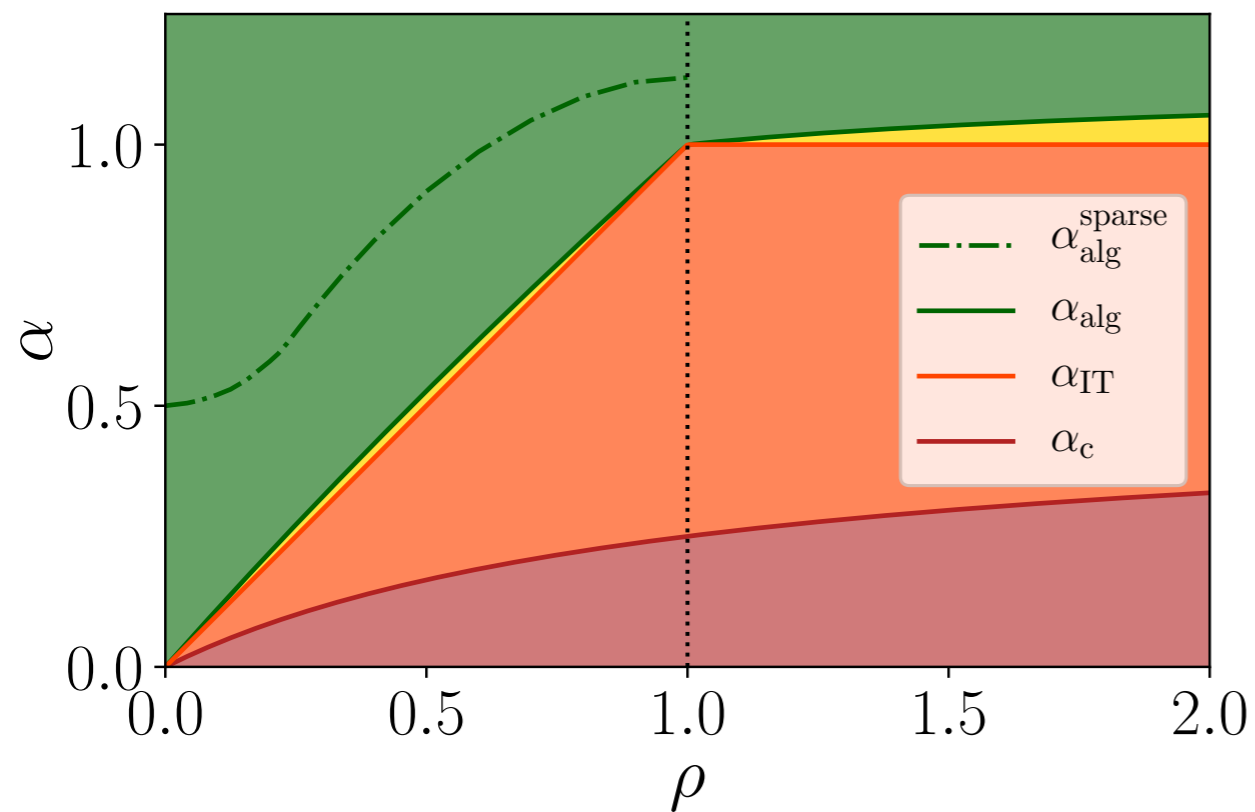
$$\sigma(x) = \text{relu}(x)$$



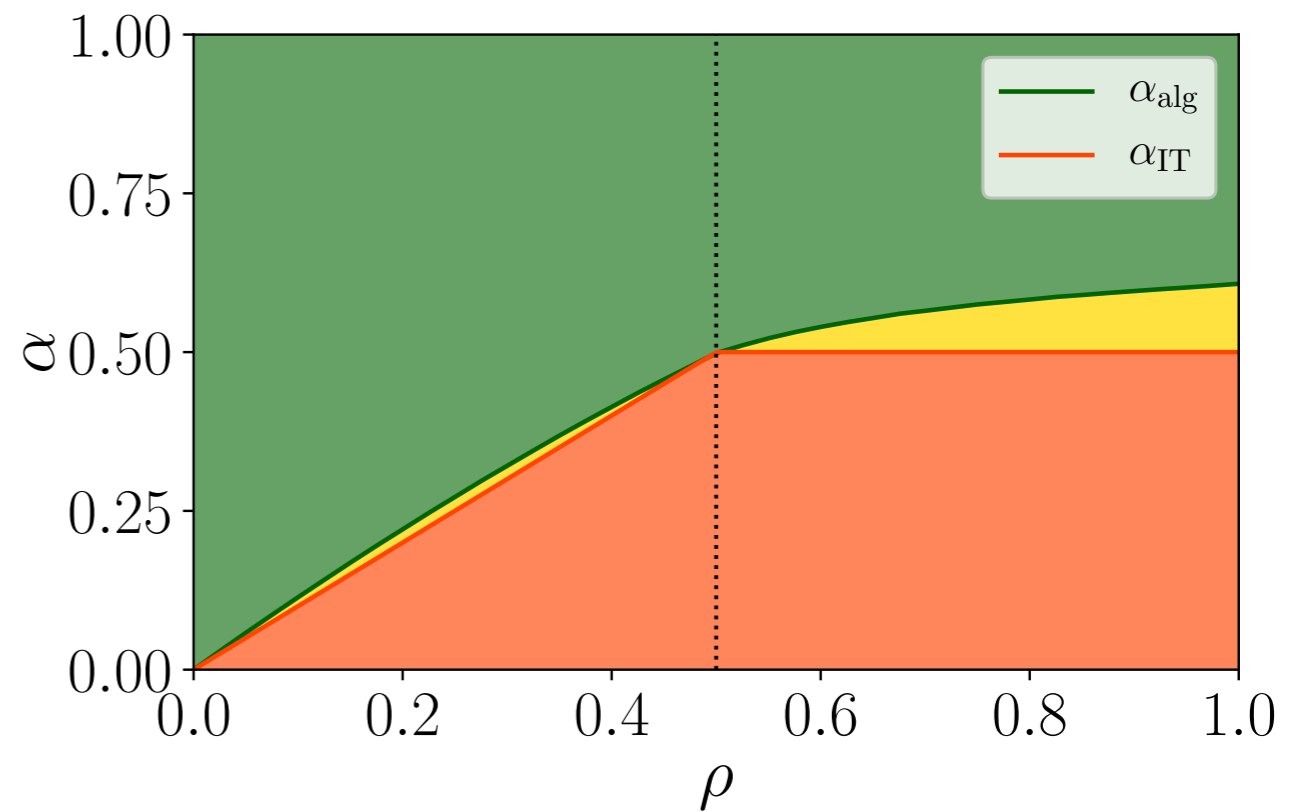
Generative network

