

Generalization in diffusion models arises from geometry-adaptive harmonic representations

Zahra Kadkhodaie
September 2024



Florentin Guth



Eero Simoncelli



Stéphane Mallat



Diffusion models embed densities

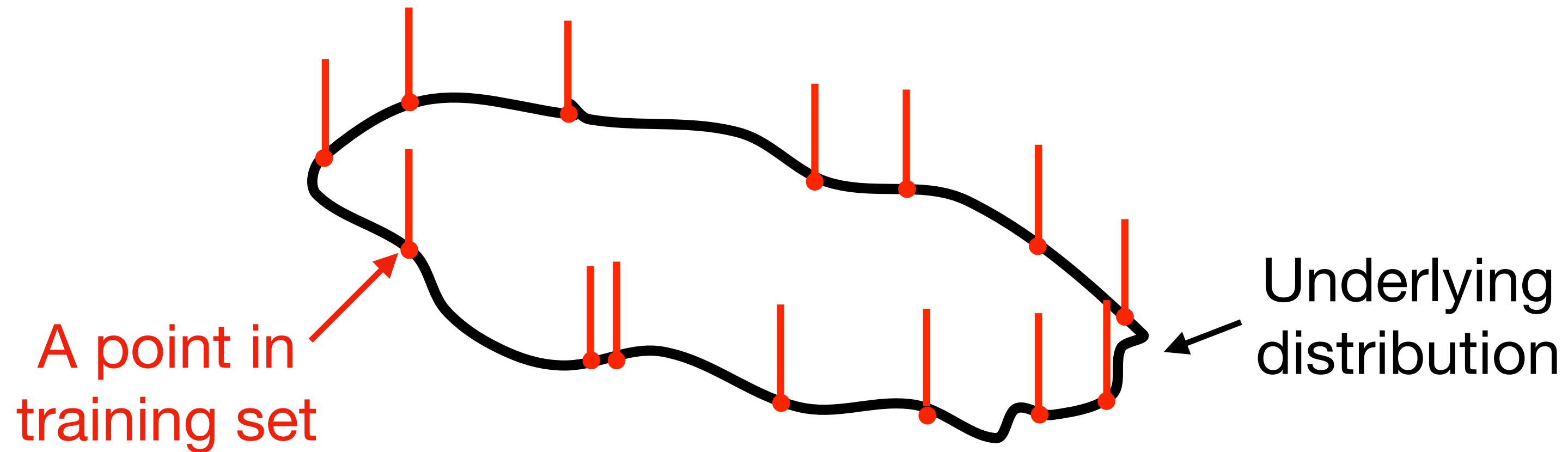
- Sample from a **learned** density
[Song & Ermon 2019; Ho et al 2020]
- How is this possible, given the
“Curse of dimensionality” ?!

generated by a diffusion model

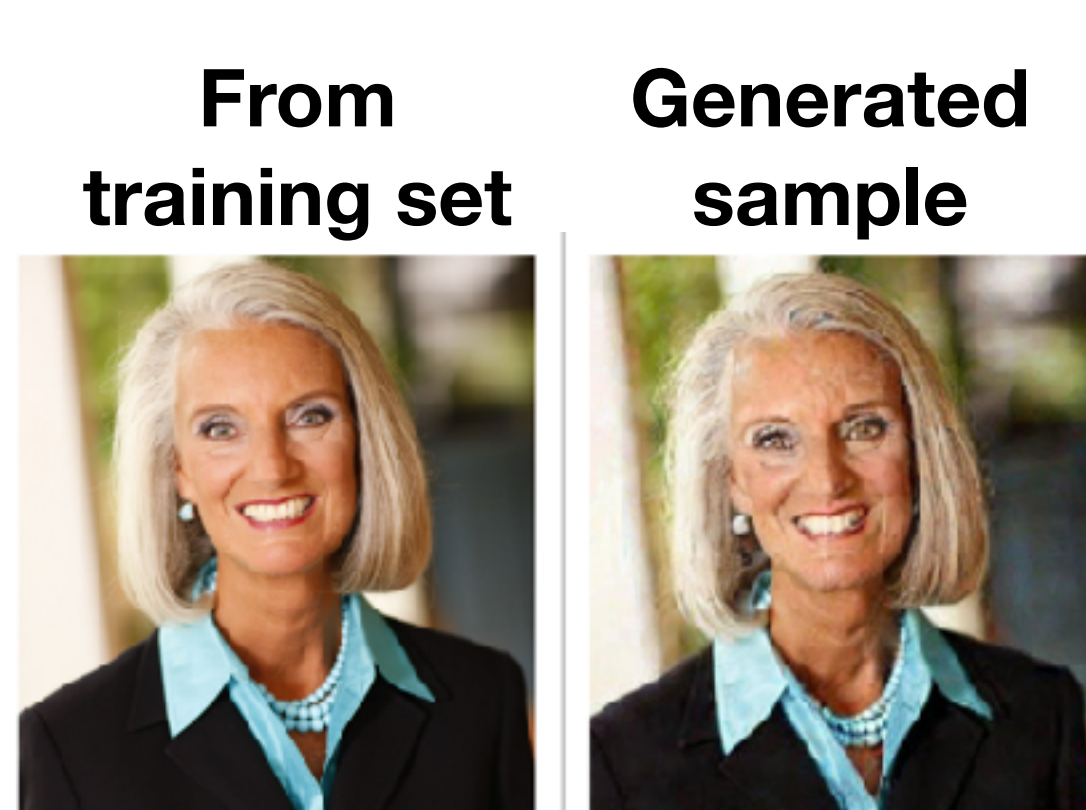


[Ho et al, 2022]

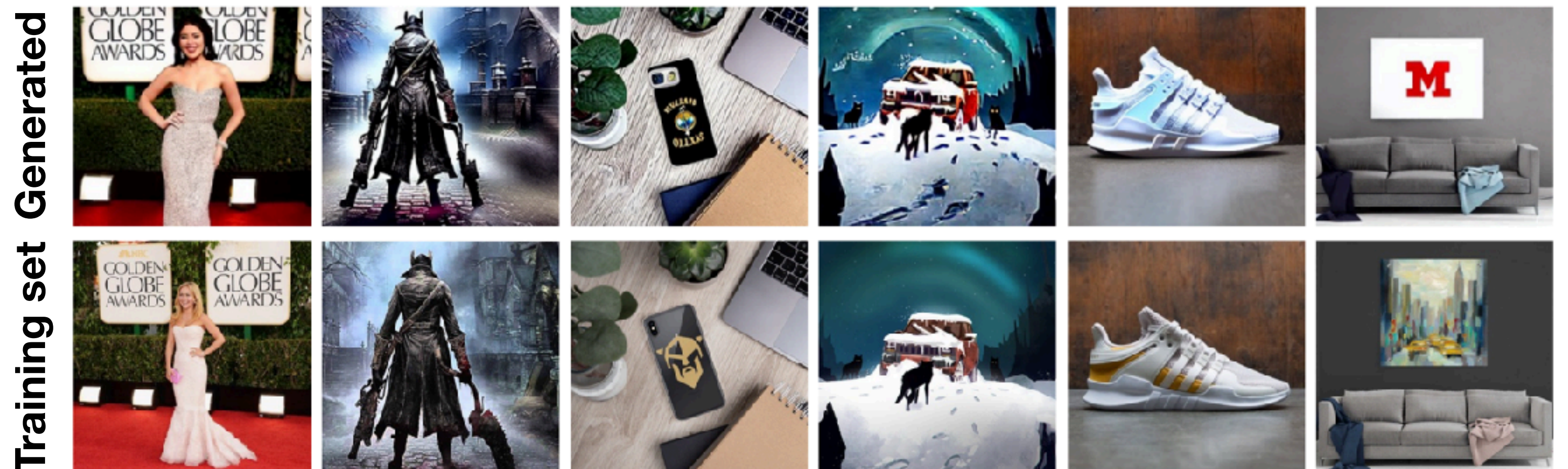
Memorization vs. generalization



A collection of delta functions

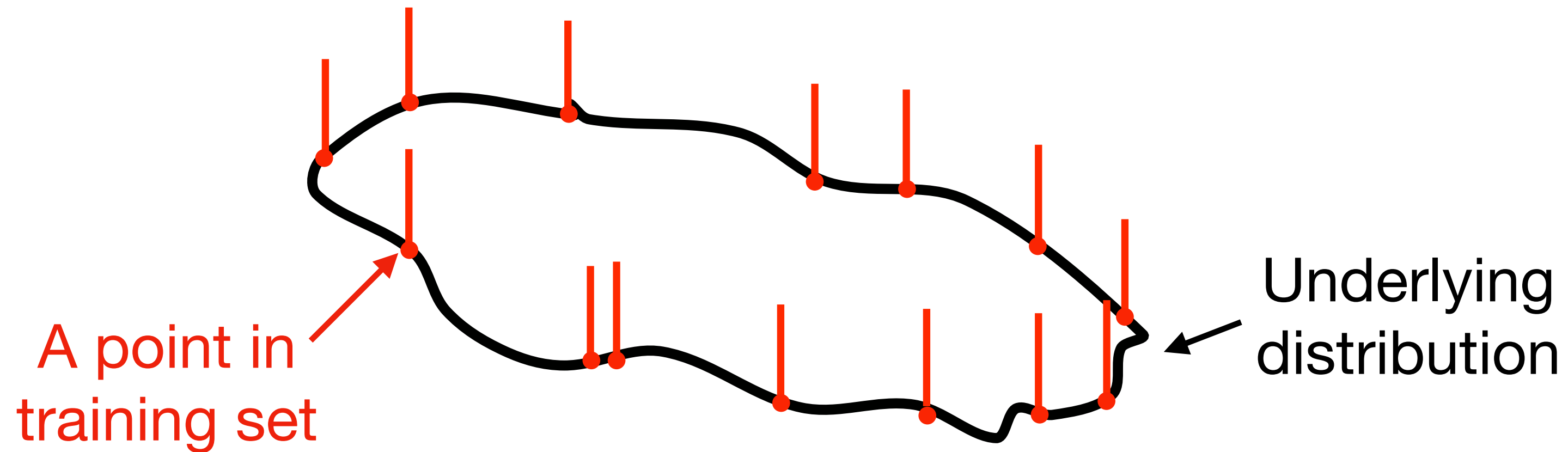


[Carlini et al, 2023]

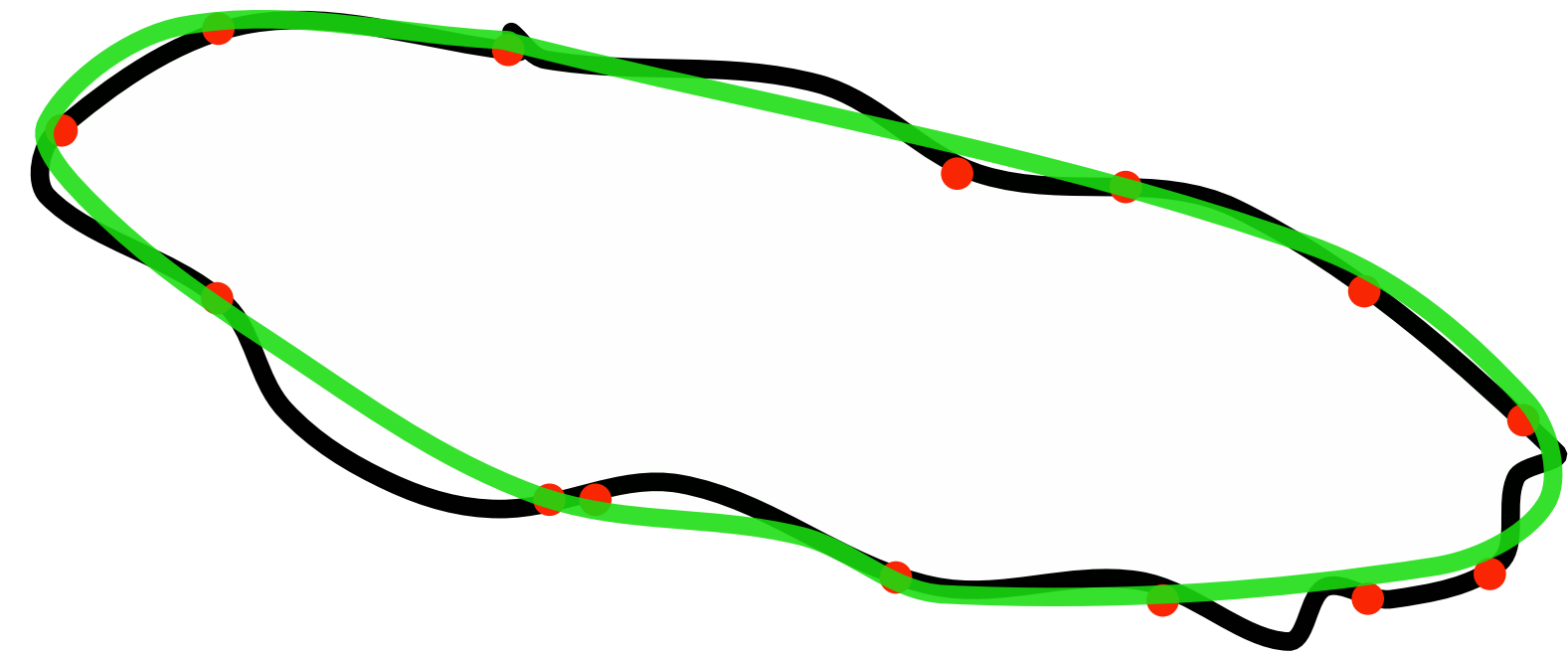


[Somepalli et al, 2023]

Memorization vs. generalization

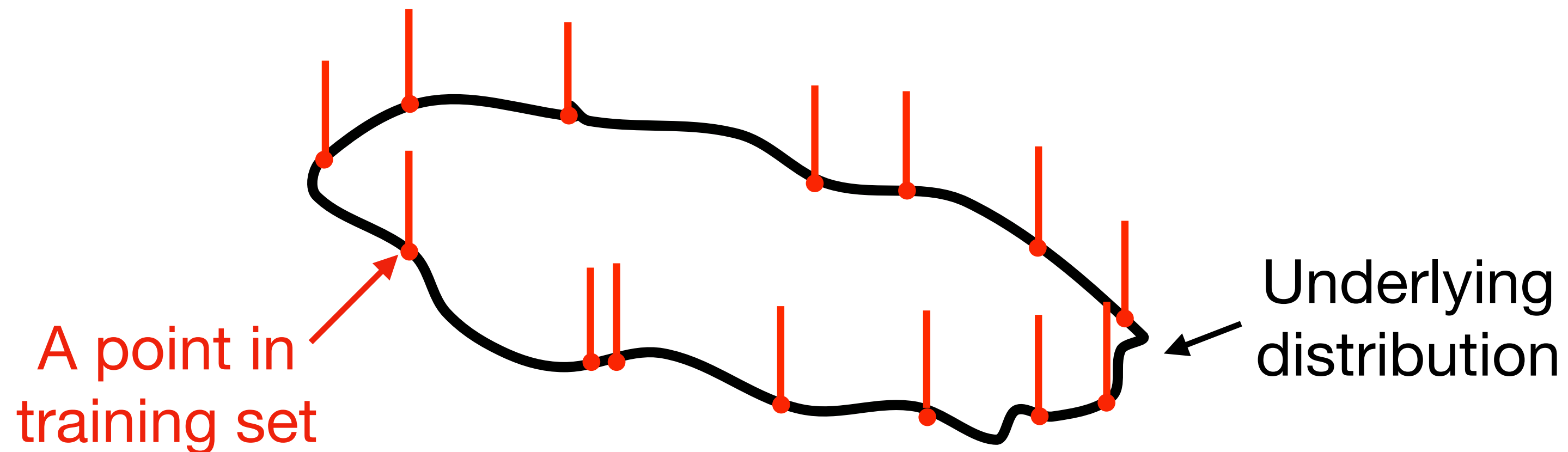


A collection of delta functions

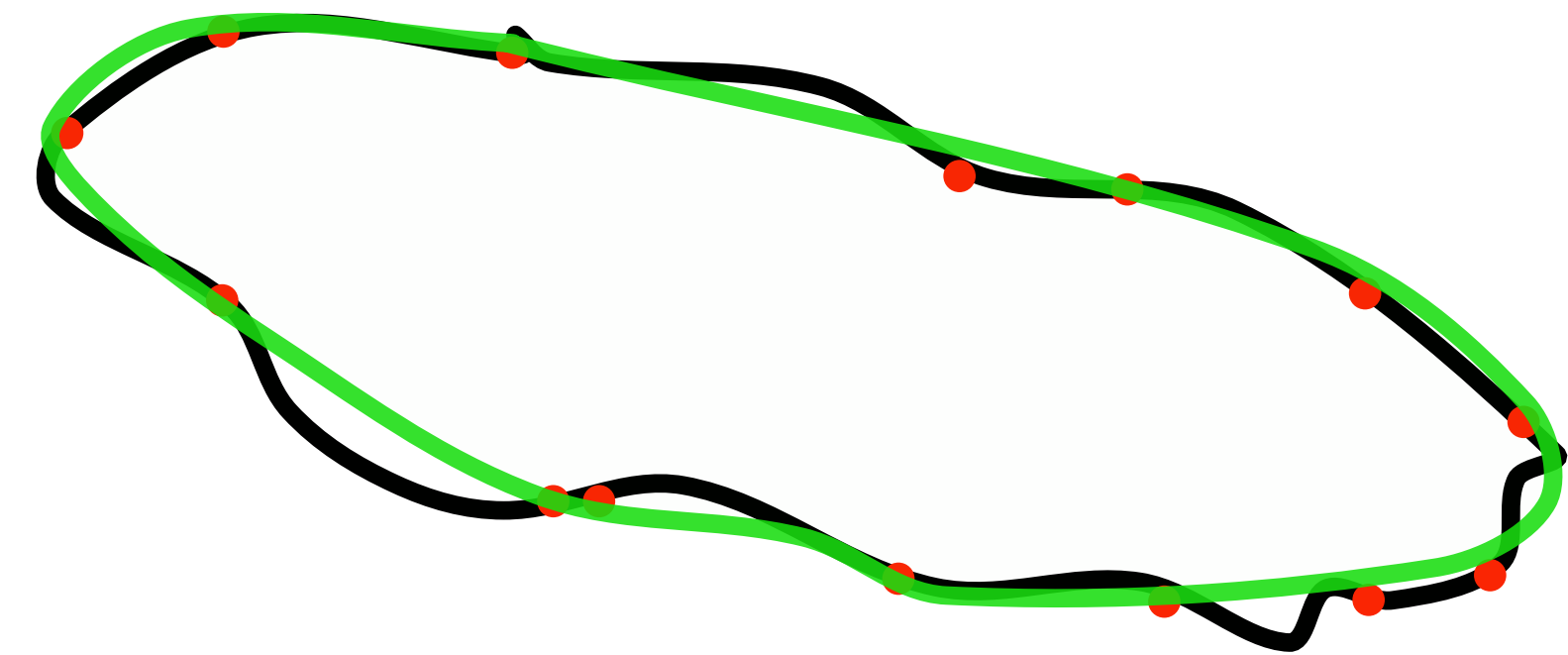


A continuous model of the underlying distribution

Memorization vs. generalization



A collection of delta functions



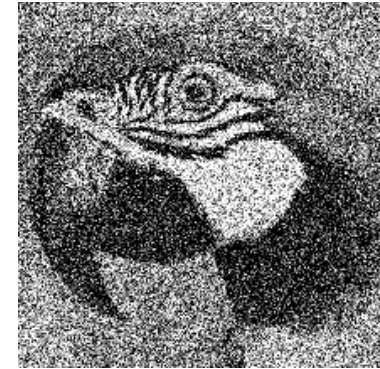
A continuous model of the underlying distribution

1. Can diffusion models generalize?
2. If so, how?

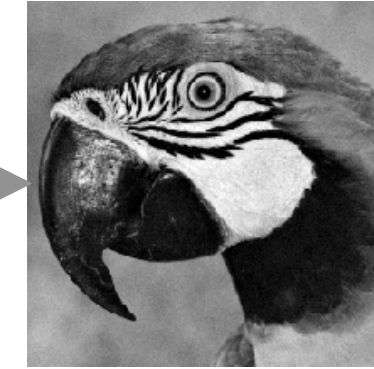
Diffusion models and denoising

noisy image $y = x + z$

$$z \sim \mathcal{N}(0, \sigma^2 \text{Id})$$

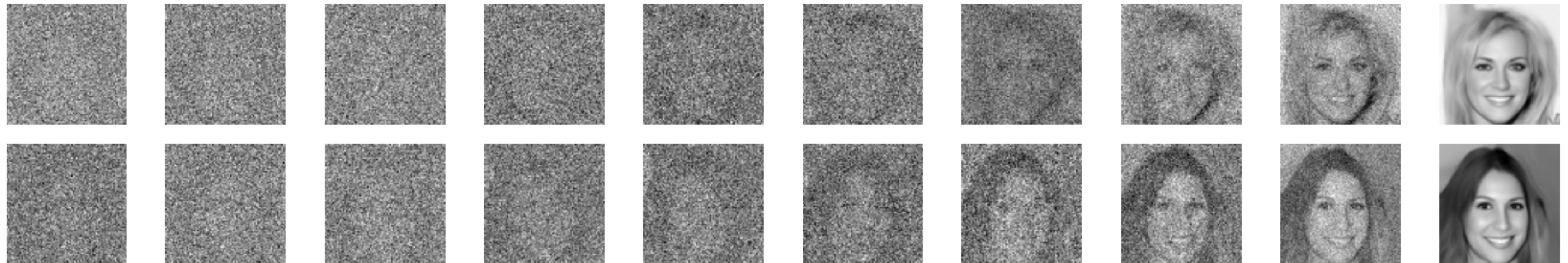


Deep Neural Network



$\hat{f}(y)$ denoised image

Denoiser is applied iteratively and partially

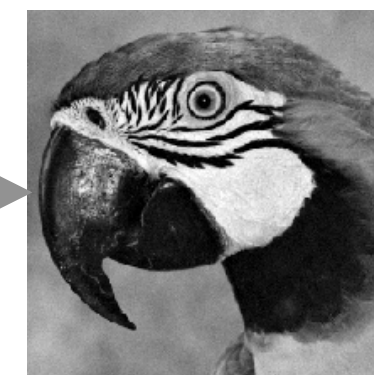
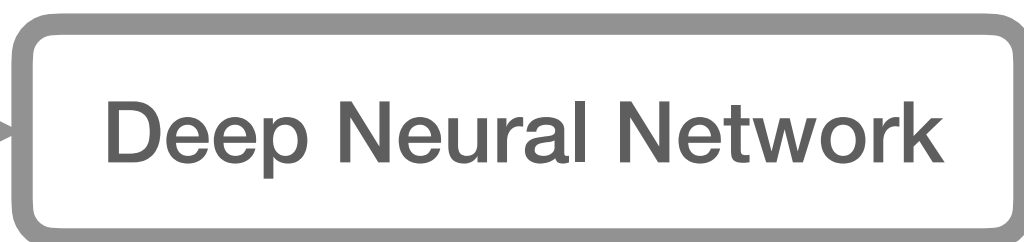
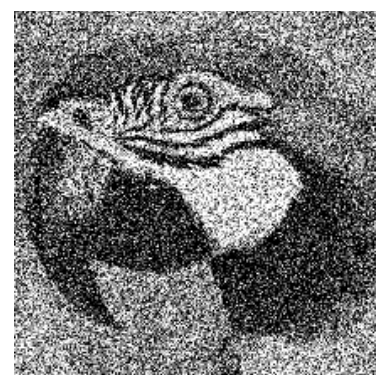