$$y = |Ax^*|$$

$$\mathbf{x}^* = \mathbf{Wz}$$



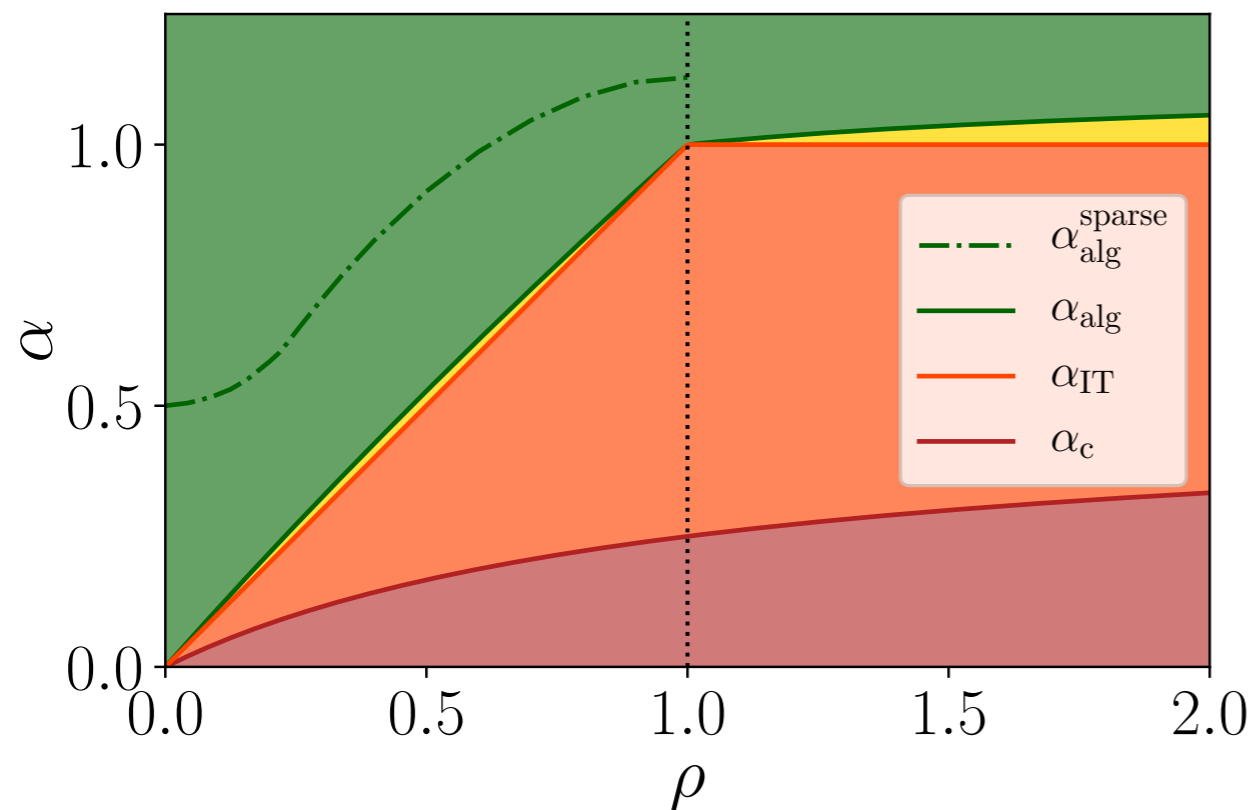
$$\mathbf{x}^* = \text{relu}(\mathbf{Wz})$$



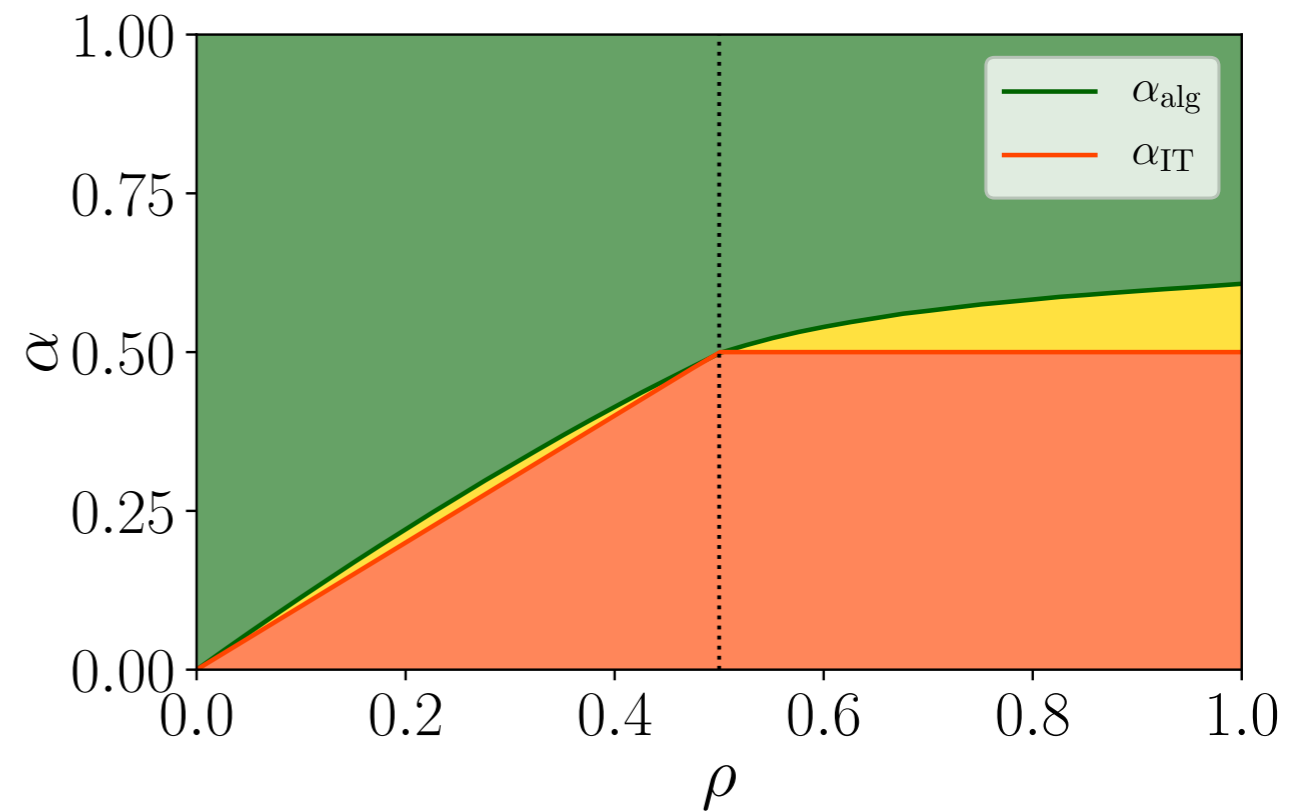
Generative network

$$y = |Ax^*|$$

$$\mathbf{x}^* = \mathbf{Wz}$$



$$\mathbf{x}^* = \text{relu}(\mathbf{Wz})$$



Generative priors have a **smaller** algorithmic gap!

c.f. [Hand, Leong, Voroninski 18']

Conclusions

Generative priors have comparatively **smaller algorithmic gaps** than their sparse counterparts

Algorithmic gap can be reduced with both **depth** and **compression**

Conclusions

Generative priors have comparatively **smaller algorithmic gaps** than their sparse counterparts

Algorithmic gap can be reduced with both **depth** and **compression**

Are generative priors the new sparsity?

or

Generative priors are **better** than sparsity?

| Thank you

Based on:
arXiv: 1905.12385,
1912.02008

Selected references

Eric W Tramel, Andre Manoel, Francesco Caltagirone, Marylou Gabrié, and Florent Krzakala.
Inferring sparsity: Compressed sensing using generalized restricted Boltzmann machines.
In 2016 IEEE Information Theory Workshop (ITW), pages 265–269. IEEE, 2016

Ashish Bora, Ajil Jalal, Eric Price, and Alexandros G Dimakis.
Compressed sensing using generative models.
In Proceedings of the 34th International Conference on Machine Learning-Volume 70,
pages 537–546. JMLR. org, 2017

Soledad Villar.
Generative models are the new sparsity?
<https://solevillar.github.io/2018/03/28/SUNLayer.html>, 2018.