[Kadkhodaie & Simoncelli arXiv2020, NeurIPS2021]

Diffusion models and denoising

noisy image $y = x + z$

$$z \sim \mathcal{N}(0, \sigma^2 \text{Id})$$



$\hat{f}(y)$ denoised image

Minimize $\mathbb{E} \left[\|x - f(y)\|^2 \right]$

[Tweedie, via Robbins, 1956;
Miyasawa, 1961]

$$f^*(y) = \mathbb{E}_x[x|y] = \int x p^*(x|y) dx = f^*(y) = y + \sigma^2 \nabla \log p_\sigma^*(y)$$

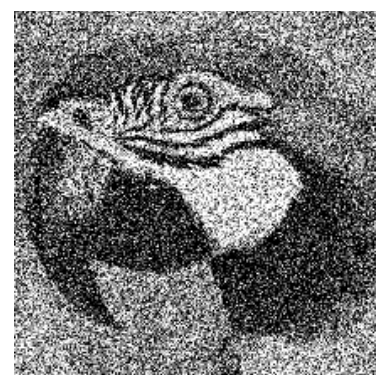
**optimal
denoiser**

mean of the posterior distribution

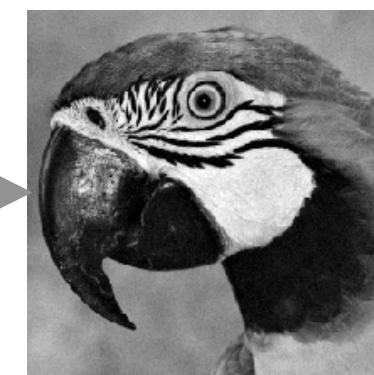
Diffusion models and denoising

noisy image $y = x + z$

$$z \sim \mathcal{N}(0, \sigma^2 \text{Id})$$



Deep Neural Network



$\hat{f}(y)$ denoised image

$$\text{Minimize } \mathbb{E} \left[\|x - f(y)\|^2 \right]$$

[Tweedie, via Robbins, 1956;
Miyasawa, 1961]

$$f^*(y) = \mathbb{E}_x[x|y] = \int x p^*(x|y) dx = \boxed{f^*(y) = y + \sigma^2 \nabla \log p_\sigma^*(y)} \quad \text{optimal denoiser}$$

mean of the posterior distribution

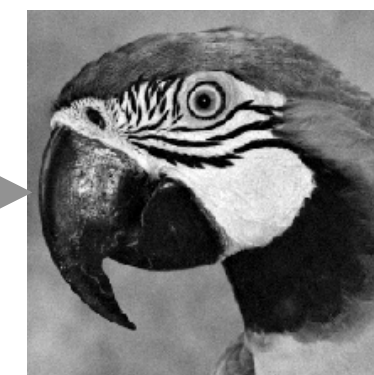
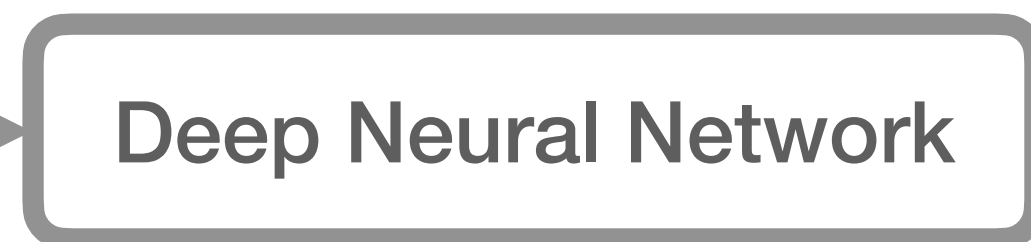
$$p_\sigma^*(y) = \int p(y|x) p^*(x) dx = \int g_\sigma(y-x) p^*(x) dx,$$

$p_\sigma(y)$ is a blurred (diffused) version of $p(x)$

Diffusion models and denoising

noisy image $y = x + z$

$$z \sim \mathcal{N}(0, \sigma^2 \text{Id})$$



$\hat{f}(y)$ denoised image

[Tweedie, via Robbins, 1956;
Miyasawa, 1961]

$$f^*(y) = y + \sigma^2 \nabla \log p_\sigma^*(y) \quad \text{optimal denoiser}$$

$$\hat{f}(y) = y + \sigma^2 \nabla \log \hat{p}_\sigma(y) \quad \text{learned denoiser}$$

$$\nabla_y \log \hat{p}_\sigma(y) \approx (\hat{f}(y) - y) / \sigma^2$$

Coarse-to-fine **gradient ascent**



Diffusion models and denoising

$$f^*(y) = y + \sigma^2 \nabla \log p_\sigma^*(y)$$

**optimal
denoiser**

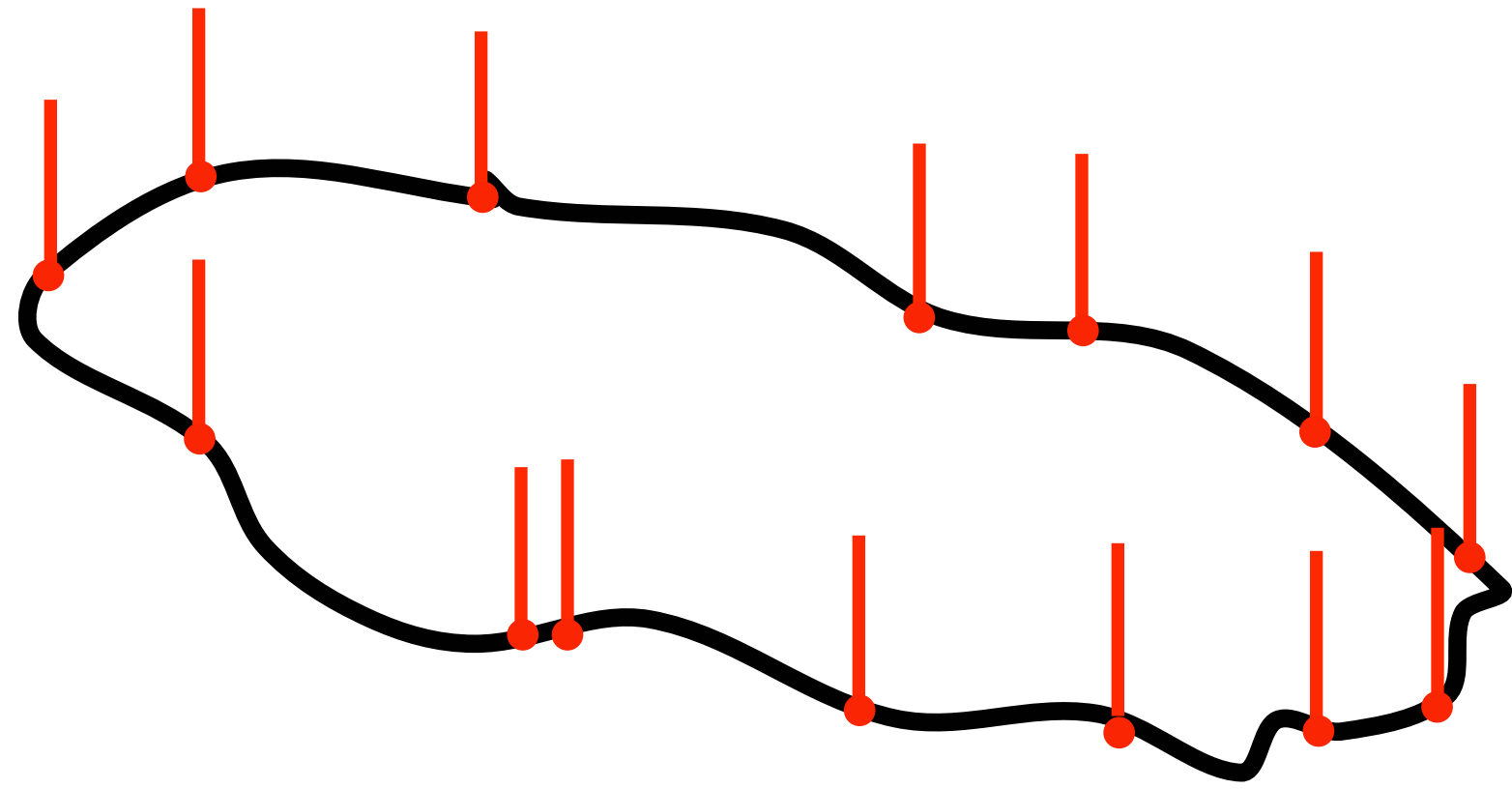
$$\hat{f}(y) = y + \sigma^2 \nabla \log \hat{p}_\sigma(y)$$

**learned
denoiser**

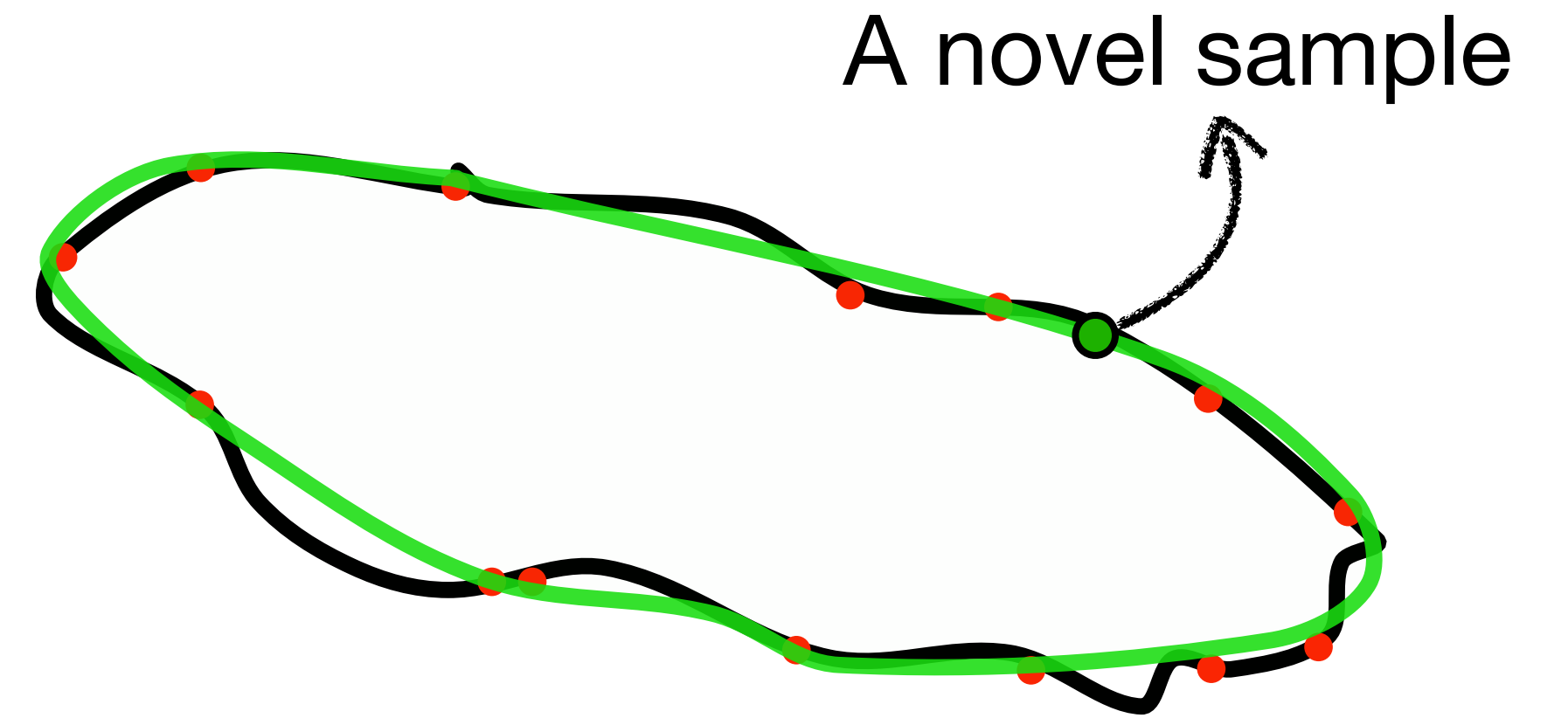


Samples from this

Diffusion models and denoising



$$\hat{p}_\sigma(y) \stackrel{?}{\approx} p_\sigma^*(y)$$



$$f^*(y) = y + \sigma^2 \nabla \log p_\sigma^*(y)$$

**optimal
denoiser**

$$\hat{f}(y) = y + \sigma^2 \nabla \log \hat{p}_\sigma(y)$$

**learned
denoiser**



Samples from this

Transition from memorization to generalization

Training set size:

1

10

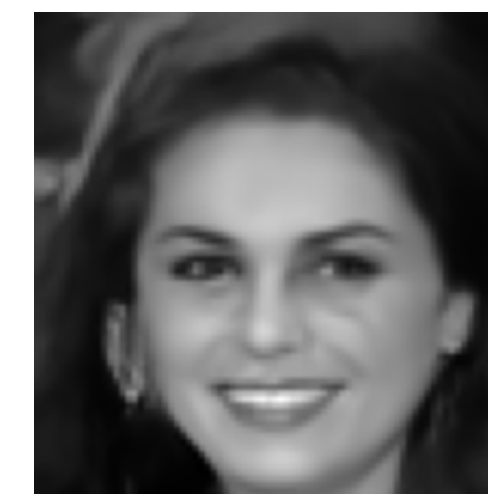
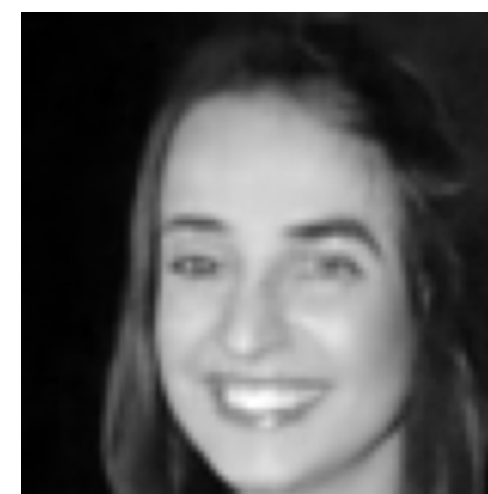
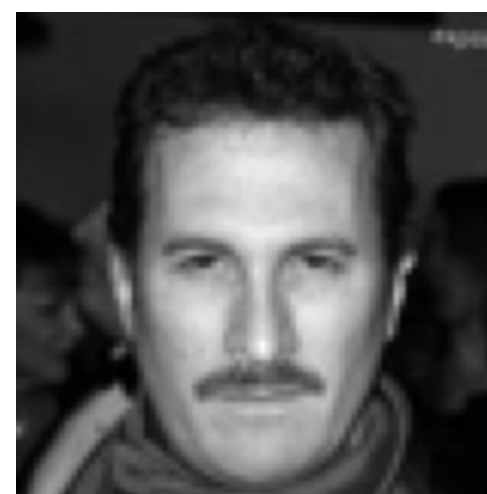
100

1,000

10,000

100,000

Samples,
model trained
on set A:



Transition from memorization to generalization

Training set size:

1

10

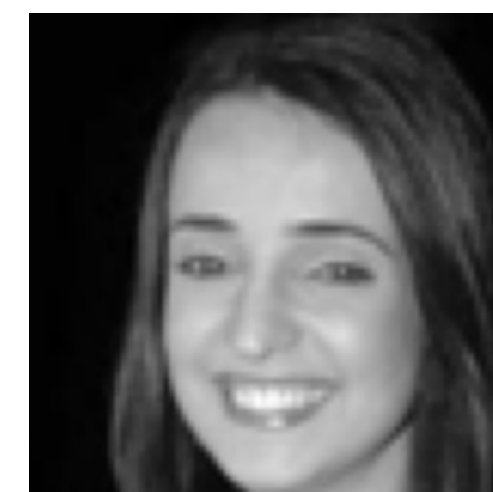
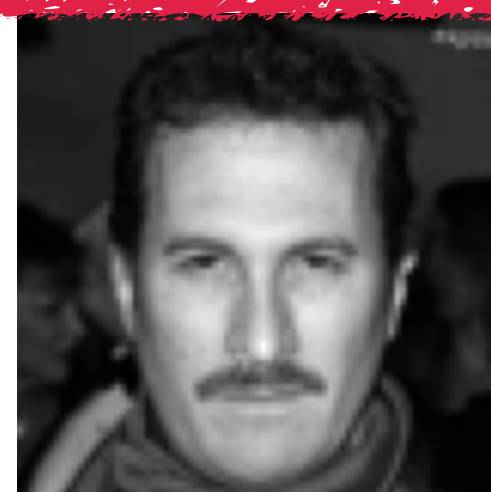
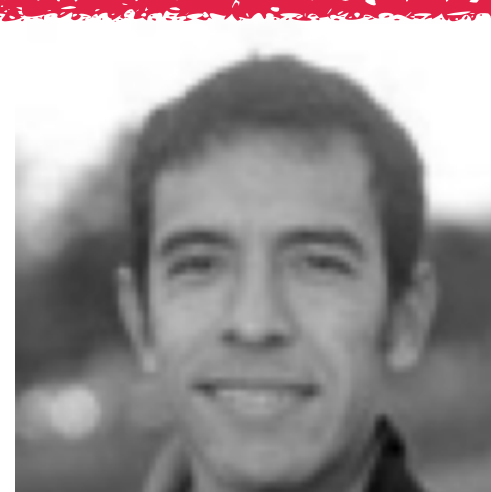
100

1,000

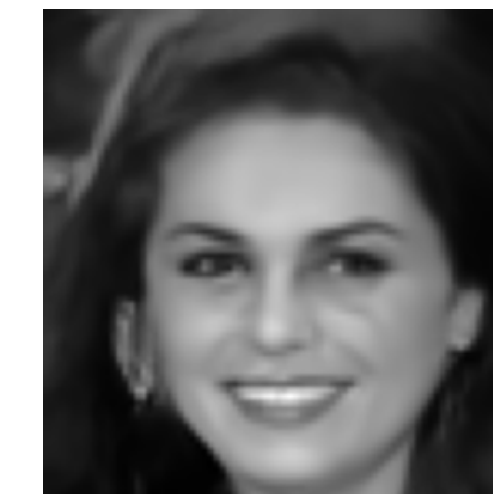
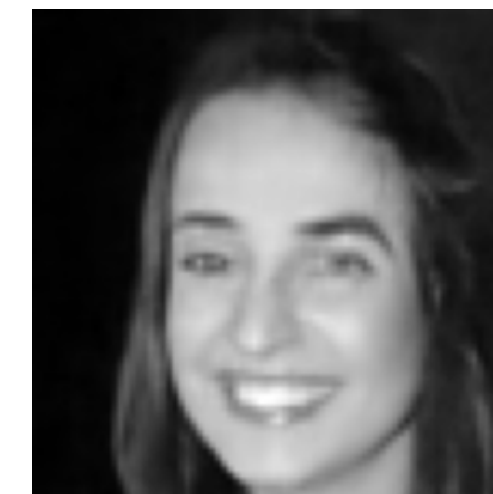
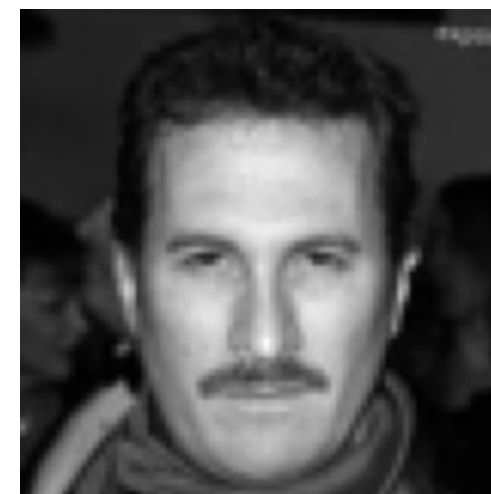
10,000

100,000

Closest training example from A:



Samples, model trained on set A:



Transition from memorization to generalization

Training set size:

1

10

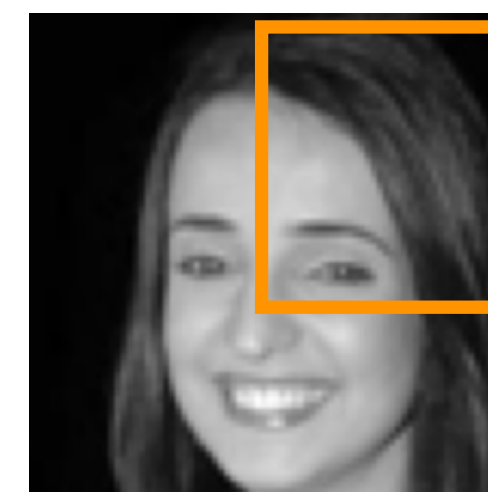
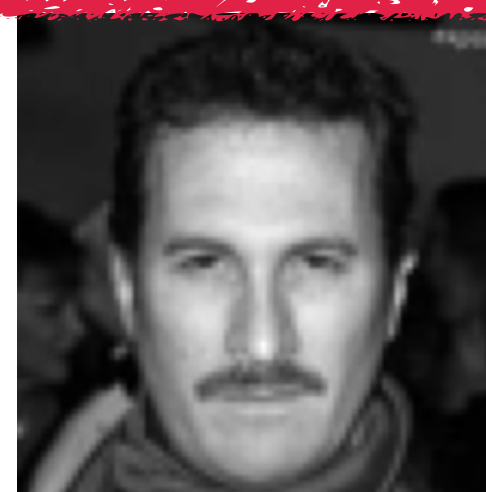
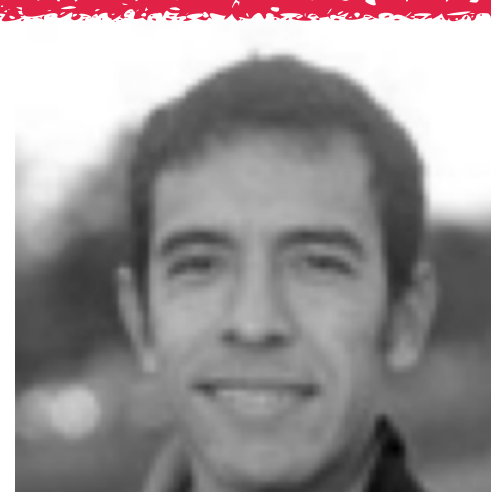
100

1,000

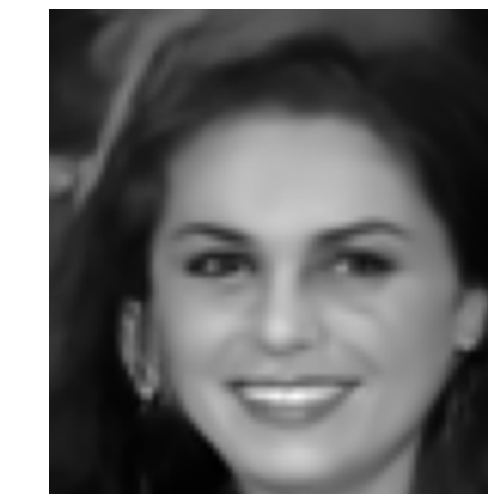
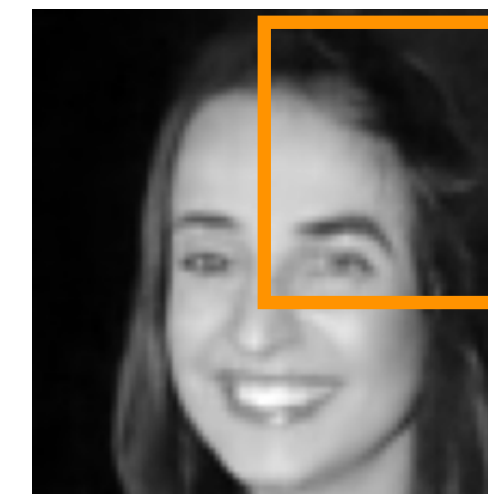
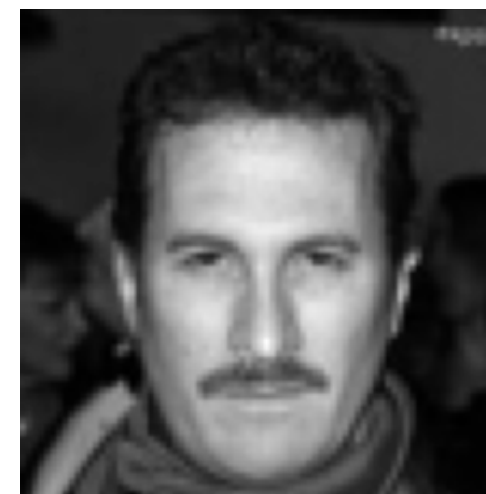
10,000

100,000

Closest training example from A:



Samples, model trained on set A:



Transition from memorization to generalization

Training set size:

1

10

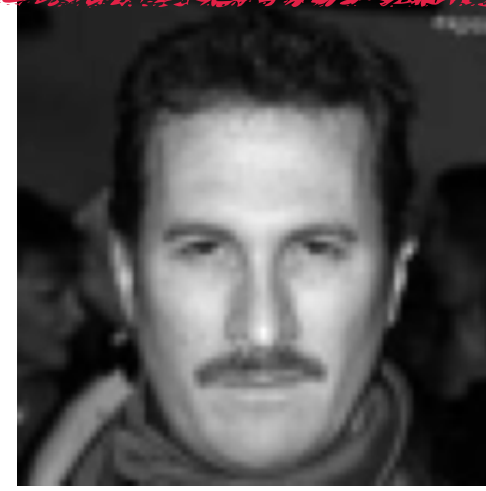
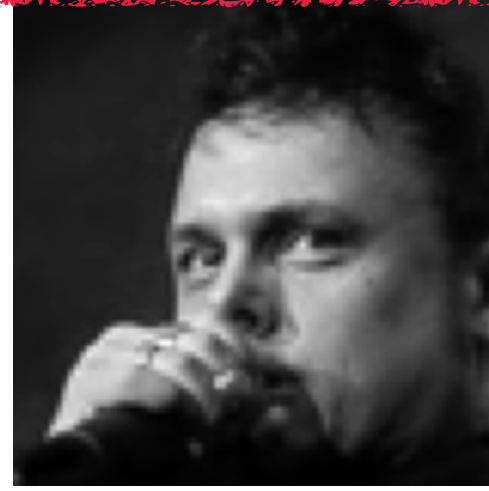
100

1,000

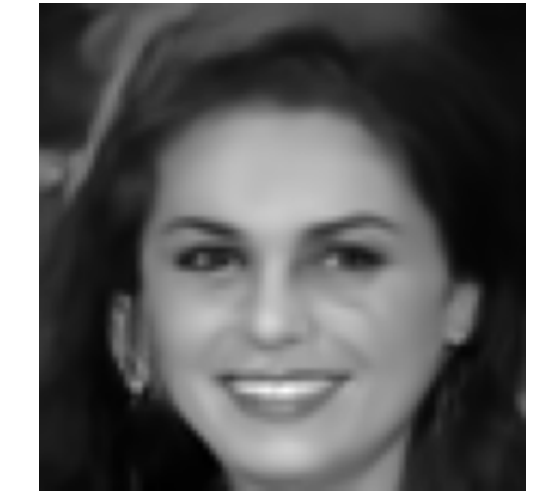
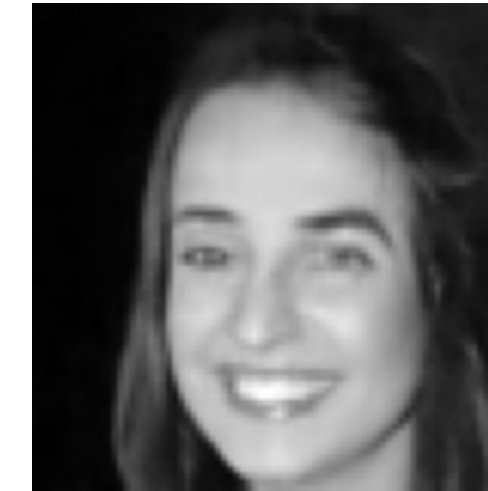
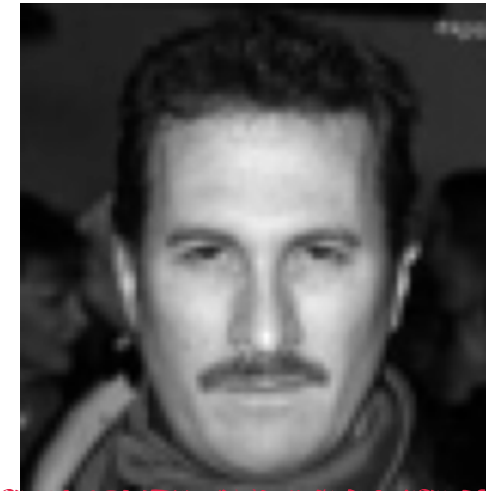
10,000

100,000

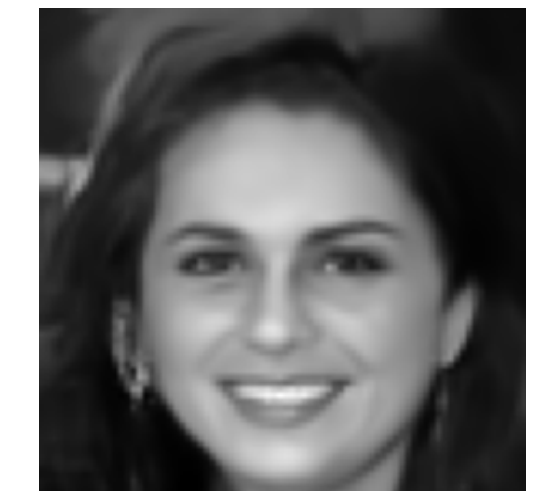
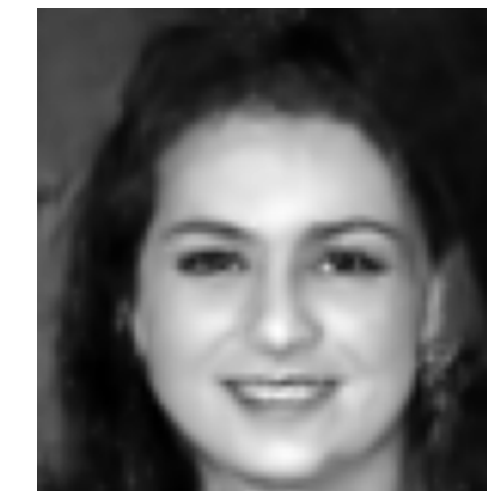
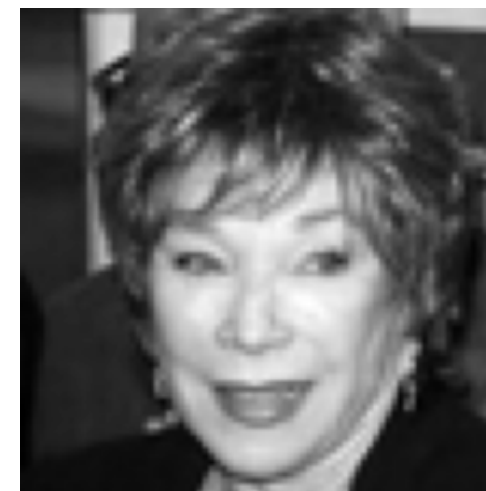
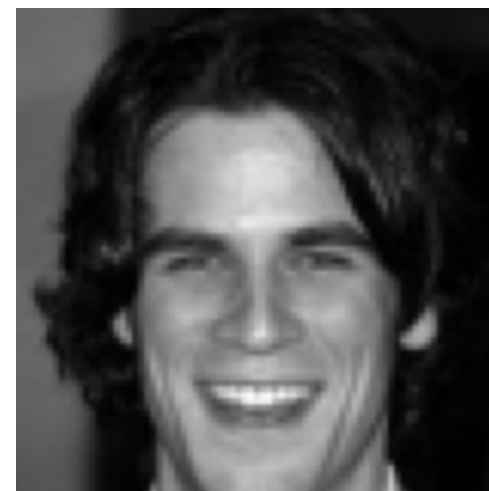
Closest training example from A:



Samples, model trained on set A:



Samples, model trained on set B (same seed):



Transition from memorization to generalization

Training set size:

1

10

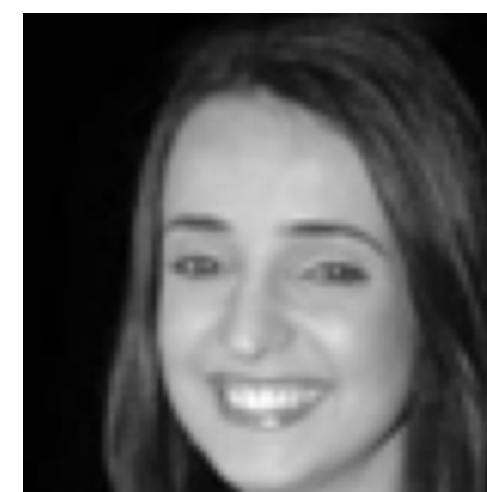
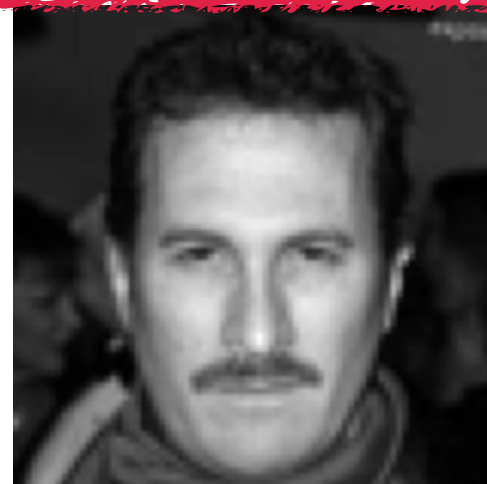
100

1,000

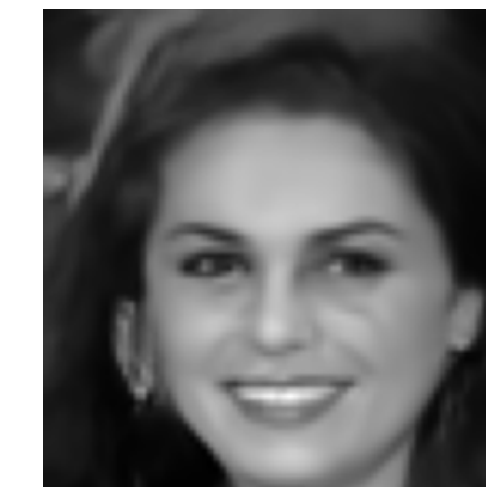
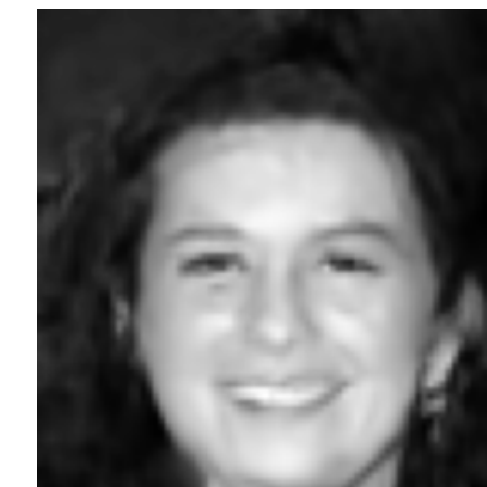
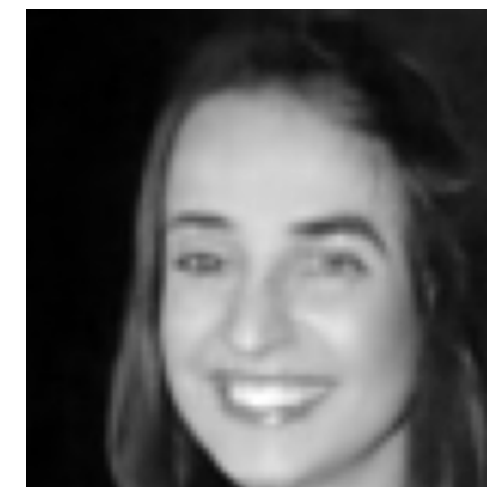
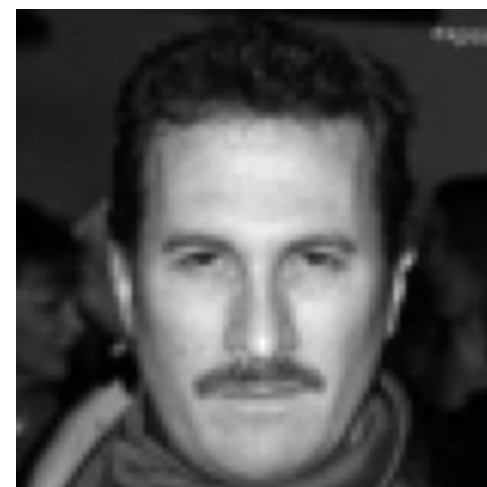
10,000

100,000

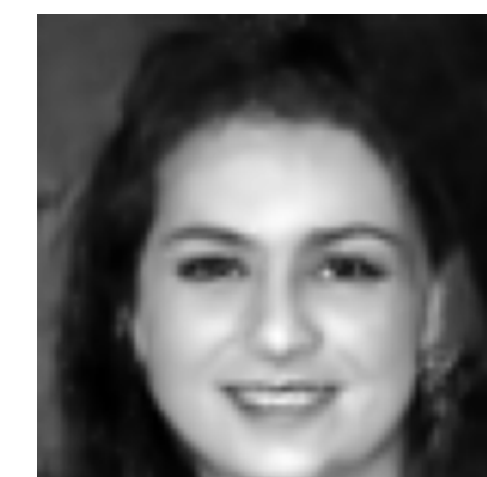
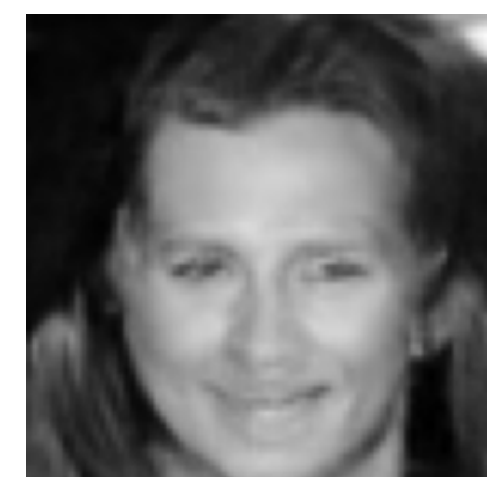
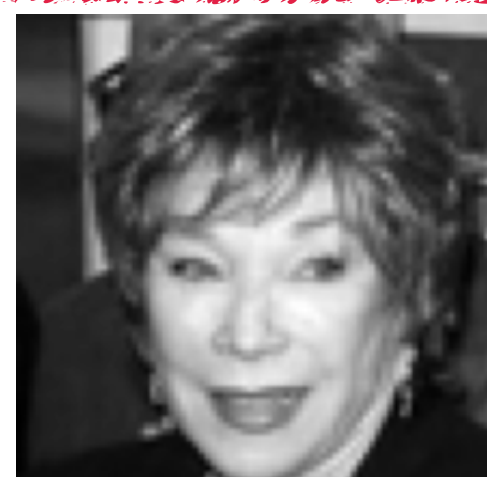
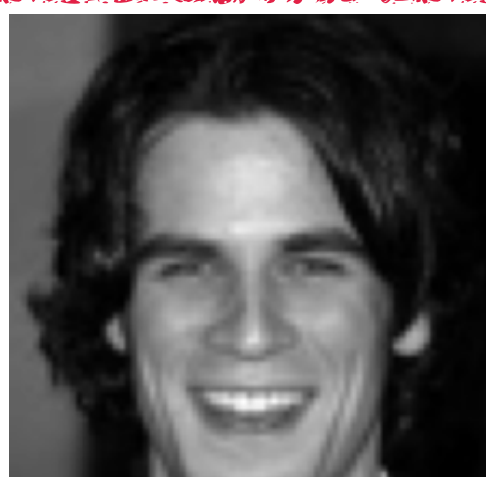
Closest training example from A:



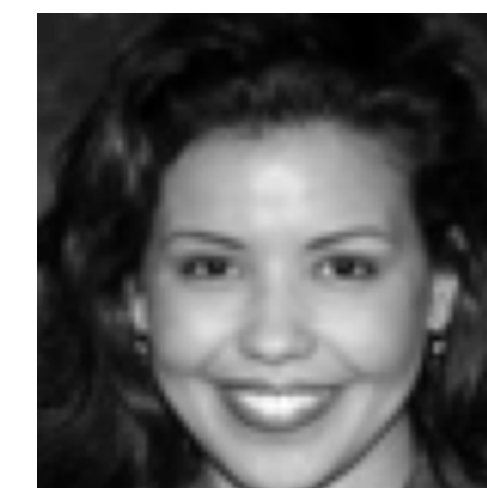
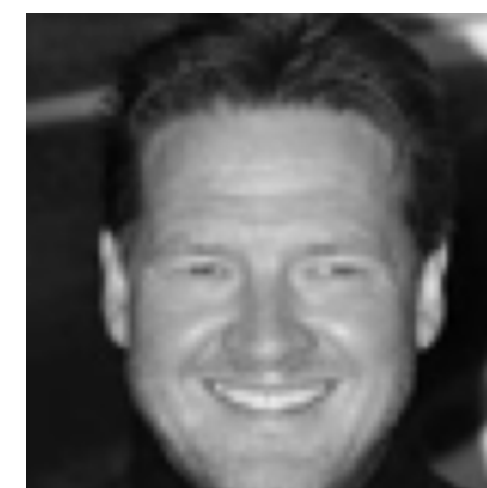
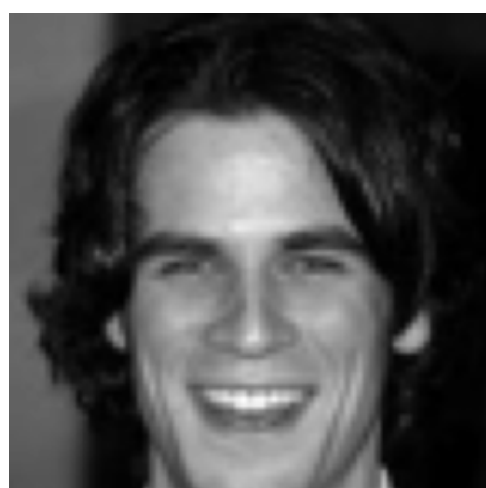
Samples, model trained on set A:



Samples, model trained on set B (same seed):



Closest training example from B:



Transition from memorization to generalization

Training set size:

1

10

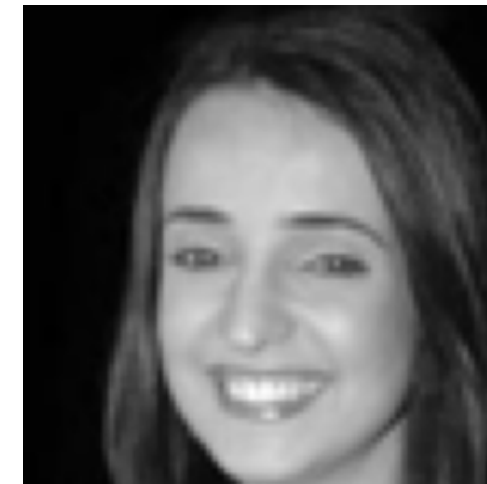
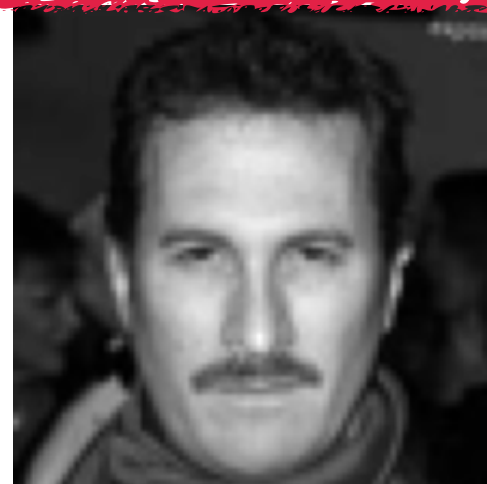
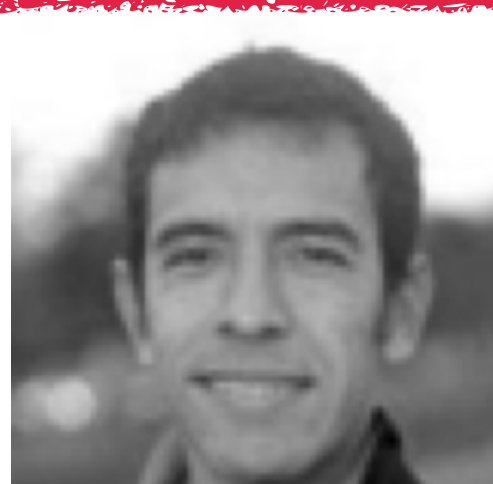
100

1,000

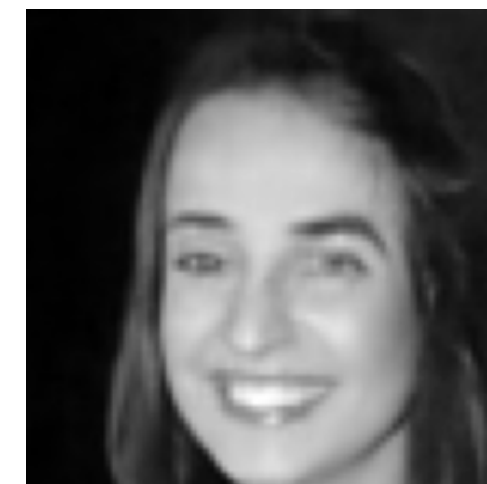
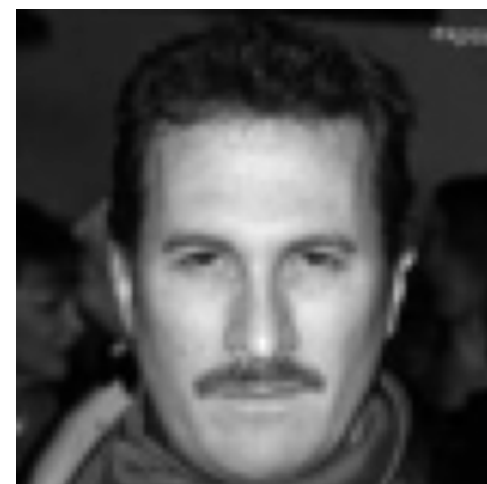
10,000

100,000

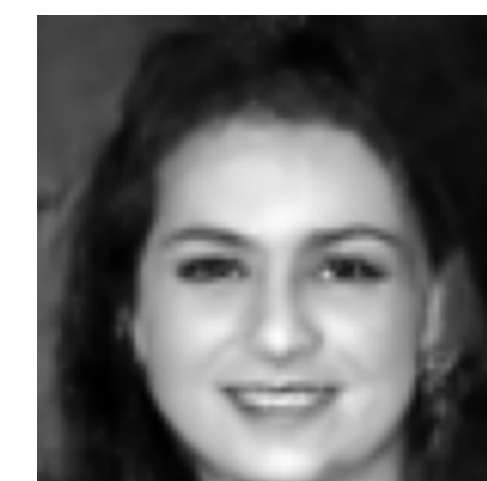
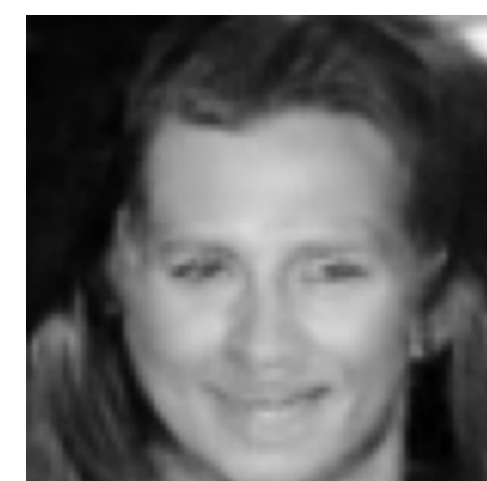
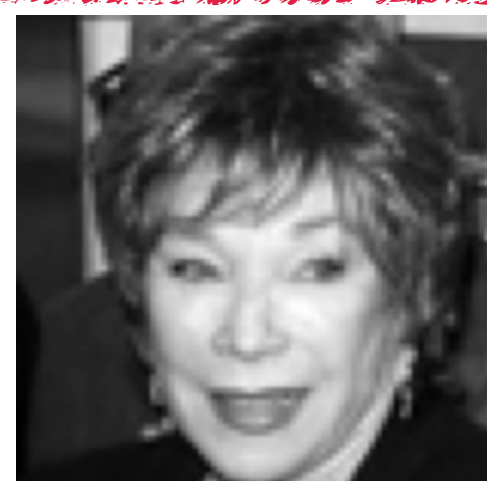
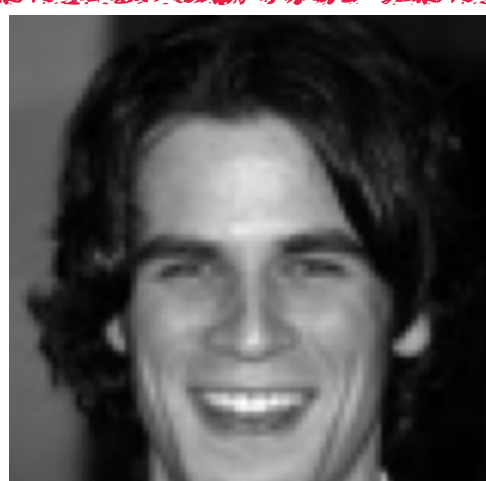
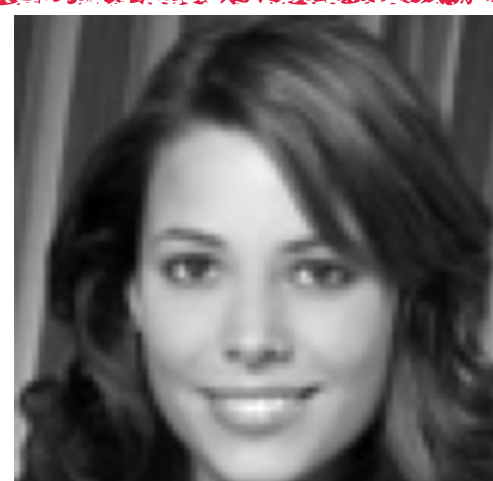
Closest training example from A:



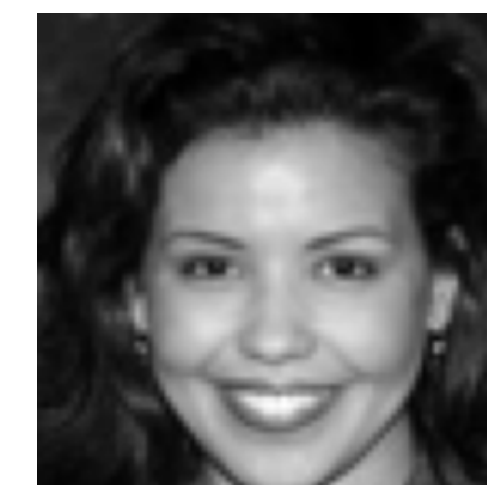
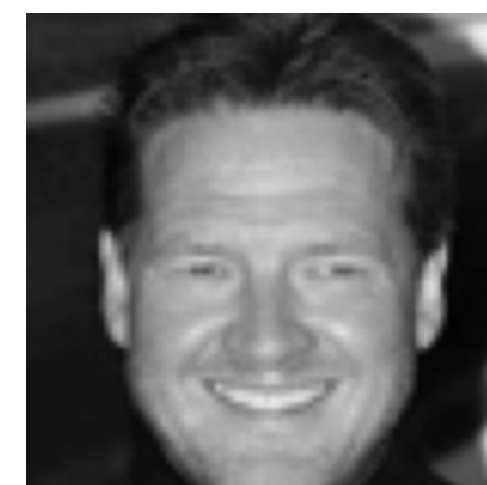
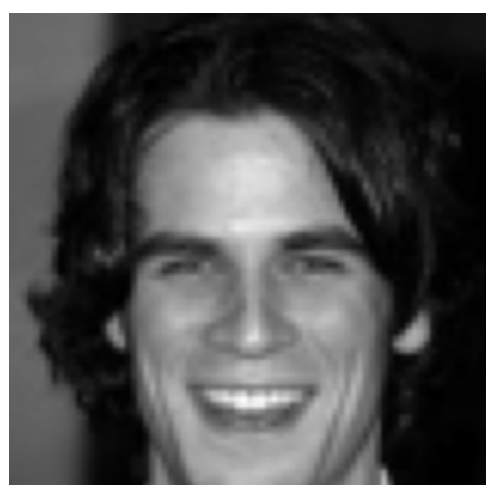
Samples, model trained on set A:



Samples, model trained on set B (same seed):

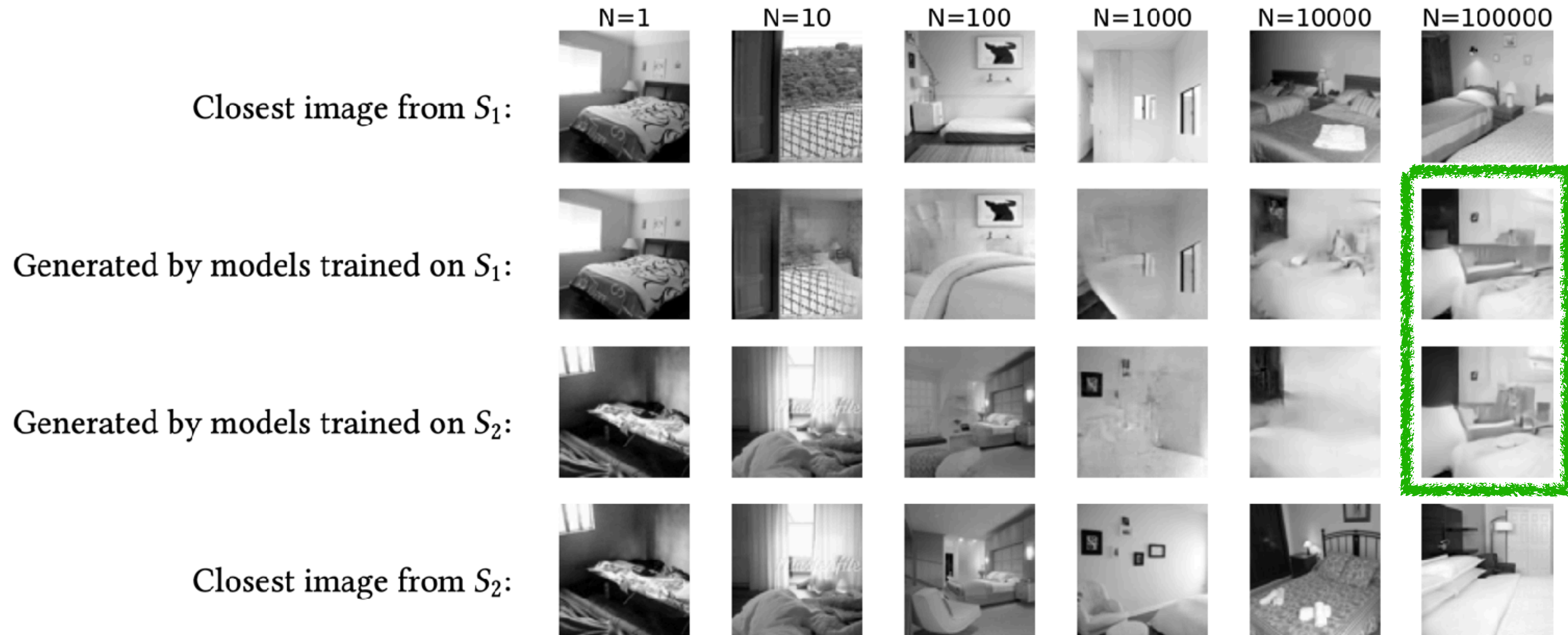


Closest training example from B:

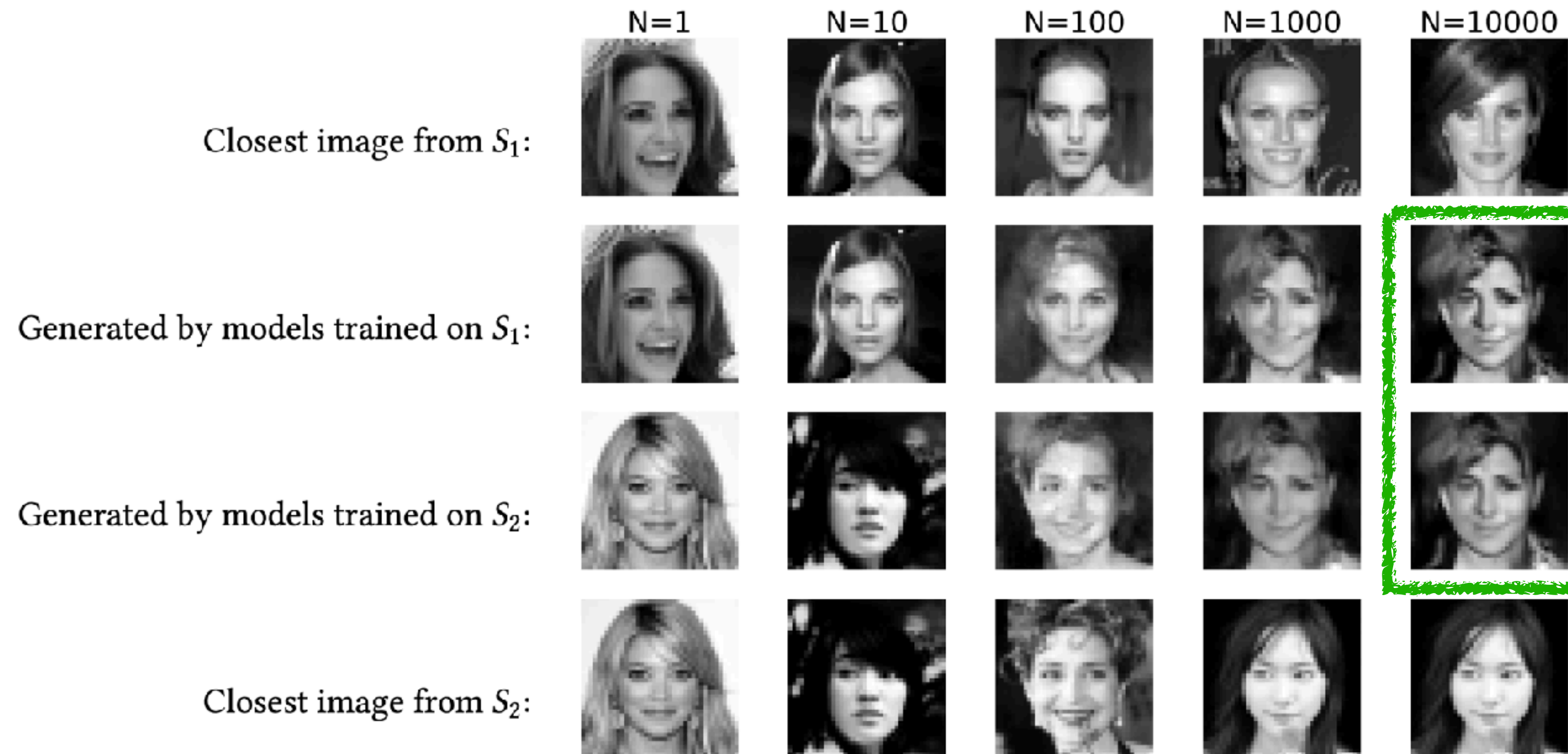


Same initialization

Strong **generalization** in LSUN bedroom dataset



Strong **generalization** in BF-CNN architecture



How do diffusion models generalize?

What are inductive biases of the denoiser?

Denoising as shrinkage in a basis

Classical framework for denoising:

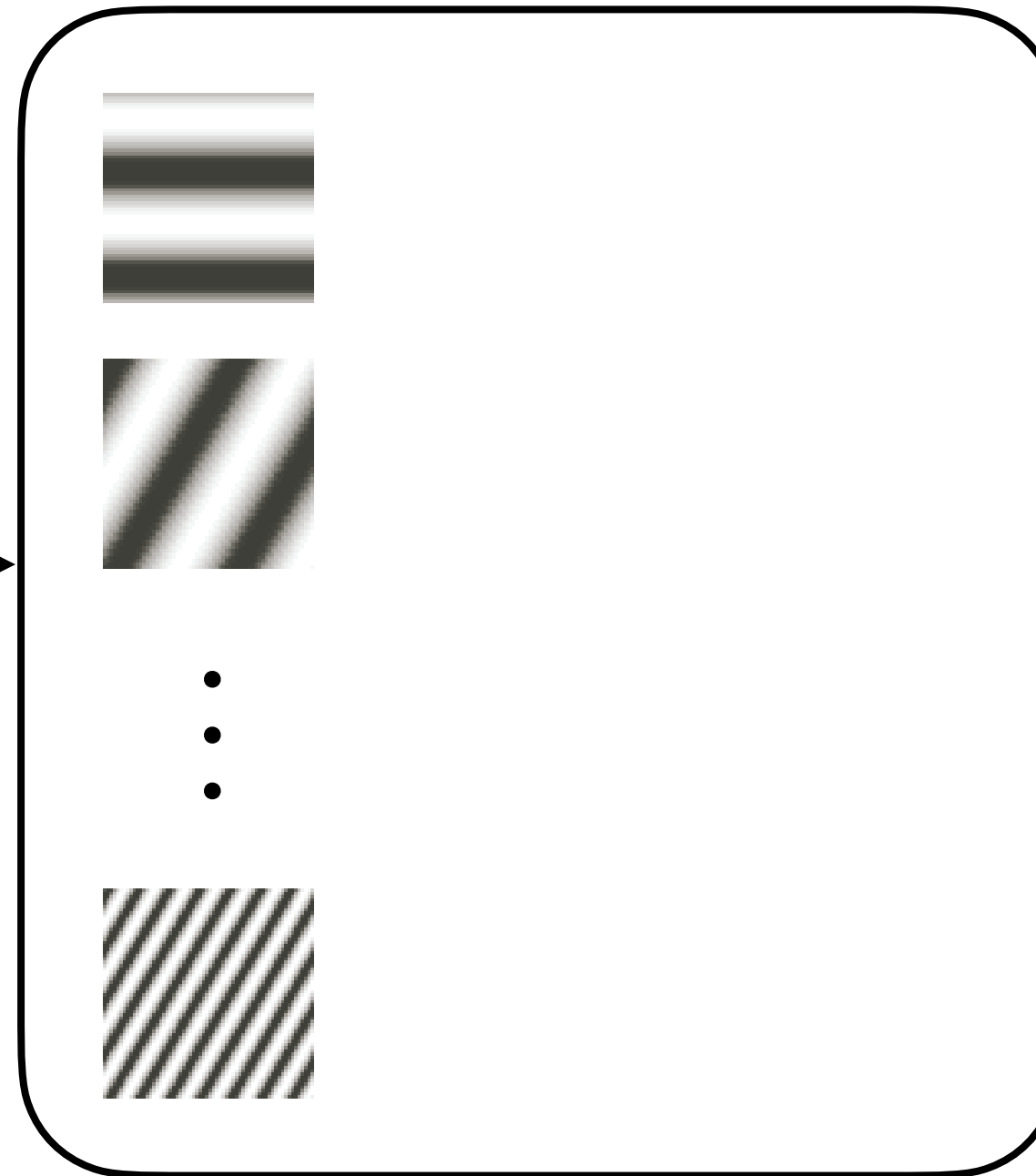
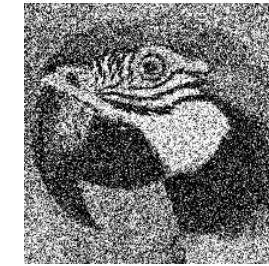
1. Transform the noisy image to a basis where noise and signal are separable
2. Suppress the noise (shrinkage)
3. Transform back to pixel domain