Jean Barbier, Mohamad Dia, Nicolas Macris, Florent Krzakala, Thibault Lesieur, and Lenka Zdeborová.
Mutual information for symmetric rank-one matrix estimation: A proof of the replica formula.
In Advances in Neural Information Processing Systems, pages 424–432, 2016

Thibault Lesieur, Florent Krzakala, and Lenka Zdeborová.
Constrained low-rank matrix estimation: phase transitions, approximate message passing and applications.
Journal of Statistical Mechanics: Theory and Experiment, 2017(7):073403, 2017.

Selected references

Jean Barbier and Nicolas Macris.

The adaptive interpolation method: a simple scheme to prove replica formulas in bayesian inference.
Probability Theory and Related Fields, pages 1–53, 2018.

Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová.

Optimal errors and phase transitions in high-dimensional generalized linear models.
Proceedings of the National Academy of Sciences, 116(12):5451–5460, 2019.

M. Gabrié, A. Manoel, C. Luneau, J. Barbier, N. Macris, F. Krzakala, and L. Zdeborová.

Entropy and mutual information in models of deep neural networks.
In Advances in Neural Information Processing Systems 31

Jinho Baik, Gérard Ben Arous, Sandrine Péché, et al.

Phase transition of the largest eigenvalue for non null complex sample covariance matrices.
The Annals of Probability, 33(5):1643–1697, 2005.

Andre Manoel, Florent Krzakala, Marc Mézard, and Lenka Zdeborová.

Multi-layer generalized linear estimation.
In 2017 IEEE International Symposium on Information Theory (ISIT),
pages 2098–2102. IEEE, 2017.

Recent development

On the Cryptographic Hardness of Learning Single Periodic Neurons

Min Jae Song^{*a}, Ilias Zadik^{*b}, and Joan Bruna^{a,b}

^aCourant Institute of Mathematical Sciences, New York University, New York

^bCenter for Data Science, New York University, New York

June 22, 2021

Moreover, we demonstrate the *necessity of noise* in the hardness result by designing a polynomial-time algorithm for learning certain families of such functions under exponentially small adversarial noise. Our proposed algorithm is not a gradient-based or an SQ algorithm, but is rather based on the celebrated Lenstra-Lenstra-Lovász (LLL) lattice basis reduction algorithm. Furthermore, in the absence of noise, this algorithm can be directly applied to solve CLWE detection (Bruna et al.'21) and phase retrieval with an optimal sample complexity of $d + 1$ samples. In the former case, this improves upon the quadratic-in- d sample complexity required in (Bruna et al.'21). In the latter case, this improves upon the state-of-the-art AMP-based algorithm, which requires approximately $1.128d$ samples (Barbier et al.'19).

See also [Andoni, Hsu, Shi, Sun 17']

AMP for one-layer prior

- 1: **Input:** $Y \in \mathbb{R}^{p \times p}$ and $W \in \mathbb{R}^{p \times k}$:
- 2: *Initialize with:* $\hat{\mathbf{v}}^{t=1} = \mathcal{N}(\mathbf{0}, \sigma^2 I_p)$, $\hat{\mathbf{z}}^{t=1} = \mathcal{N}(\mathbf{0}, \sigma^2 I_k)$, and $\hat{\mathbf{c}}_v^{t=1} = I_p$, $\hat{\mathbf{c}}_z^{t=1} = I_k$,
 $t = 1$.
- 3: **repeat**
- 4: **Spiked layer denoising:**
- 5: $\mathbf{B}_v^t = \frac{1}{\Delta} \frac{Y}{\sqrt{p}} \hat{\mathbf{v}}^t - \frac{1}{\Delta} \frac{(I_p^\top \hat{\mathbf{c}}_v^t)}{p} \hat{\mathbf{v}}^{t-1}$ and $A_v^t = \frac{1}{\Delta p} (\|\hat{\mathbf{v}}^t\|_2)^2 I_p$.
- 6: **Generative layer denoising:**
- 7: $V^t = \frac{1}{k} (I_k^\top \hat{\mathbf{c}}_z^t) I_p$, $\boldsymbol{\omega}^t = \frac{1}{\sqrt{k}} W \hat{\mathbf{z}}^t - V^t \mathbf{g}^{t-1}$ and $\mathbf{g}^t = f_{\text{out}}(\mathbf{B}_v^t, A_v^t, \boldsymbol{\omega}^t, V^t)$,
- 8: $\Lambda^t = \frac{1}{k} \|\mathbf{g}^t\|_2^2 t I_k$ and $\gamma^t = \frac{1}{\sqrt{k}} W^\top \mathbf{g}^t + \Lambda^t \hat{\mathbf{z}}^t$.
- 9: **Marginals estimation:**
- 10: $\hat{\mathbf{v}}^{t+1} = f_v(\mathbf{B}_v^t, A_v^t, \boldsymbol{\omega}^t, V^t)$ and $\hat{\mathbf{c}}_v^{t+1} = \partial_{\mathbf{B}} f_v(\mathbf{B}_v^t, A_v^t, \boldsymbol{\omega}^t, V^t)$,
- 11: $\hat{\mathbf{z}}^{t+1} = f_z(\gamma^t, \Lambda^t)$ and $\hat{\mathbf{c}}_z^{t+1} = \partial_{\gamma} f_z(\gamma^t, \Lambda^t)$,
- 12: $t = t + 1$.
- 13: **until** Convergence.
- 14: **Output:** $\hat{\mathbf{v}}, \hat{\mathbf{z}}$.