Denoising as shrinkage in a basis

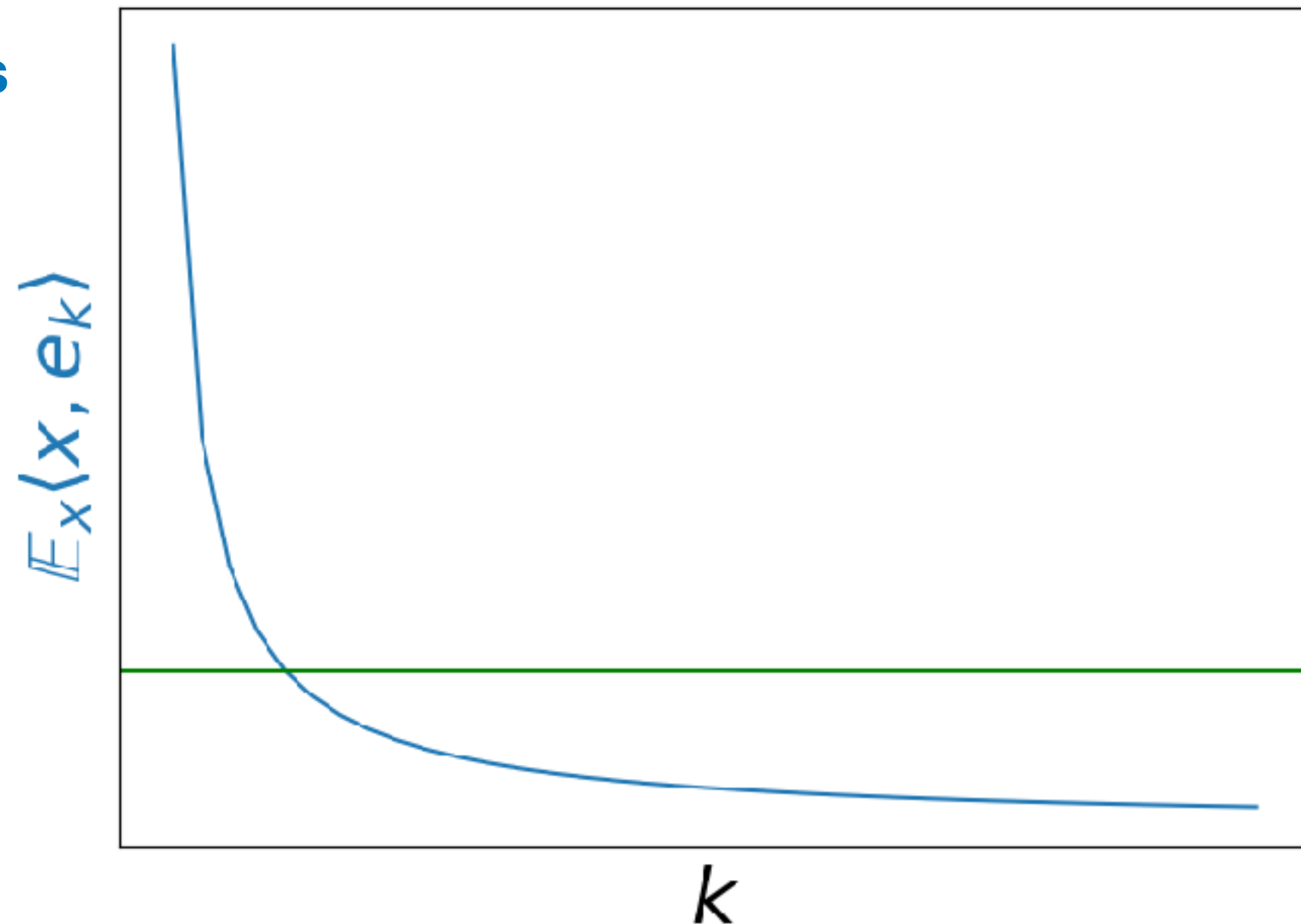
Fixed basis, fixed shrinkage

$$f(y) = \sum_k \lambda_k \langle y, e_k \rangle e_k$$

Fourier basis



Signal falls
as $\frac{1}{k}$ on
average



$E_z \langle z, e_k \rangle$ same power in all frequencies

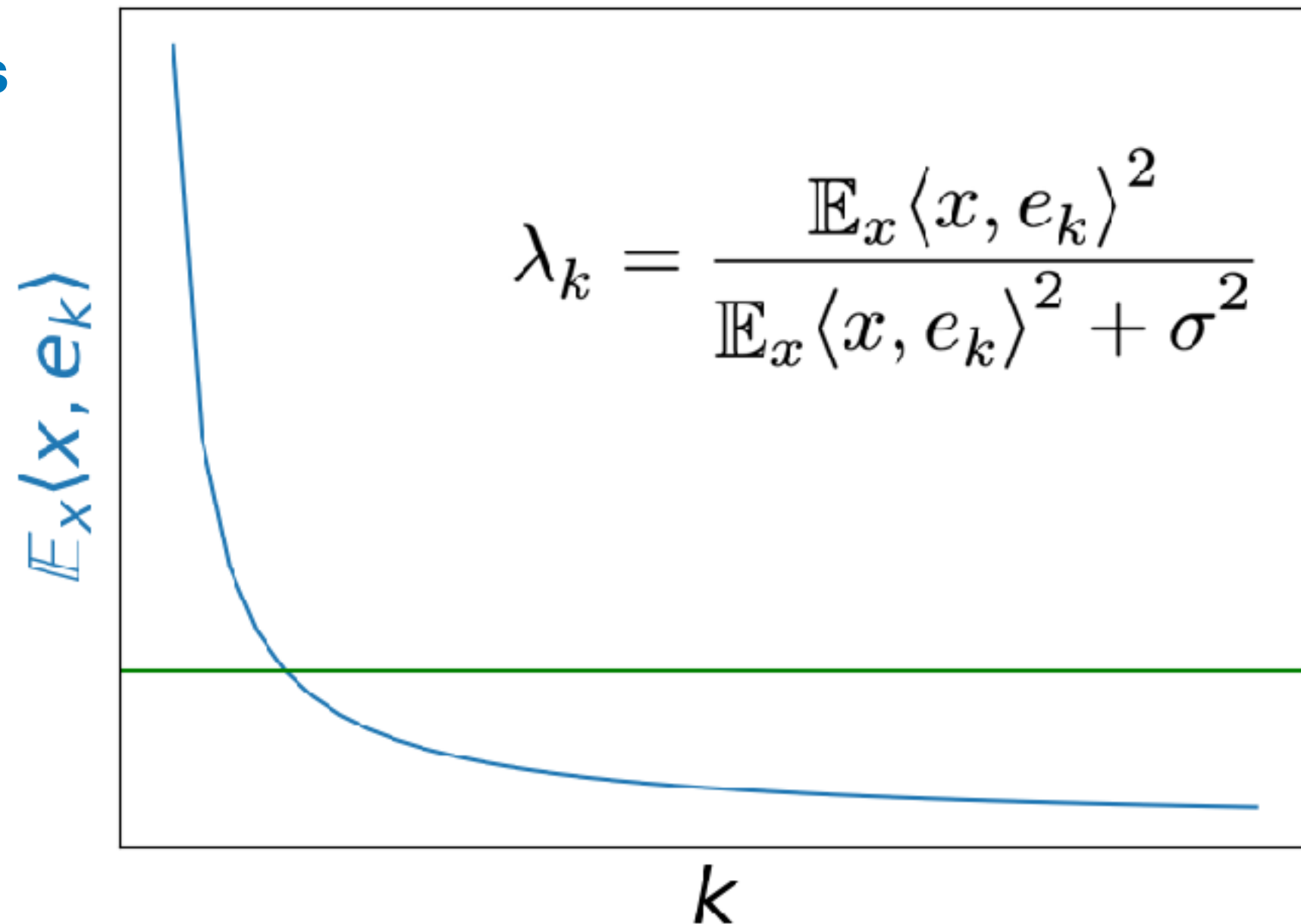
Denoising as shrinkage in a basis

Fixed basis, fixed shrinkage

$$f(y) = \sum_k \lambda_k \langle y, e_k \rangle e_k$$

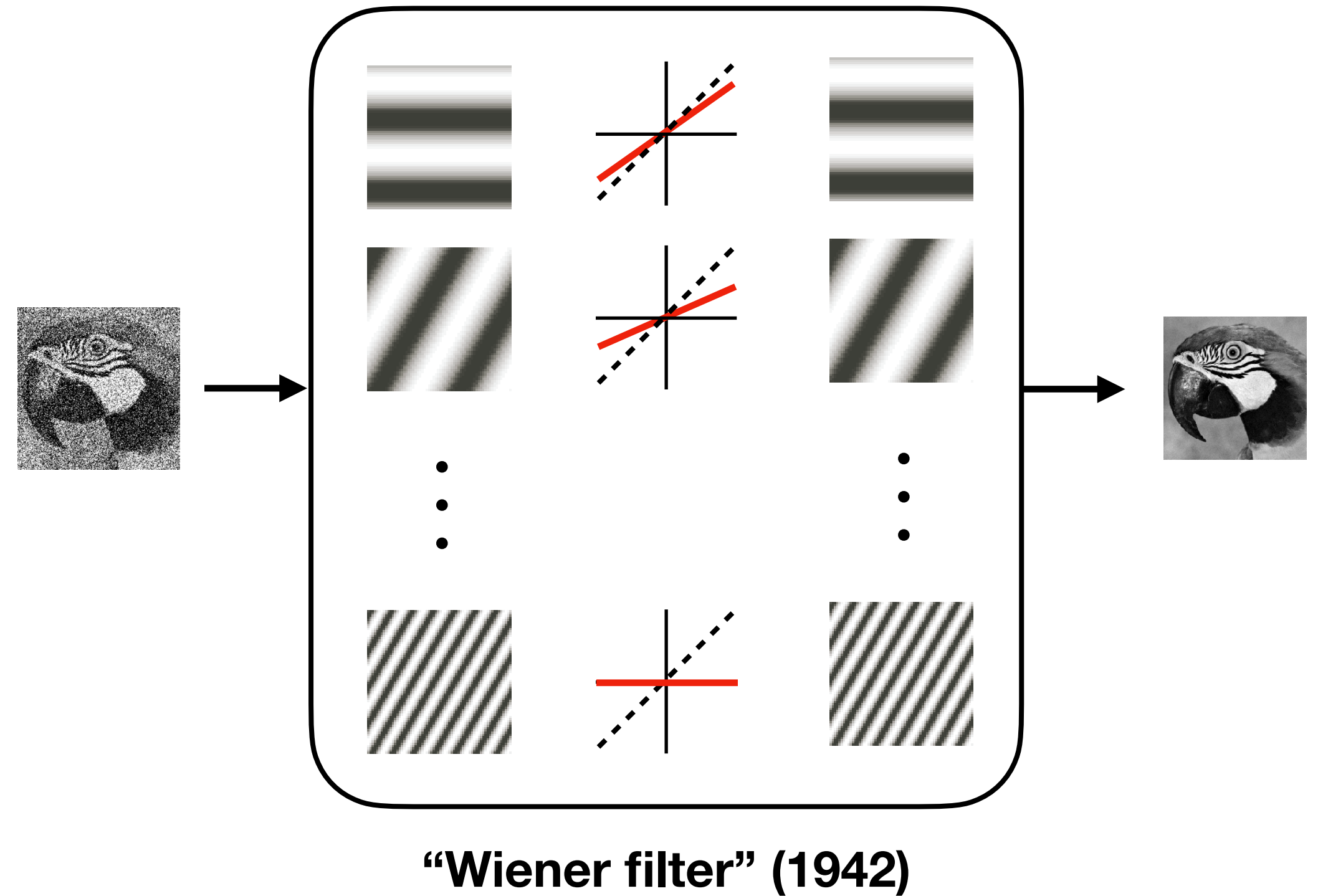
Fourier basis

Signal falls
as $\frac{1}{k}$ on
average



$$\lambda_k = \frac{\mathbb{E}_x \langle x, e_k \rangle^2}{\mathbb{E}_x \langle x, e_k \rangle^2 + \sigma^2}$$

$\mathbb{E}_z \langle z, e_k \rangle^2$ same power in all frequencies

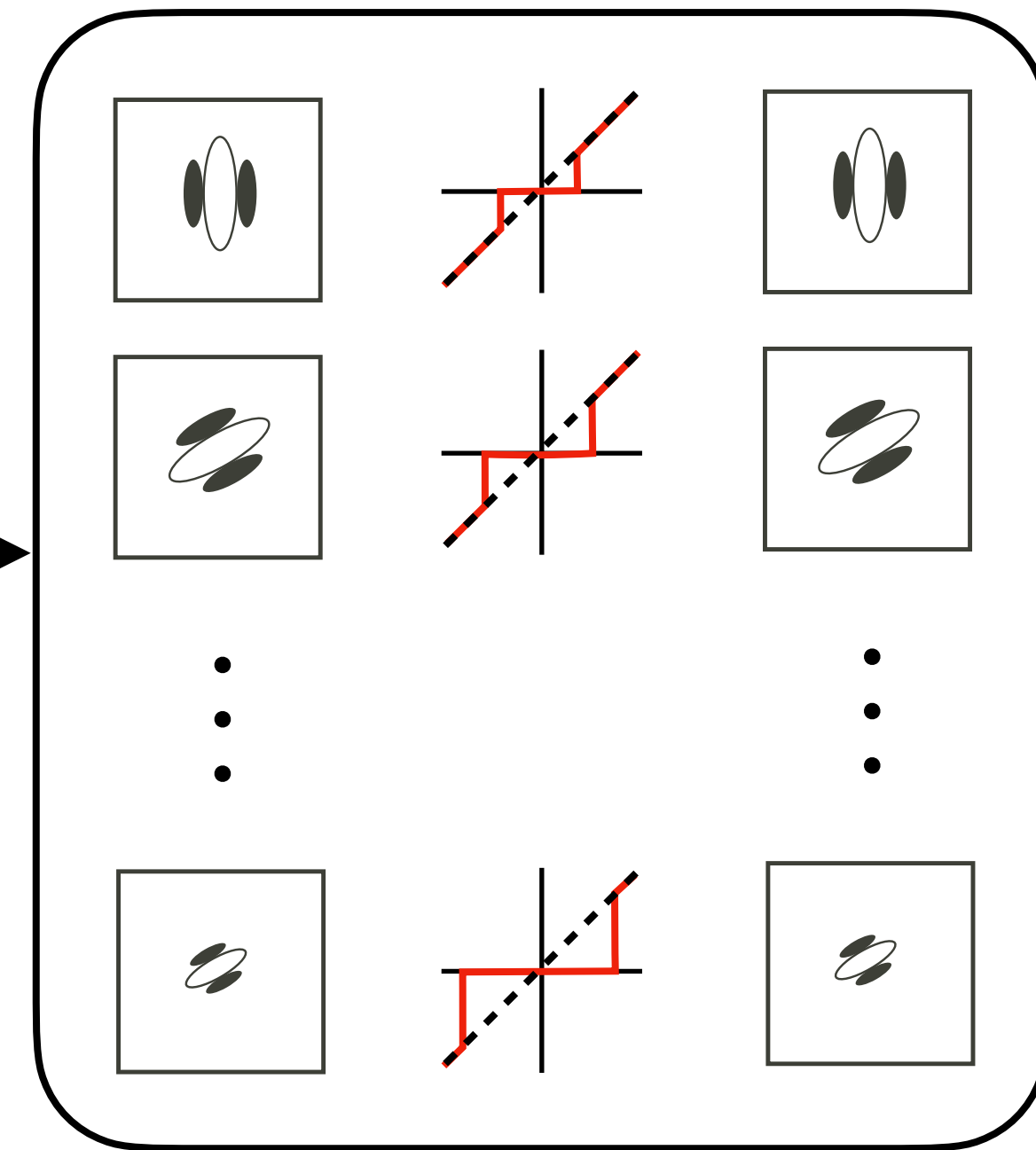


Denoising as shrinkage in a basis

Fixed basis, adaptive shrinkage

$$f(y) = \sum_k \lambda_k(y) \langle y, e_k \rangle e_k$$

Wavelet basis



“wavelet thresholding”

[Donoho & Johnstone 94]

- ☑ Coefficients fall faster in wavelet basis. More compact representation of signal.
Easier separation between noise and signal with sparse signal

Denoising as shrinkage in a basis

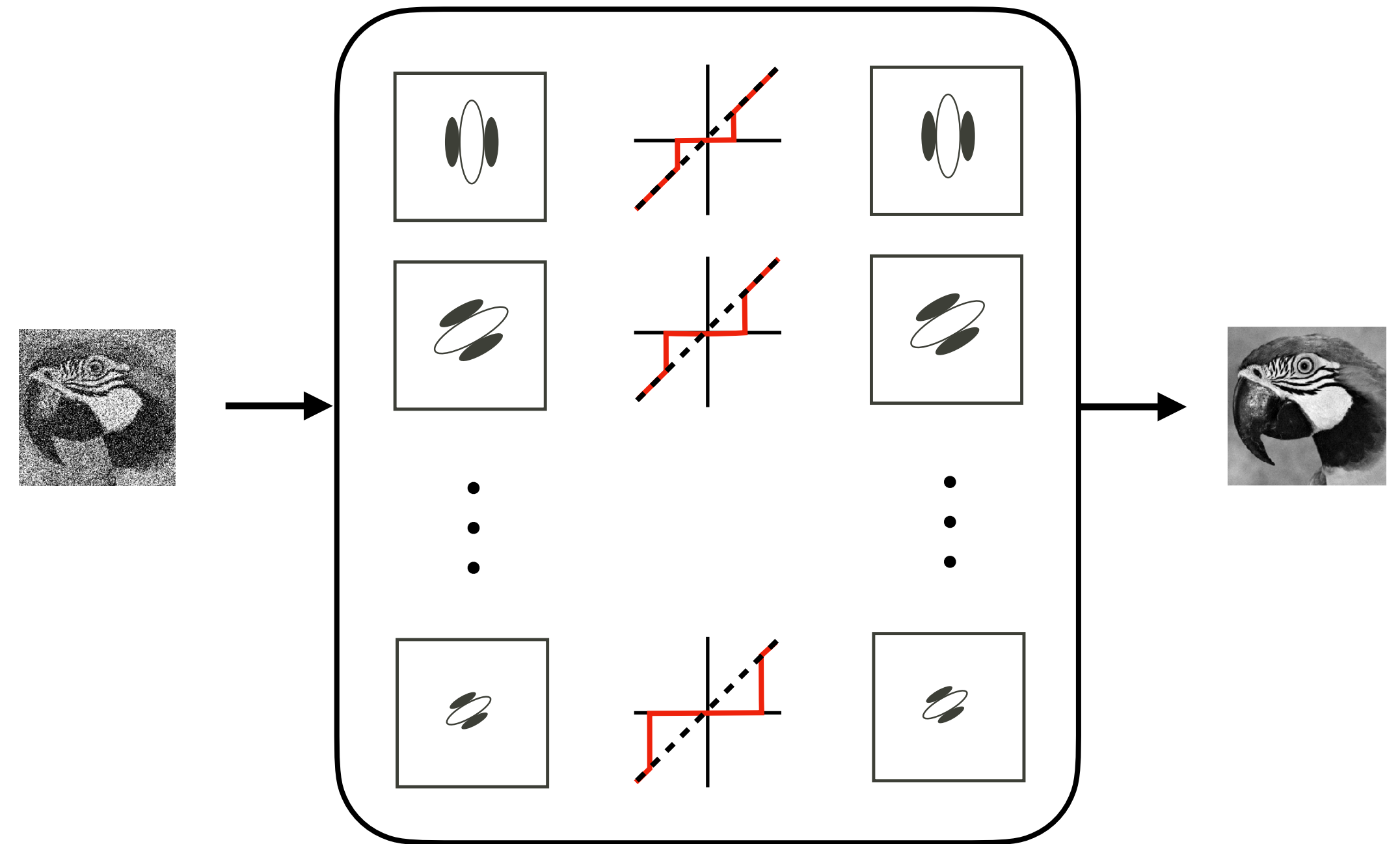
Fixed basis, adaptive shrinkage

$$f(y) = \sum_k \lambda_k(y) \langle y, e_k \rangle e_k$$

$\lambda_k(y)$ \downarrow Adaptive thresholding

$\langle y, e_k \rangle$ \downarrow Wavelet basis

$$\lambda_k(y) = \begin{cases} 1 & |\langle y, e_k \rangle| > \alpha\sigma \\ 0 & \text{Otherwise} \end{cases}$$



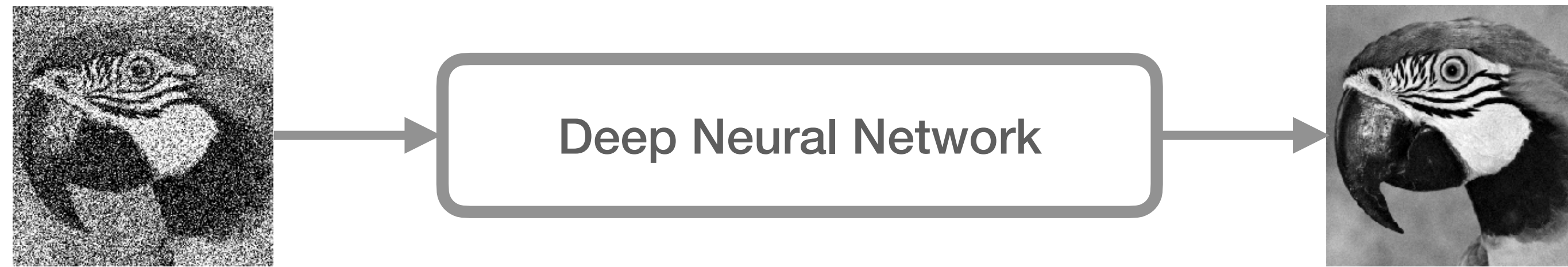
“wavelet thresholding”

[Donoho & Johnstone 94]

- ☑ Coefficients fall faster in wavelet basis. More compact representation of signal.
Easier separation between noise and signal with sparse signal

Denoising as shrinkage in a basis

Adaptive basis, adaptive shrinkage



Locally linear
function:

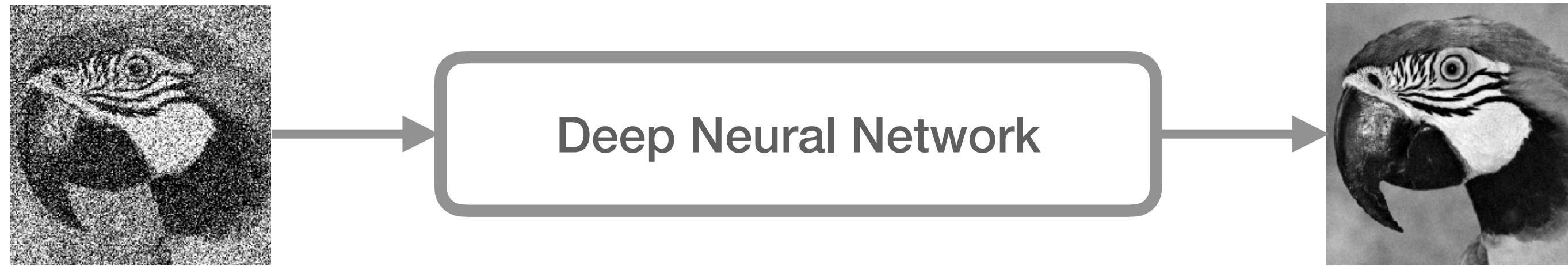
$$\hat{f}(y) = \underbrace{\nabla \hat{f}(y)}_y$$