Replica symmetric formula

For a L-layer deep, feed-forward, random generative prior, the *free entropy*

$$\Phi = \text{extr}_{q_x, \hat{q}_x, \{q_l, \hat{q}_l\}} \left\{ -\frac{1}{2} \hat{q}_x q_x - \frac{\rho}{2} \sum_{l=1}^L \beta_l q_l \hat{q}_l + \alpha \Psi_y(q_x) + \rho \sum_{l=2}^L \beta_l \Psi_{\text{out}}^{(l)}(\hat{q}_l, q_{l-1}) + \Psi_{\text{out}}^{(L+1)}(\hat{q}_x, q_L) + \rho \Psi_z(\hat{q}_z) \right\}$$

$$\Psi_{\text{out}}^{(l)}(r, s) = \mathbb{E}_{\xi, \eta} \left[\mathcal{L}_{\text{out}}^{(l)}(\sqrt{r}\xi, r, \sqrt{s}\xi, \rho_{l-1} - s) \log \mathcal{L}_{\text{out}}^{(l)}(\sqrt{r}\xi, r, \sqrt{s}\xi, \rho_{l-1} - s) \right]$$

$$\Psi_z(t) = \mathbb{E}_{\xi} \left[\mathcal{L}_z(\sqrt{t}\xi, t) \log \mathcal{L}_z(\sqrt{t}\xi, t) \right]$$

$$Q_{\text{out}}^{(l)}(x, z; B, A, \omega, V) = \frac{e^{-\frac{A}{2}x^2 + Bx}}{\mathcal{L}_{\text{out}}(B, A, \omega, V)} \frac{e^{-\frac{1}{2V}(z - \omega)^2}}{\sqrt{2\pi V}} P_{\text{out}}^{(l)}(x | z)$$

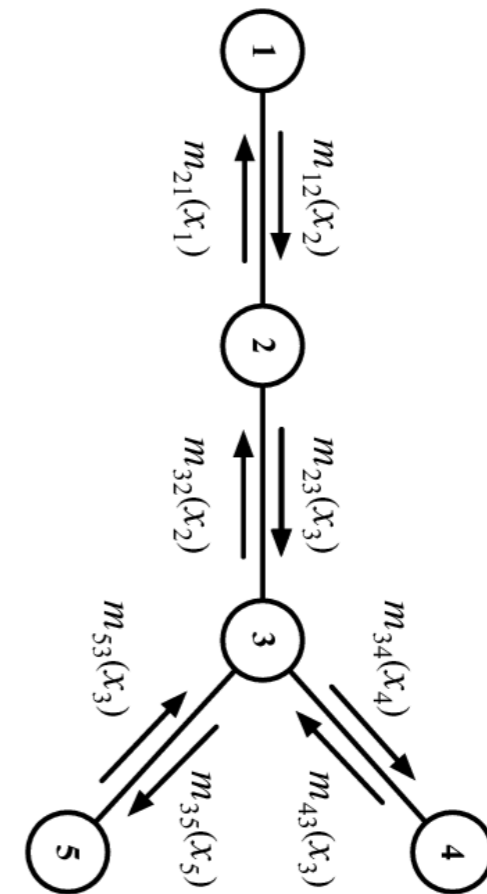
$$Q_z(z; B, A) = \frac{e^{-\frac{A}{2}z^2 + Bz}}{\mathcal{L}_z(B, A)} P_z(z)$$

Message Passing

This problem screams for a Approximated Message Passing algorithm!

Key facts about AMP:

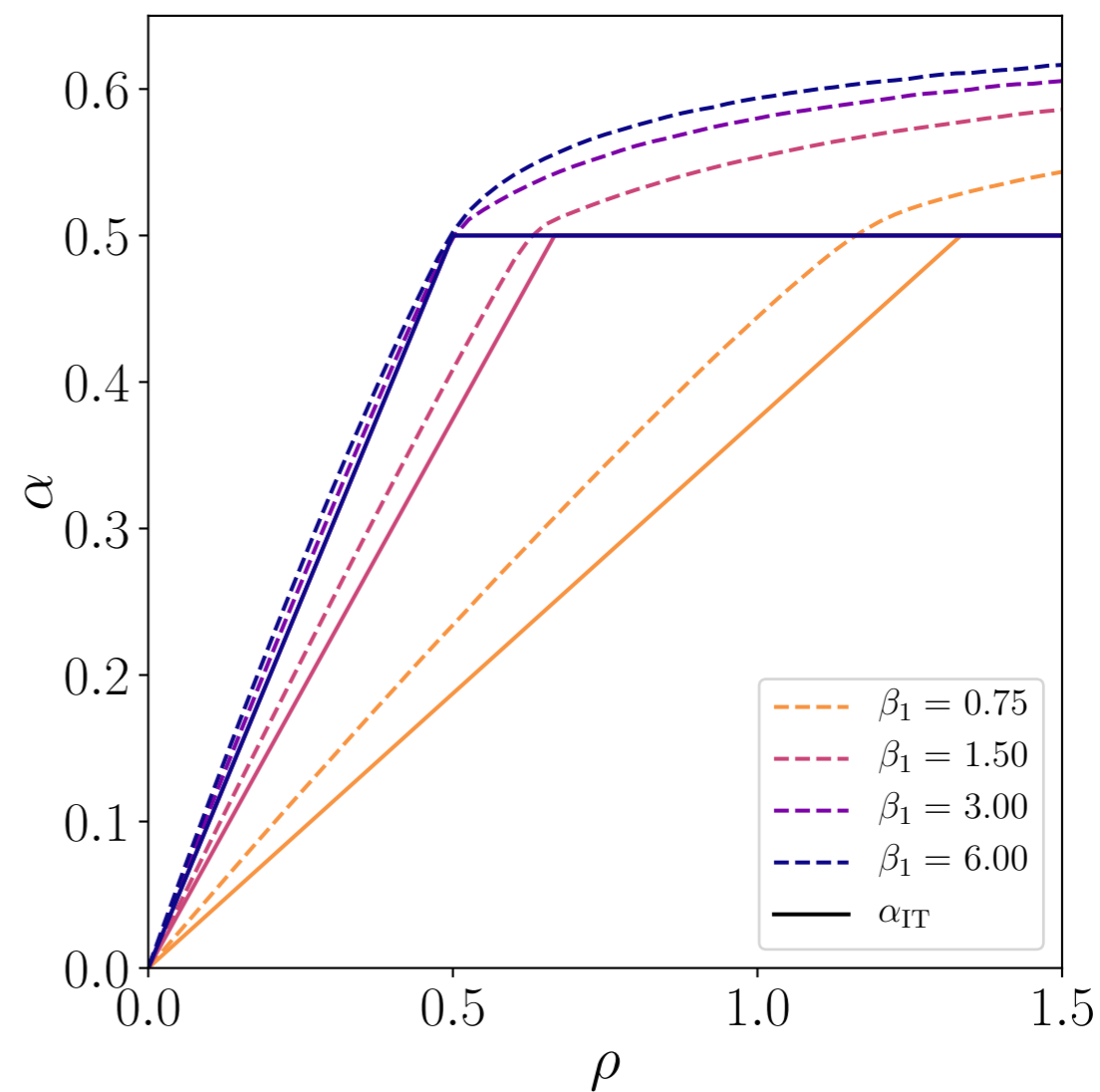
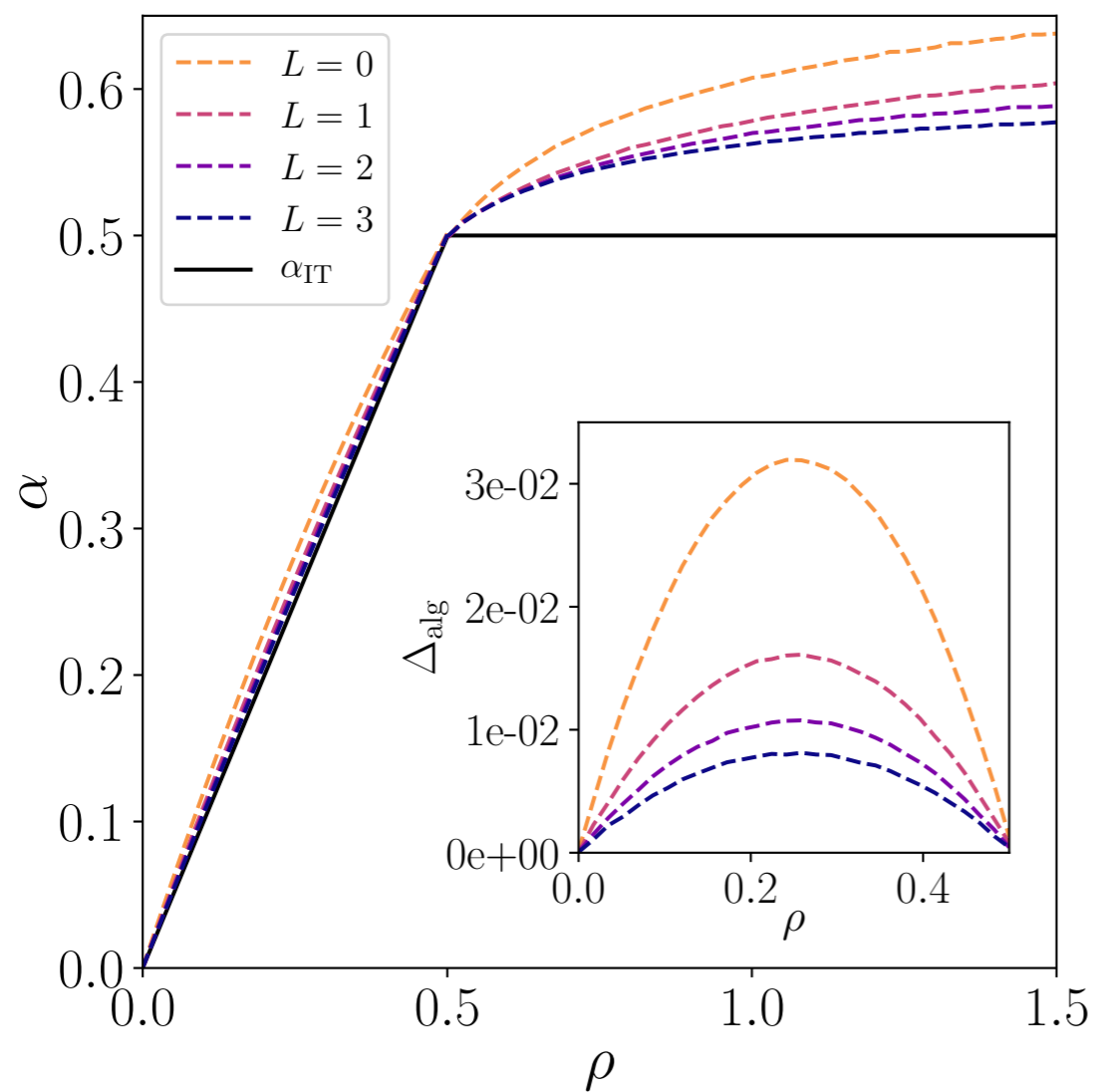
- Polynomial time iterative algorithm that approximate marginals from the posterior
- Can be derived **systematically** from Belief propagation.
- Its performance can be **tracked analytically**.
- **Best well-known polynomial time algorithm** for a large class of random problems.