Jacobian w.r.t. Input y

Nearly symmetric

Denoising as shrinkage in a basis

Adaptive basis, adaptive shrinkage



Eigen decomposition of Jacobian

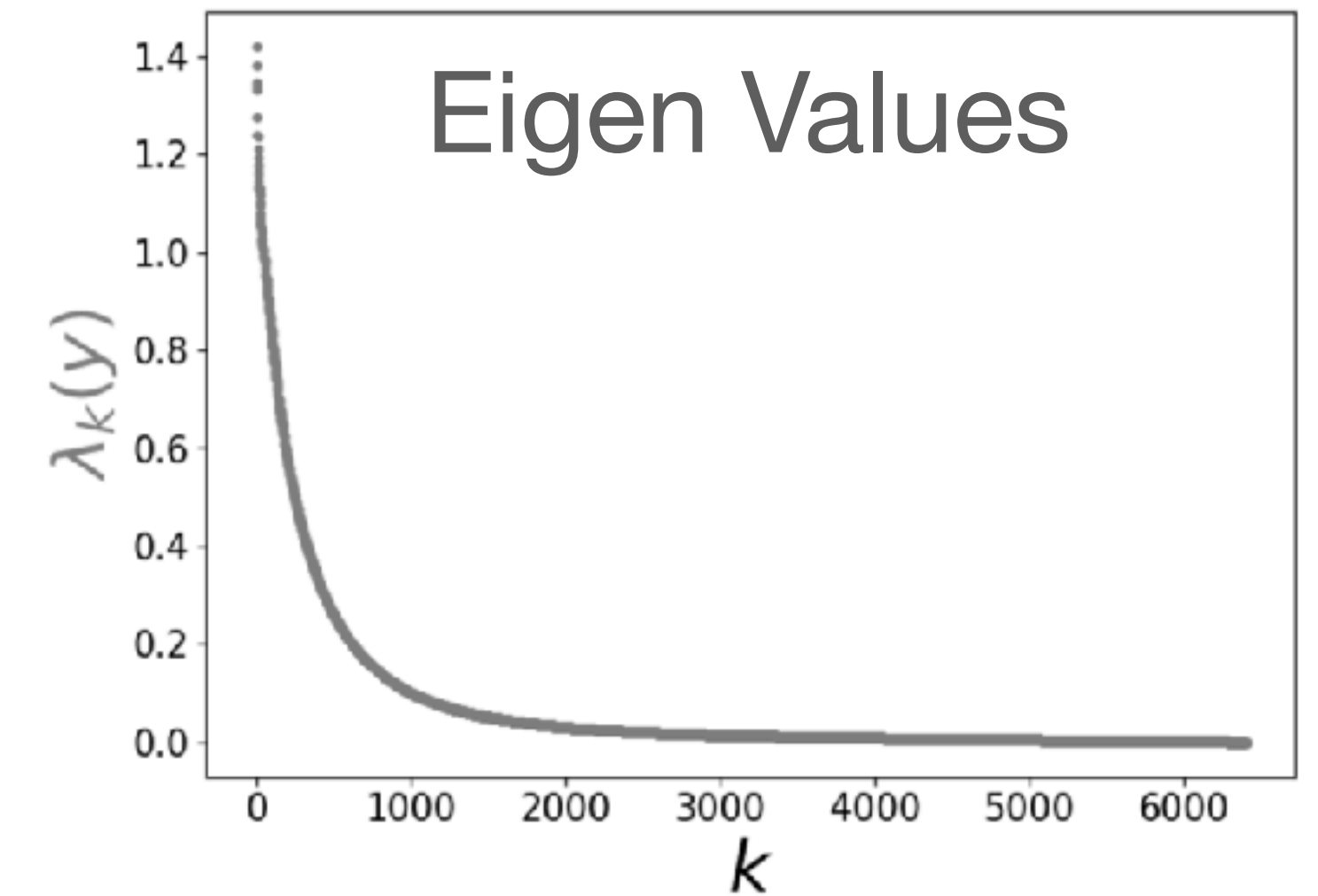
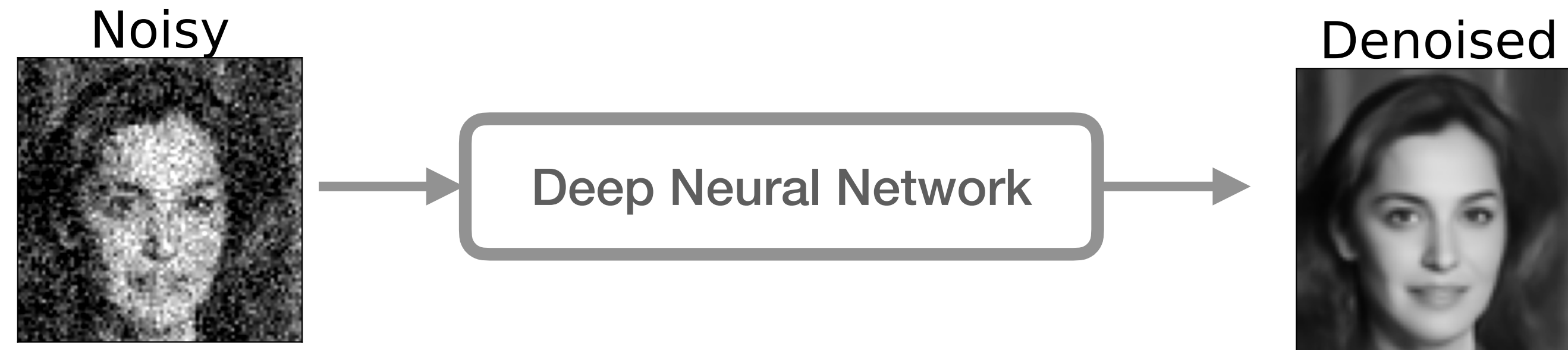
Locally linear function: $\hat{f}(y) = \nabla \hat{f}(y) y = \sum_k \lambda_k(y) \langle y, e_k(y) \rangle e_k(y)$

Shrinkage factors Eigen basis

The equation shows the locally linear function $\hat{f}(y)$ as a sum over k of shrinkage factors $\lambda_k(y)$ multiplied by the inner product of the input y and the eigen basis vectors $e_k(y)$, which are then multiplied by the eigen basis vectors themselves. Arrows point from $\lambda_k(y)$ to 'Shrinkage factors' and from $e_k(y)$ to 'Eigen basis'.

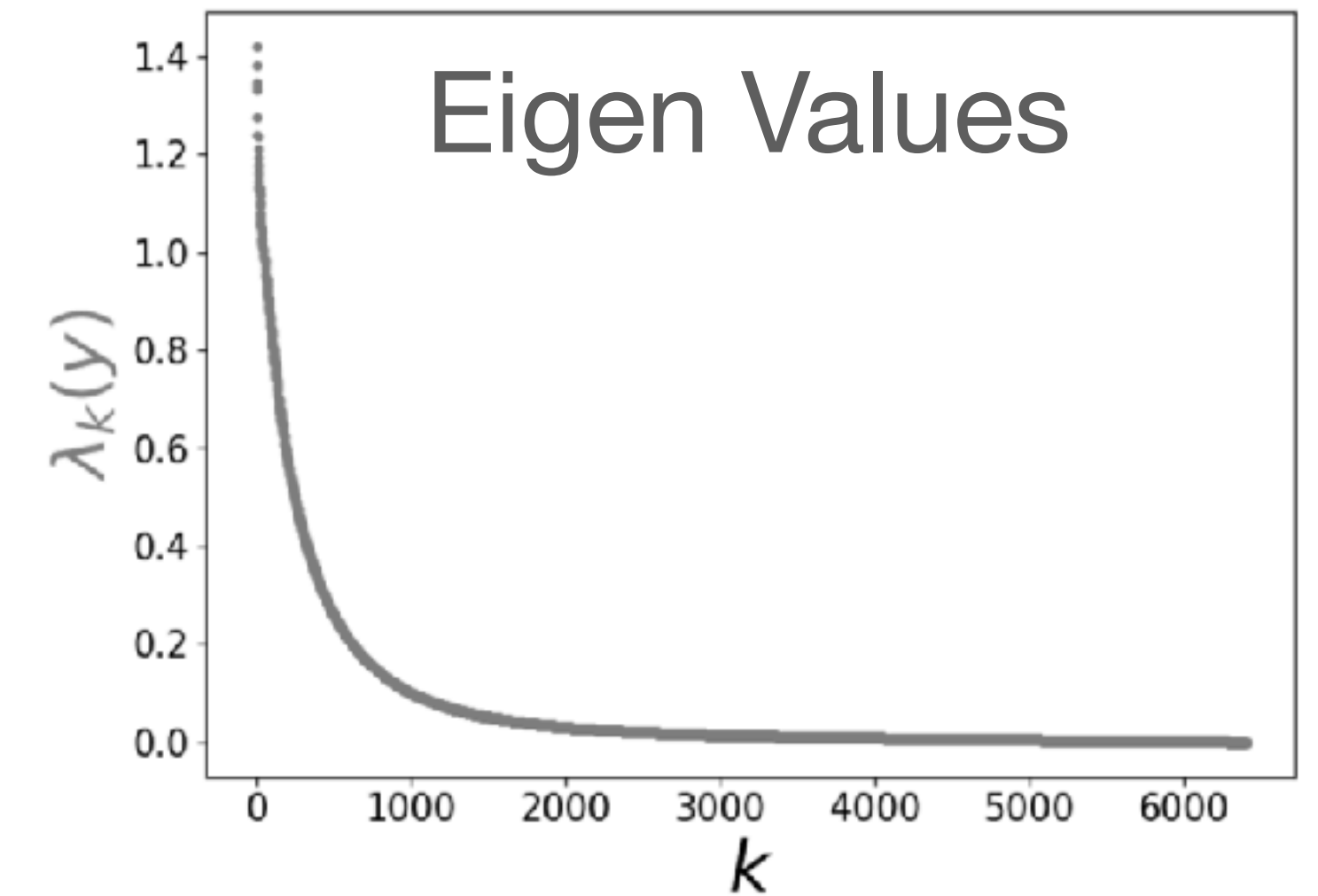
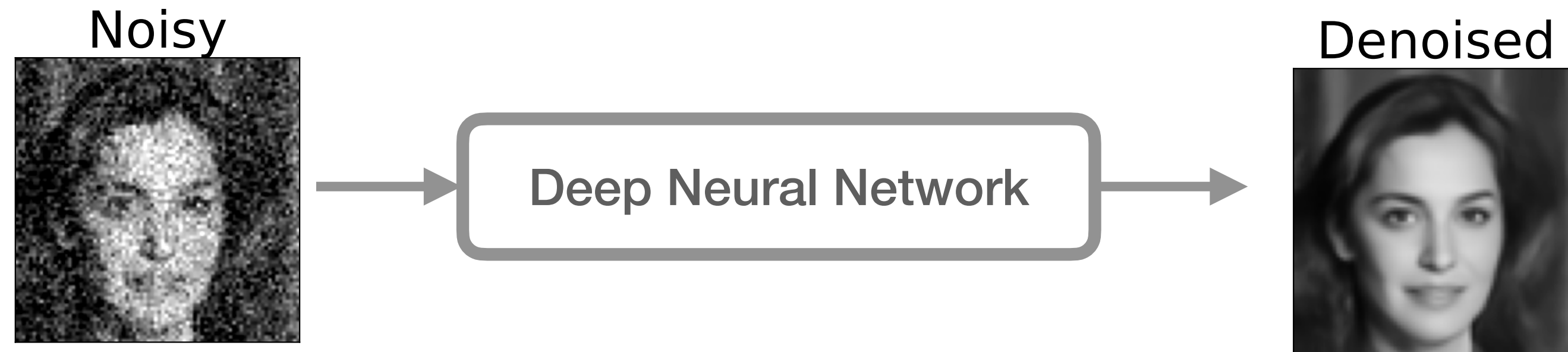
Denoising as shrinkage in an adaptive basis

Adaptive basis, adaptive shrinkage



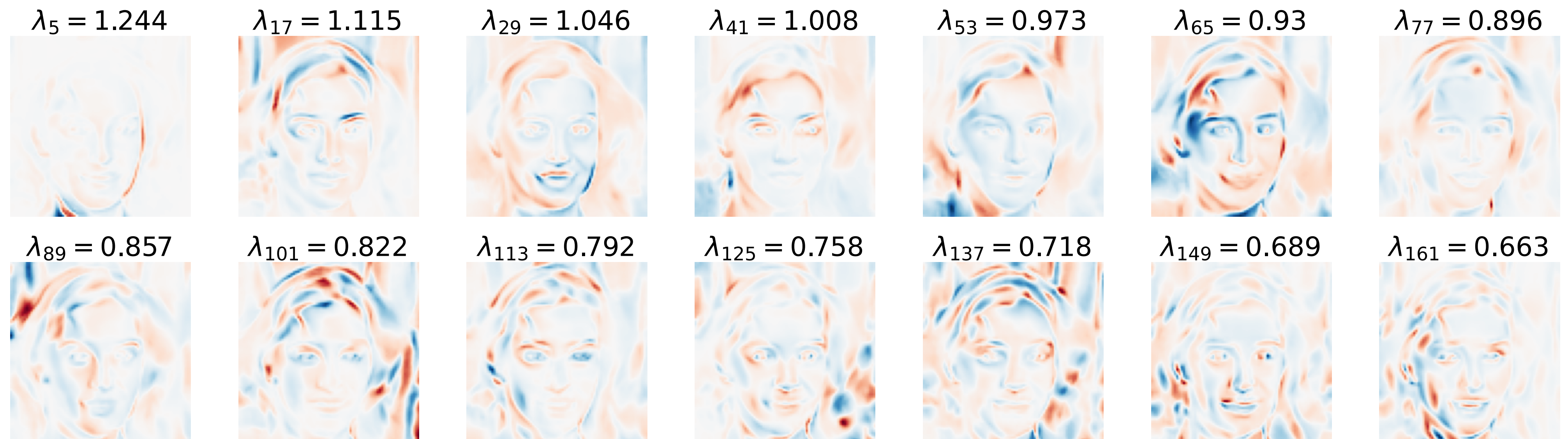
Denoising as shrinkage in an adaptive basis

Geometry Adaptive Harmonic Basis (GAHBs)



Some top Eigenvectors

- 1. Adaptive
- 2. oscillatory



Denoising as shrinkage in an adaptive basis

Geometry Adaptive Harmonic Basis (GAHBs)

hypothesis:

**DNN denoisers have inductive biases towards learning
GAHBs**

Denoising as shrinkage in an adaptive basis

Geometry Adaptive Harmonic Basis (GAHBs)

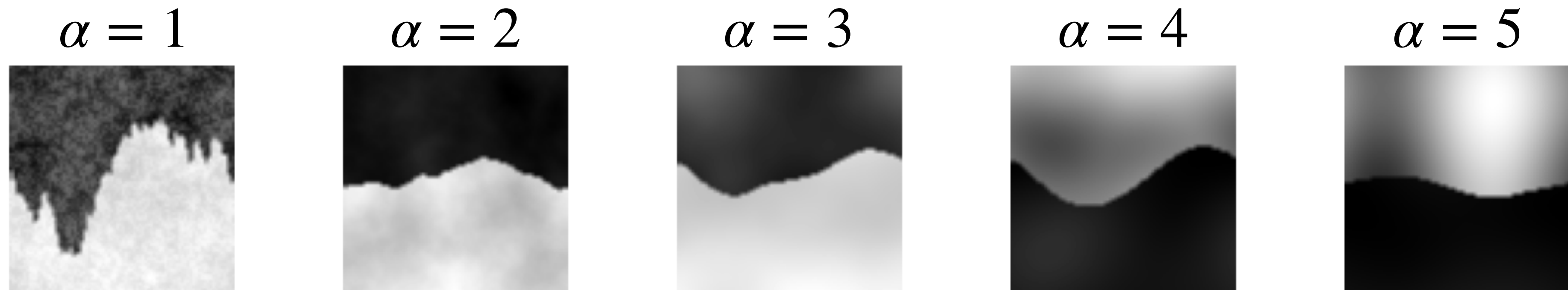
hypothesis:

**DNN denoisers have inductive biases towards learning
GAHBs**

How to test this?

Synthetic images

Geometric C^α images

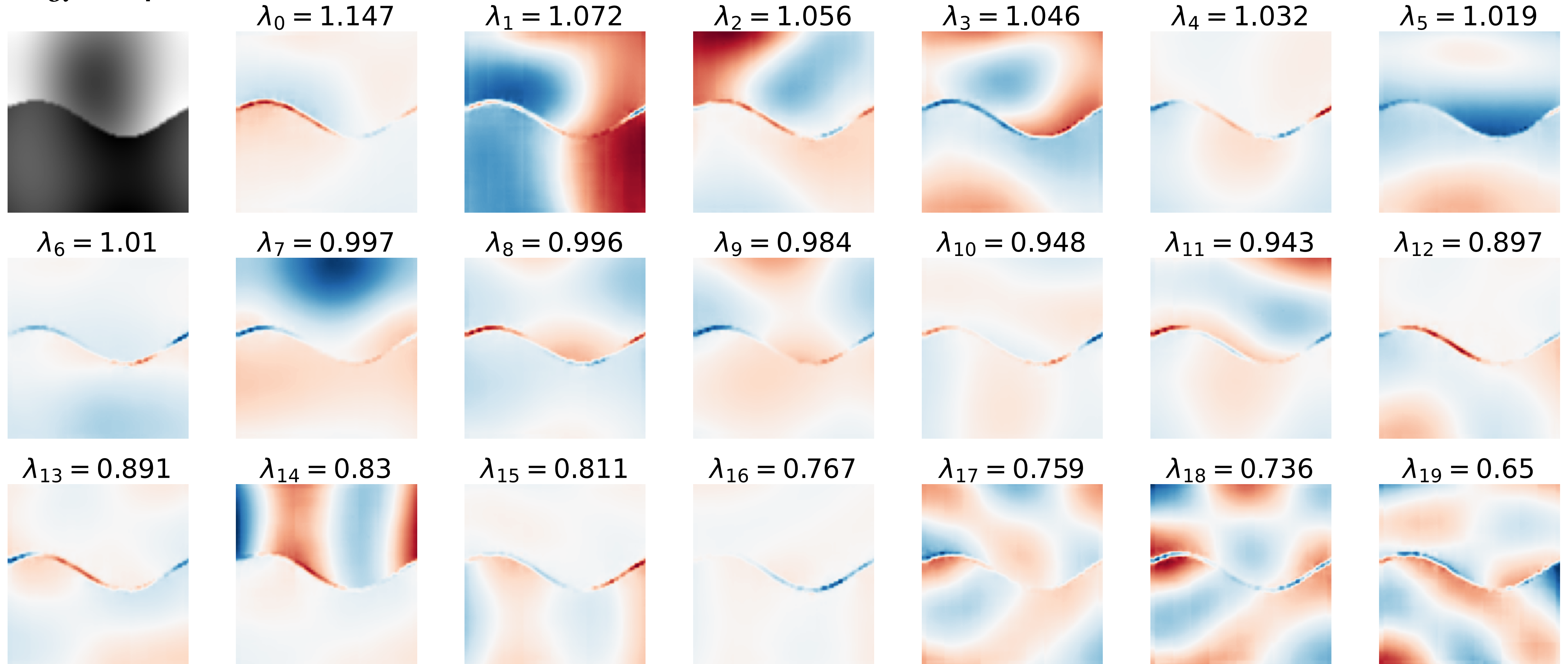


GAHBs are optimal for denoising these

[Korostelev & Tsybakov, 1993; Donoho, 1999; Peyré & Mallat, 2008]

Geometric C^α images

$\alpha = 4$



geometry adaptive harmonic basis

Geometric C^α images

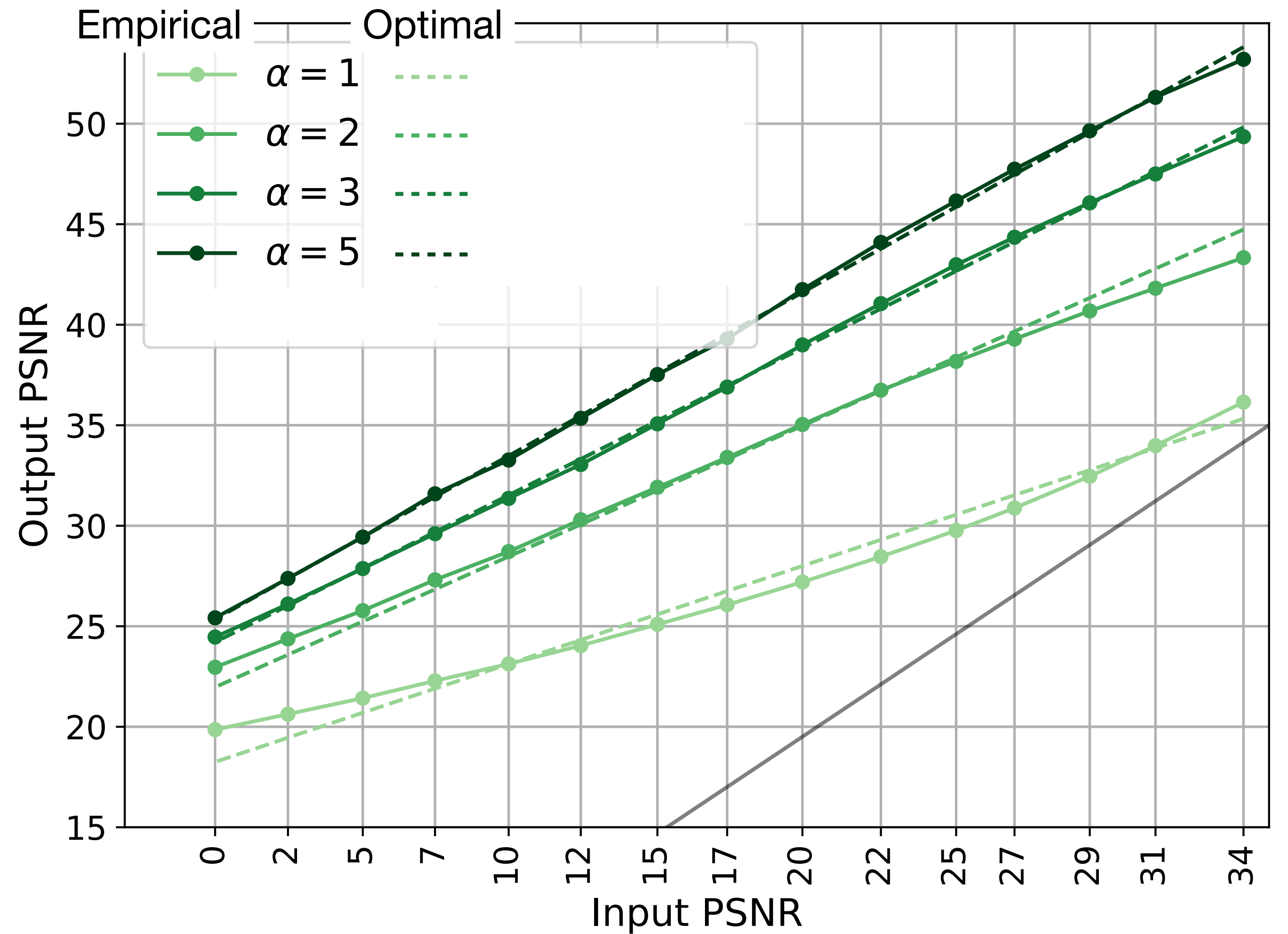
Optimal denoiser on C^α images

has slope $\frac{\alpha}{\alpha + 1}$.

[Korostelev & Tsybakov, 1993]

[Peyré & Mallat, 2008]

Denoising performance



Geometric C^α images

Optimal denoiser on C^α images

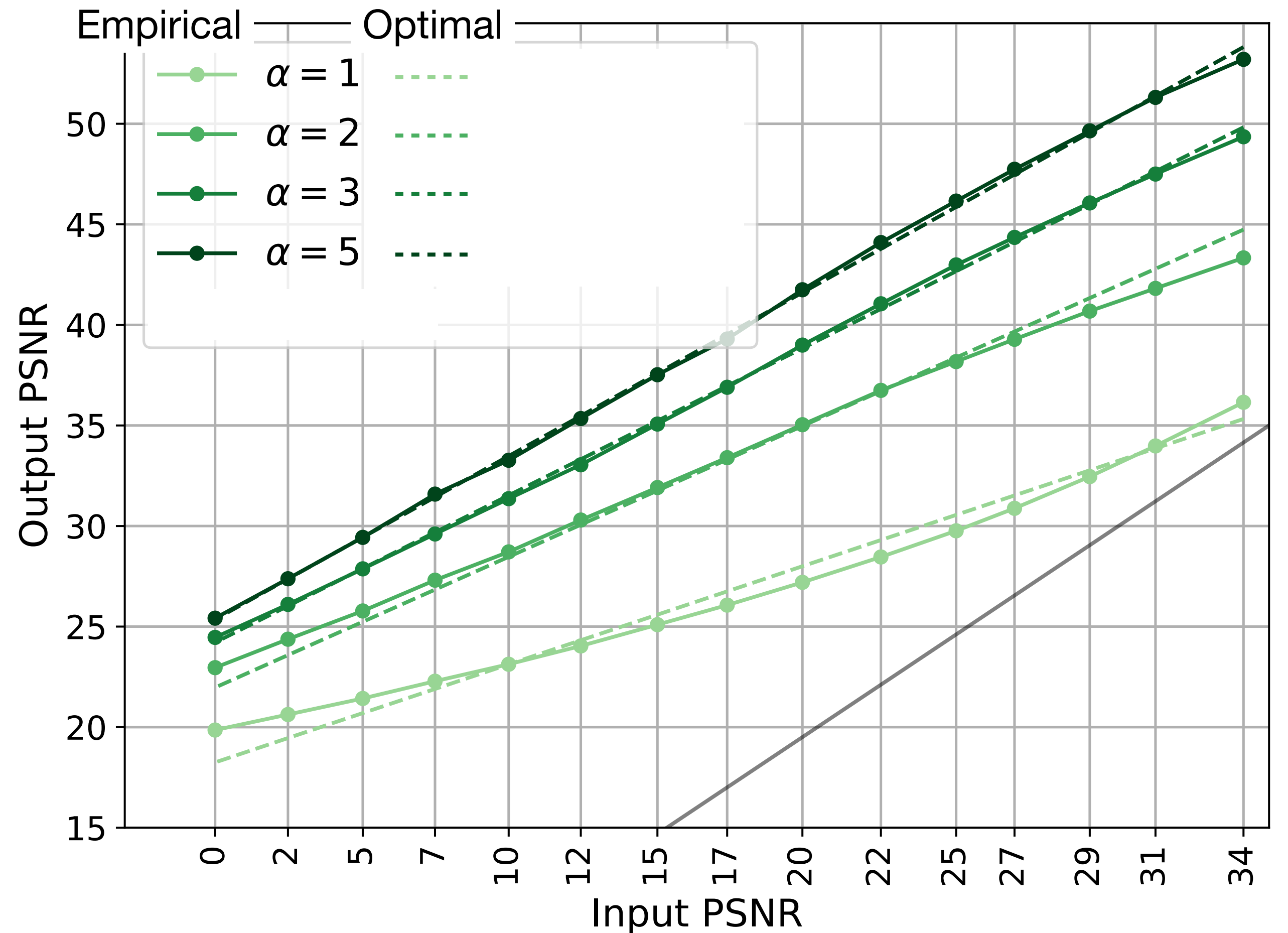
has slope $\frac{\alpha}{\alpha + 1}$.

[Korostelev & Tsybakov, 1993]

[Peyré & Mallat, 2008]