Role of depth / width

$$y = |Ax^*|$$

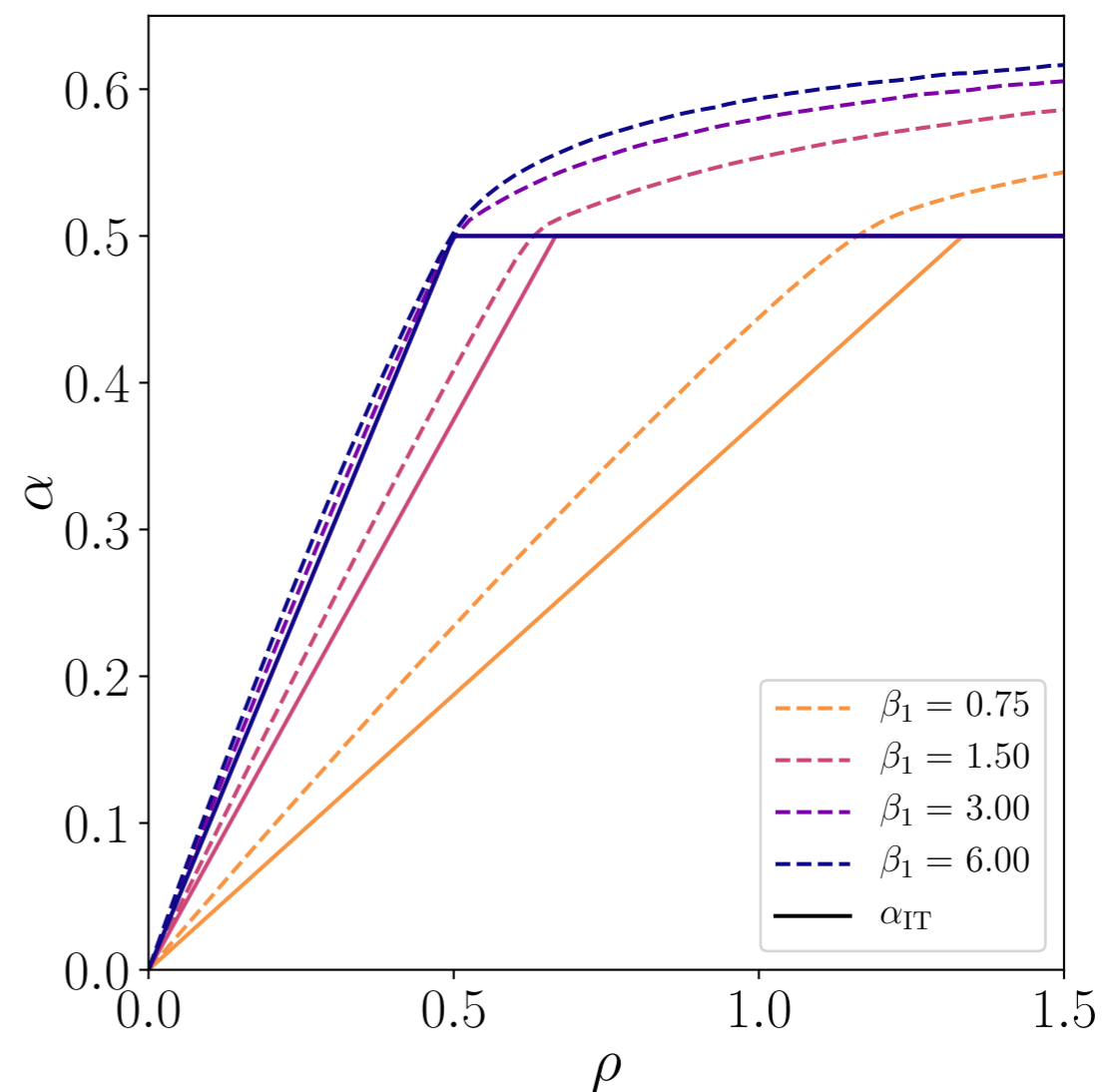
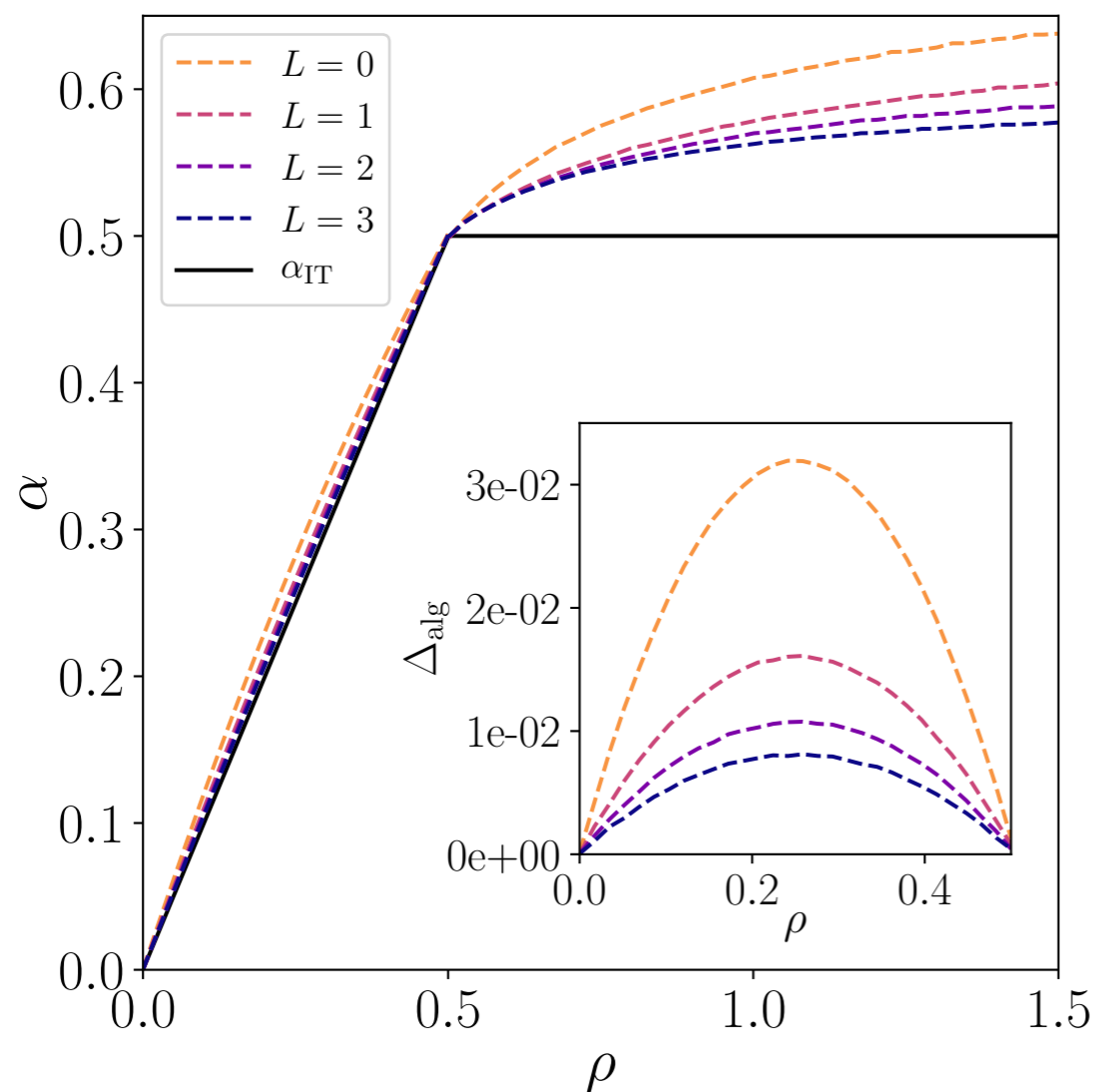
$$\mathbf{x}^* = \text{relu}(\mathbf{Wz})$$



Role of depth / width

$$y = |Ax^*|$$

$$\mathbf{x}^* = \text{relu}(\mathbf{Wz})$$



Gap can be made smaller with **increasing depth and compression**

[Bora et al 17'] [Hand et al 18']