✓ **Deep nets learn GAHB for denoising when it's optimal**

Denoising performance



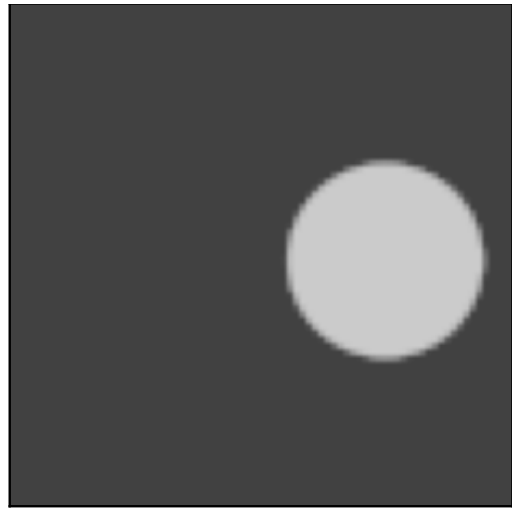
Manifold of disks



five-dimensional *curved* manifold

1. Vertical position
2. Horizontal position
3. Radius/size
4. Foreground intensity
5. Background intensity

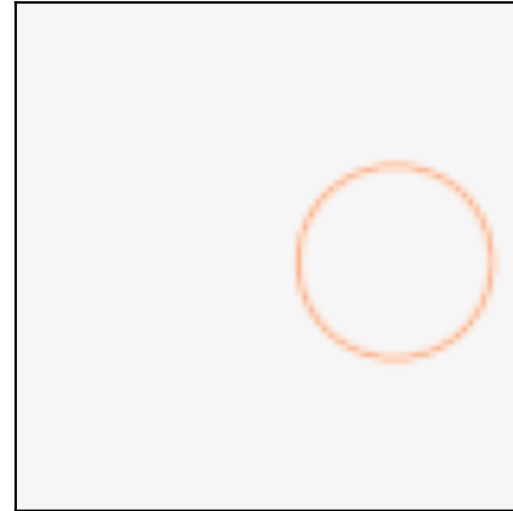
Manifold of disks



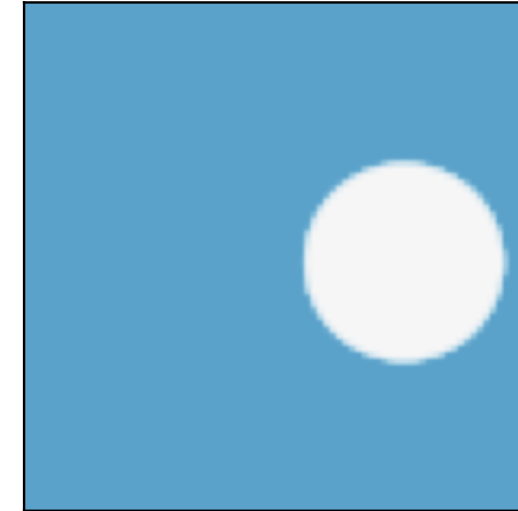
**Vertical
Translation**



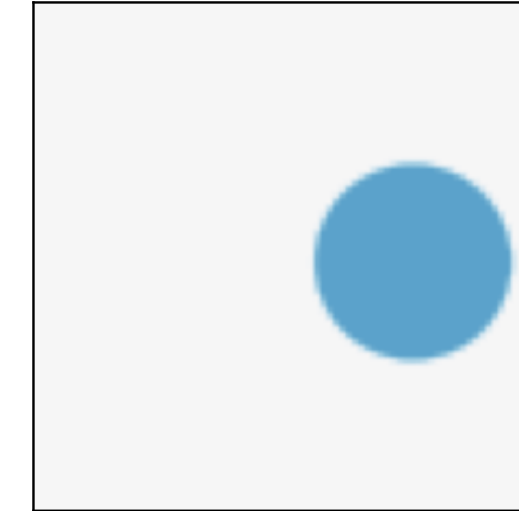
**Radius
Change**



**Background
Change**



**Foreground
Change**

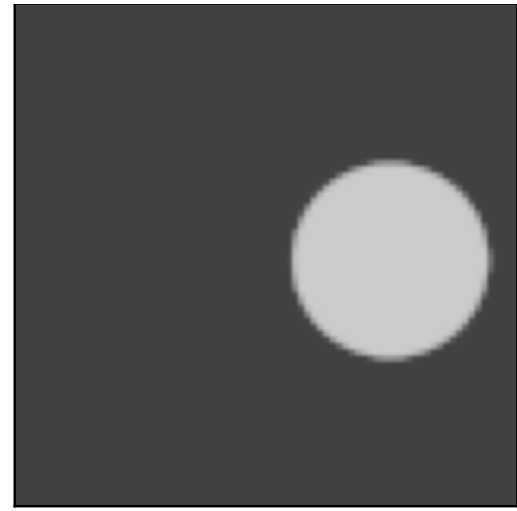


**Horizontal
Translation**



Optimal

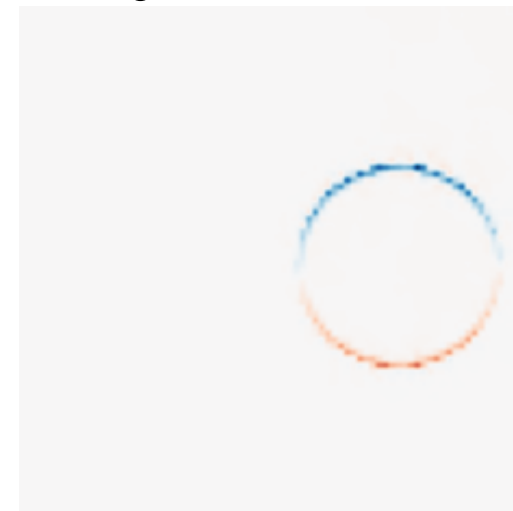
Manifold of disks



**Vertical
Translation**



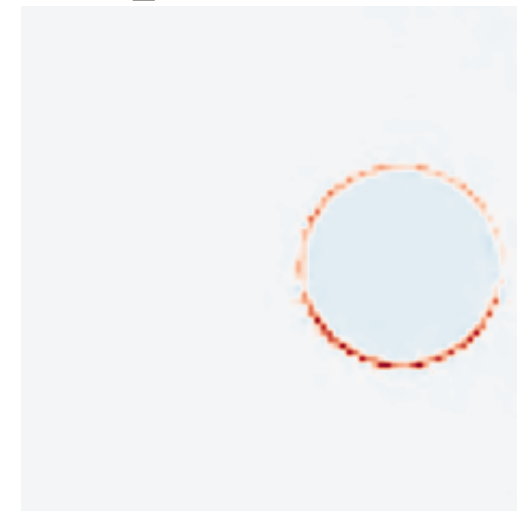
$$\lambda_0 = 1.177$$



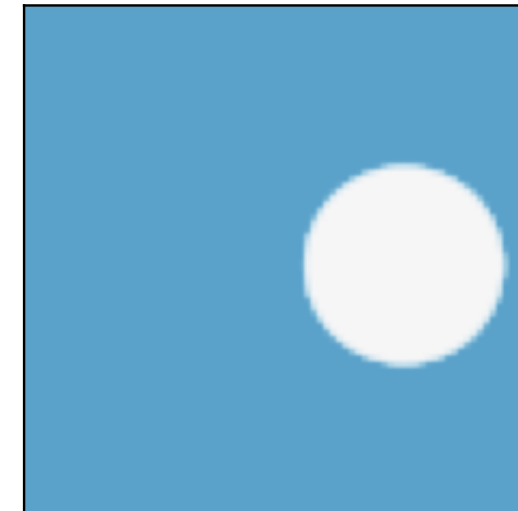
**Radius
Change**



$$\lambda_1 = 1.067$$



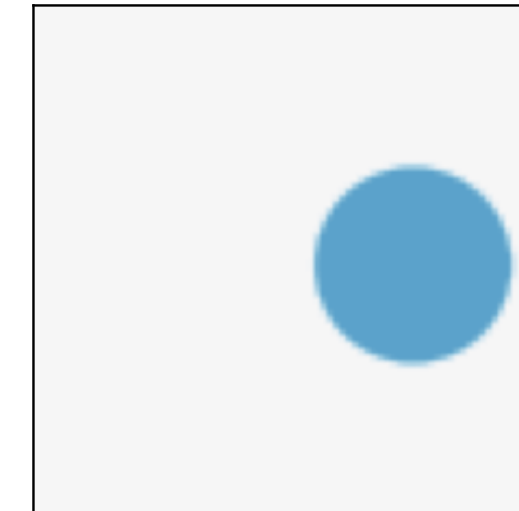
**Background
Change**



$$\lambda_2 = 1.004$$



**Foreground
Change**



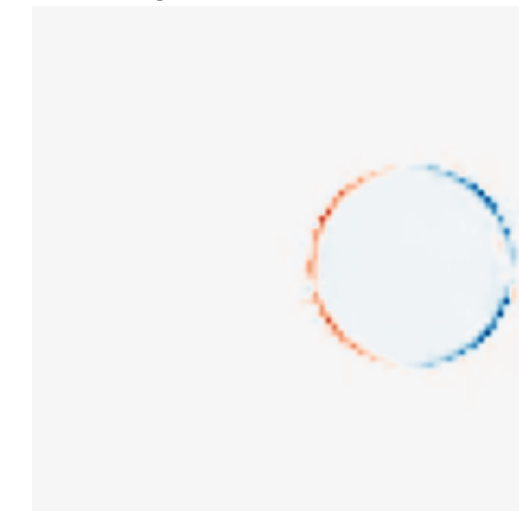
$$\lambda_3 = 0.999$$



**Horizontal
Translation**



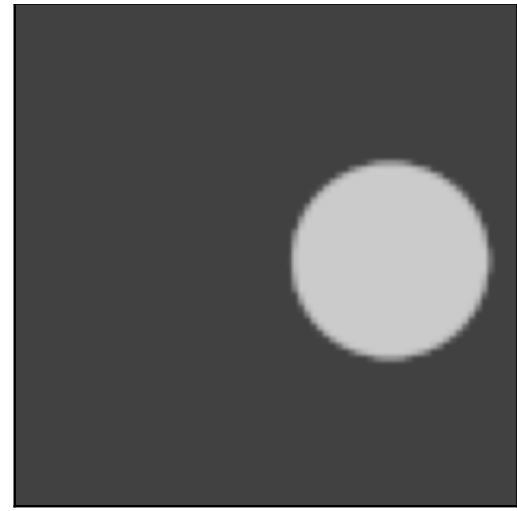
$$\lambda_4 = 0.945$$



Optimal

Empirical

Manifold of disks



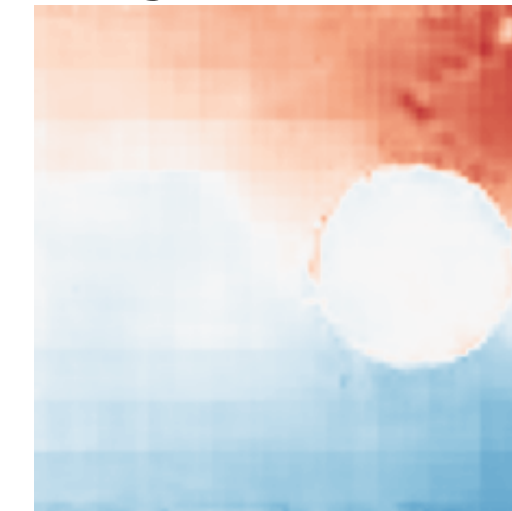
Vertical Translation



$\lambda_0 = 1.177$



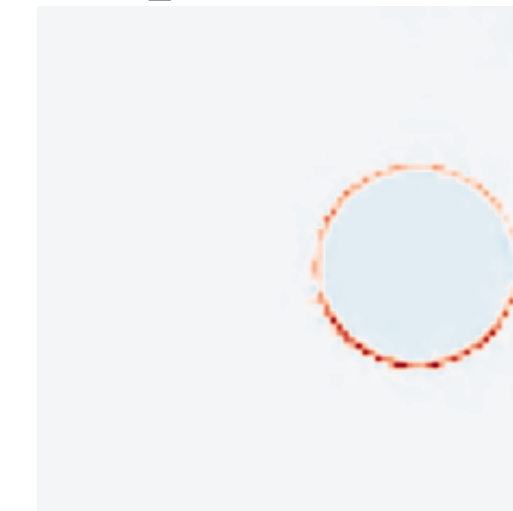
$\lambda_5 = 0.674$



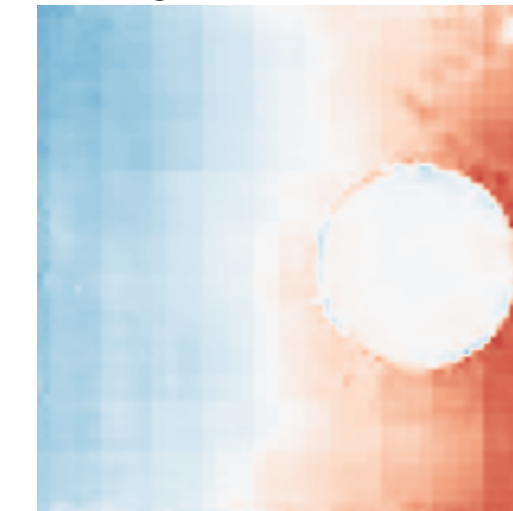
Radius Change



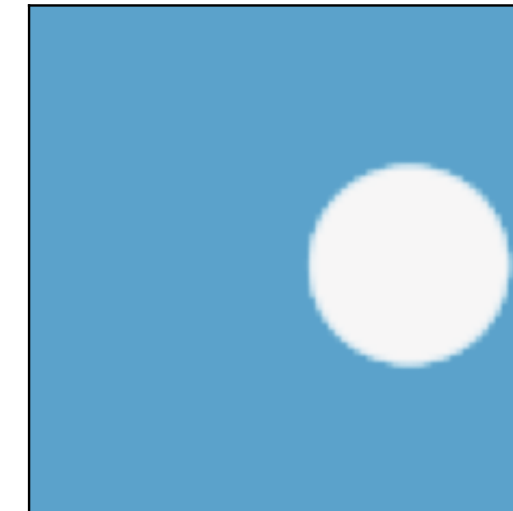
$\lambda_1 = 1.067$



$\lambda_6 = 0.552$



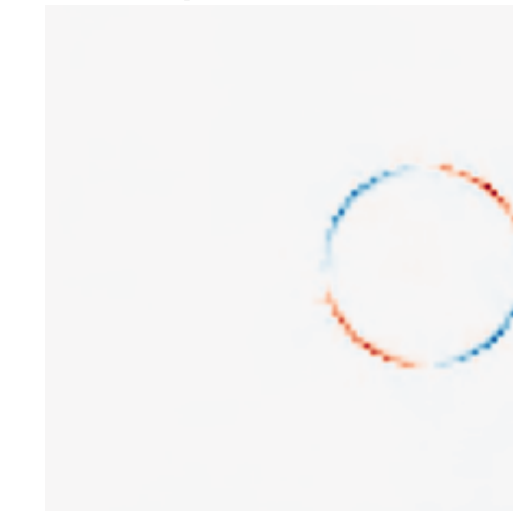
Background Change



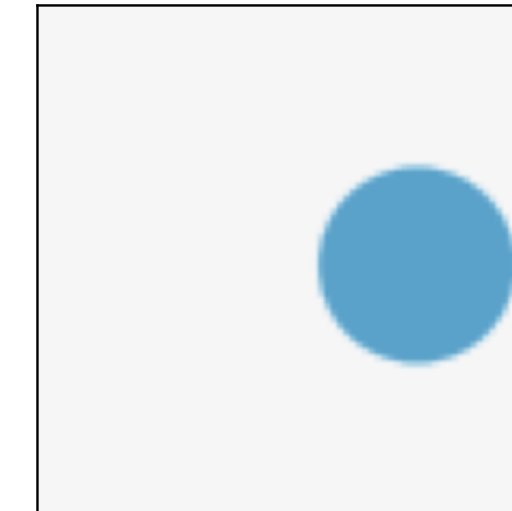
$\lambda_2 = 1.004$



$\lambda_7 = 0.39$



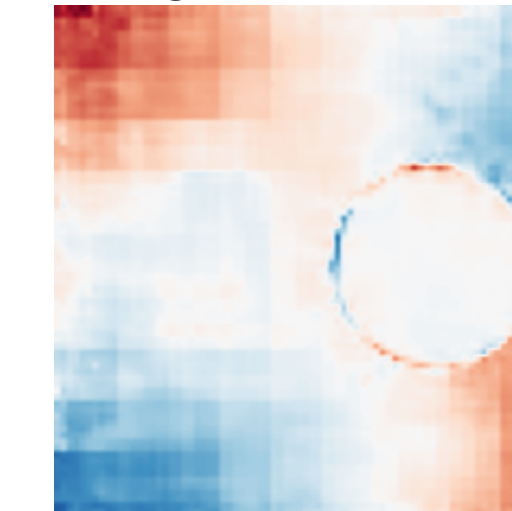
Foreground Change



$\lambda_3 = 0.999$



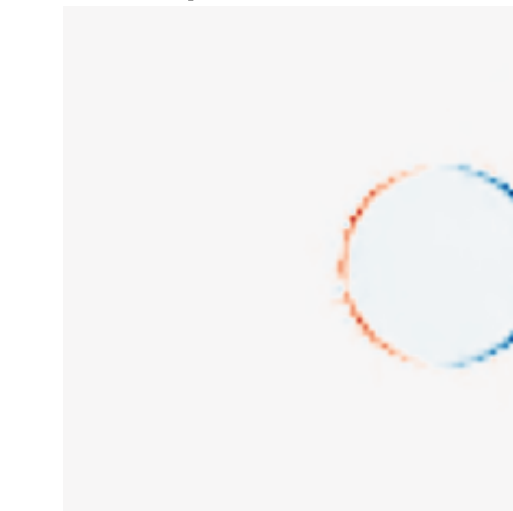
$\lambda_8 = 0.304$



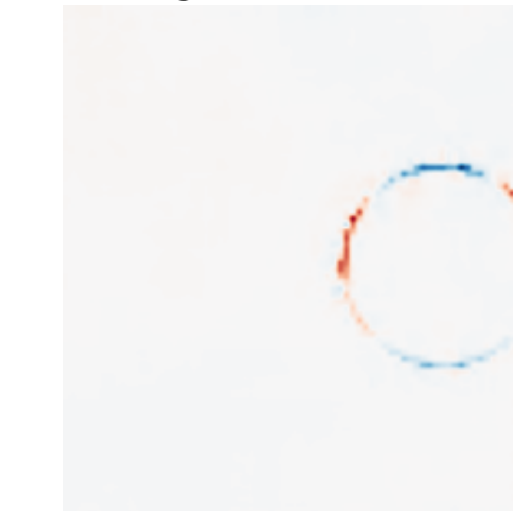
Horizontal Translation



$\lambda_4 = 0.945$



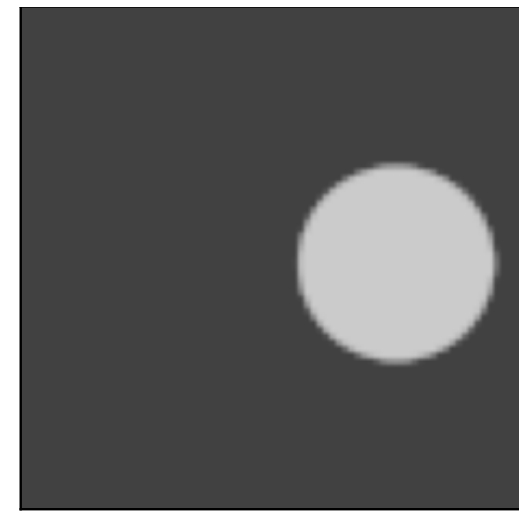
$\lambda_9 = 0.278$



Optimal

Empirical

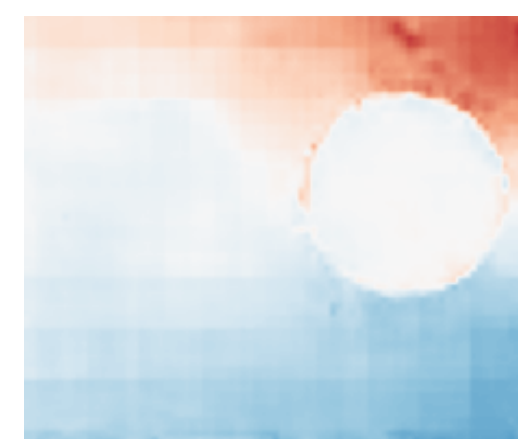
Manifold of disks



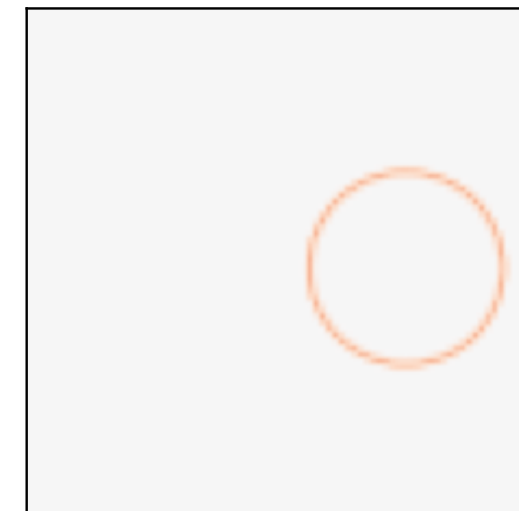
Vertical Translation



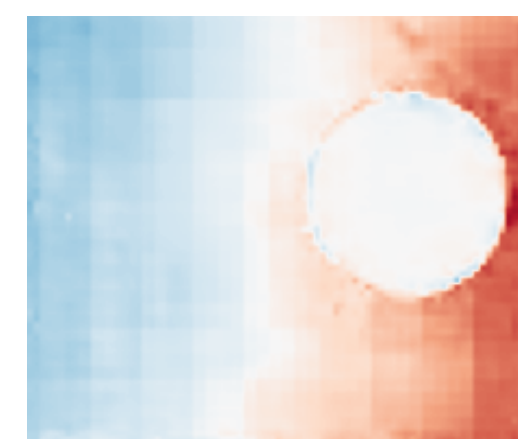
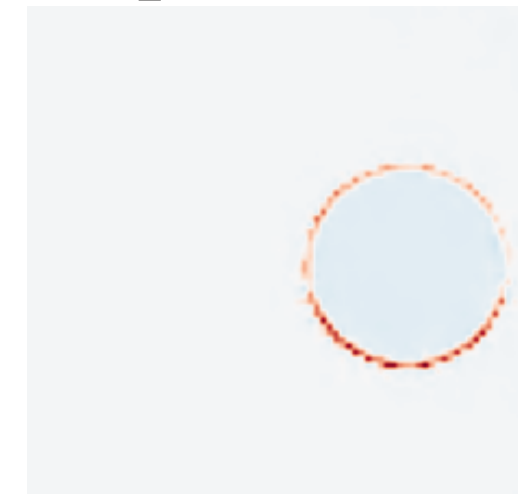
$$\lambda_0 = 1.177$$



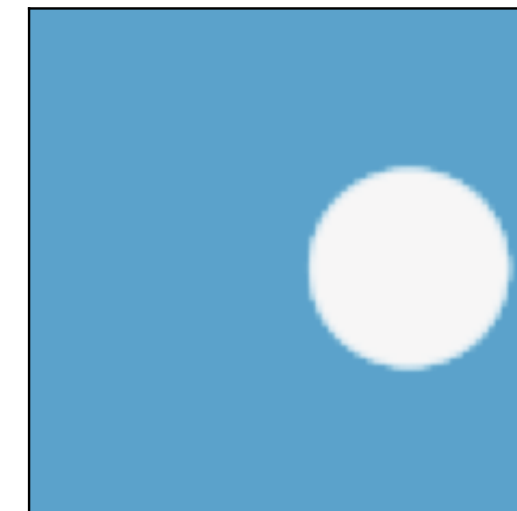
Radius Change



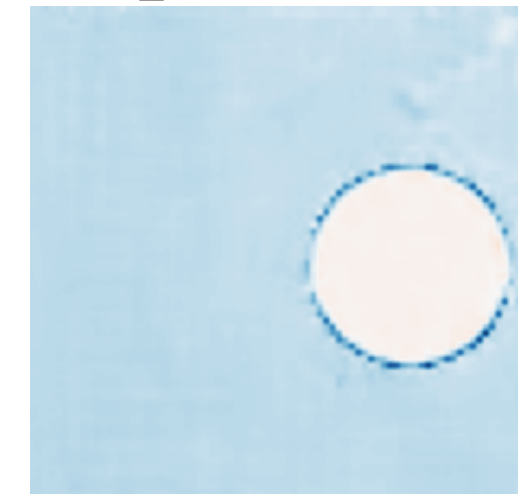
$$\lambda_1 = 1.067$$



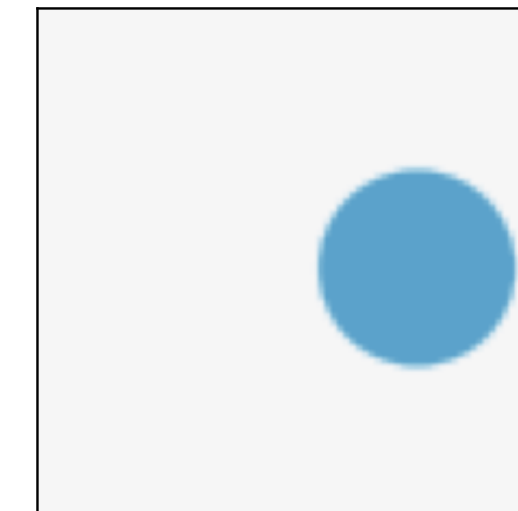
Background Change



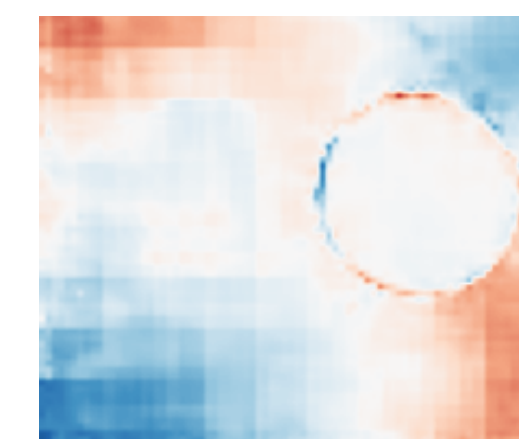
$$\lambda_2 = 1.004$$



Foreground Change



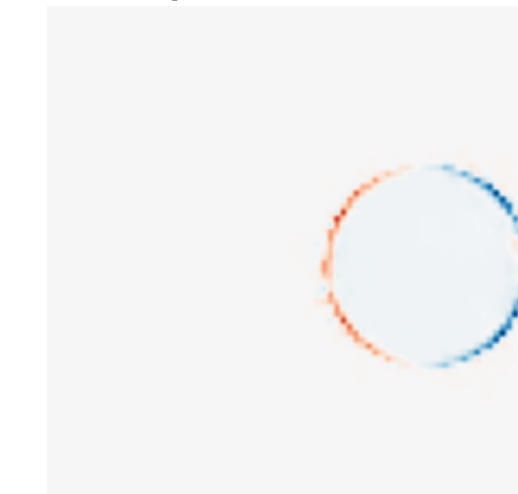
$$\lambda_3 = 0.999$$



Horizontal Translation



$$\lambda_4 = 0.945$$



Optimal

Empirical

Deep nets learn GAHB even when it's sub-optimal



Interim summary

- Diffusion models can transition from **memorization** to **generalization** with large enough training set size
- Generalization is **strong**: two denoisers trained on non-overlapping training sets converge to nearly the same function
- Generalization due to an **inductive bias** corresponding to shrinkage in a Geometry Adaptive Harmonic Basis (GAHB)

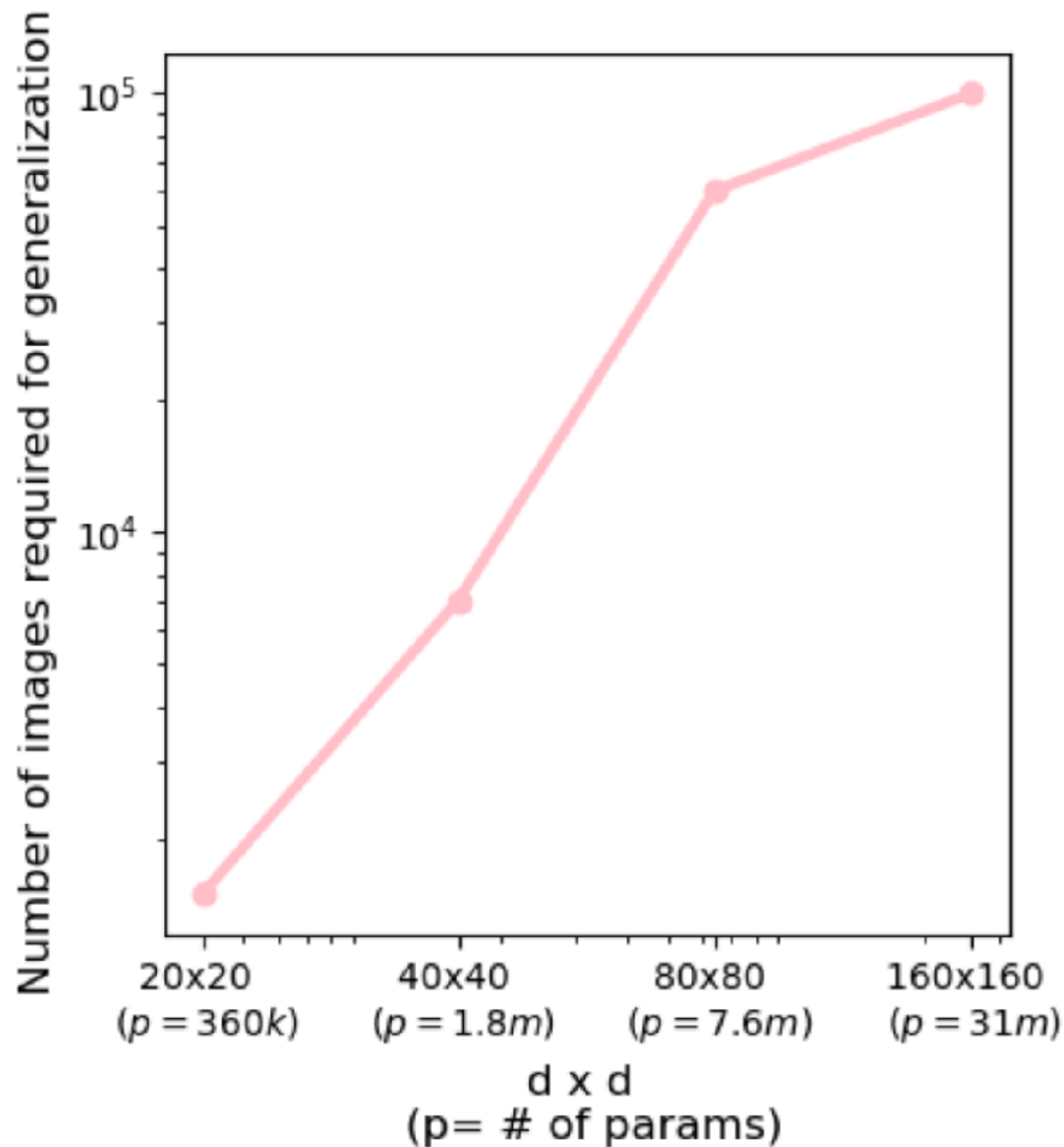
Kadkhodaie Z, Guth F, Simoncelli EP, Mallat S.

“Generalization in diffusion models arises from geometry-adaptive harmonic representation”. ICLR 2024.

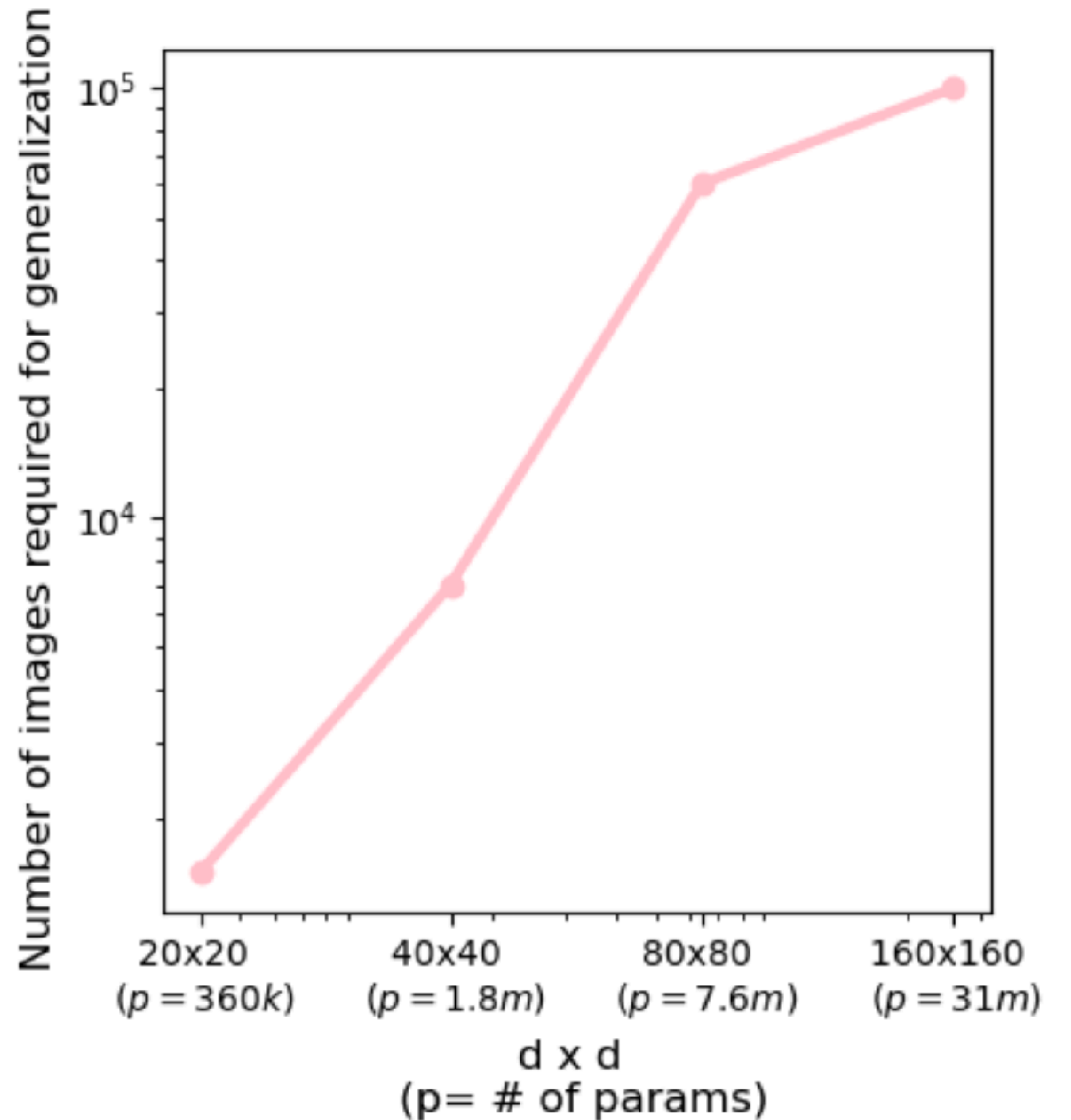
- We are in good shape with learning densities
- Can we reduce the size of training set required for generalization?

- Image resolution
- Network size
- Complexity of image dataset

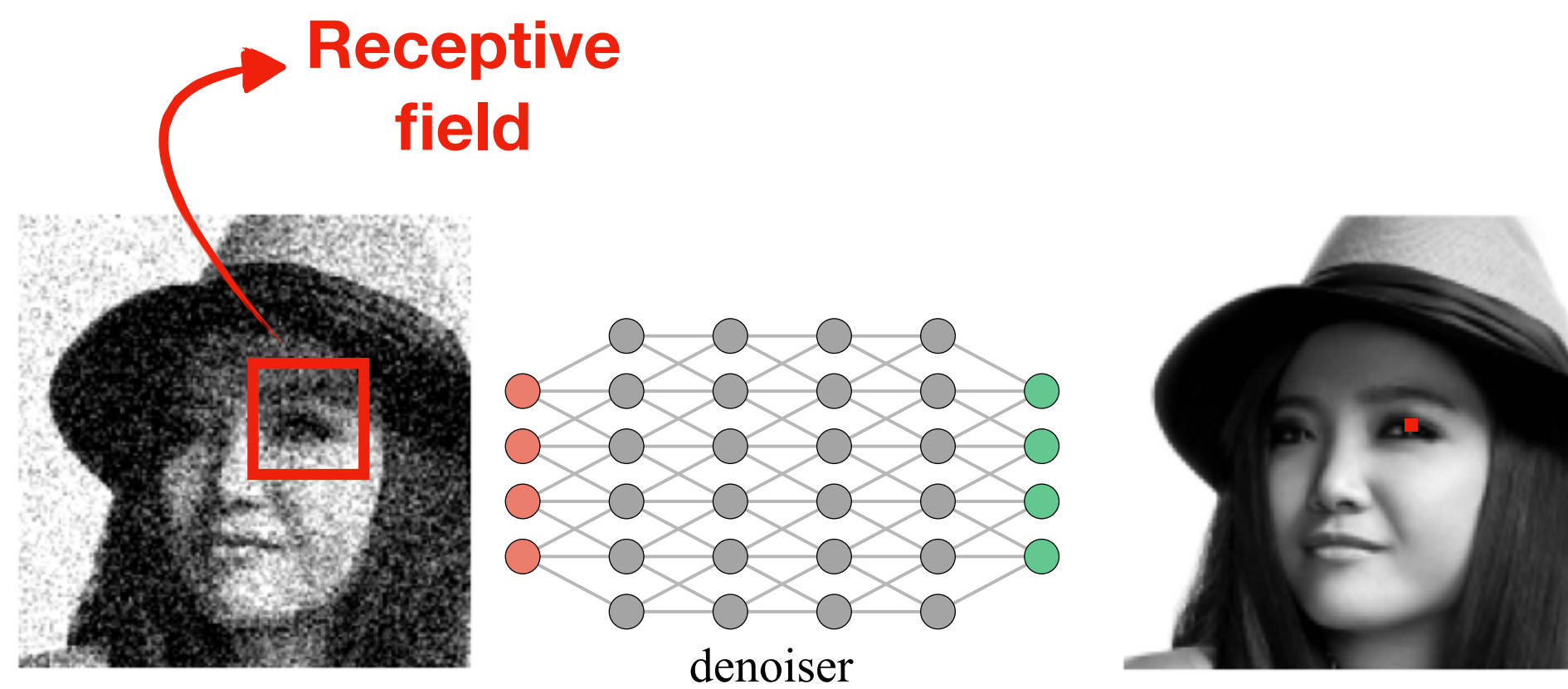
- Image resolution
- Network size
- Complexity of image dataset



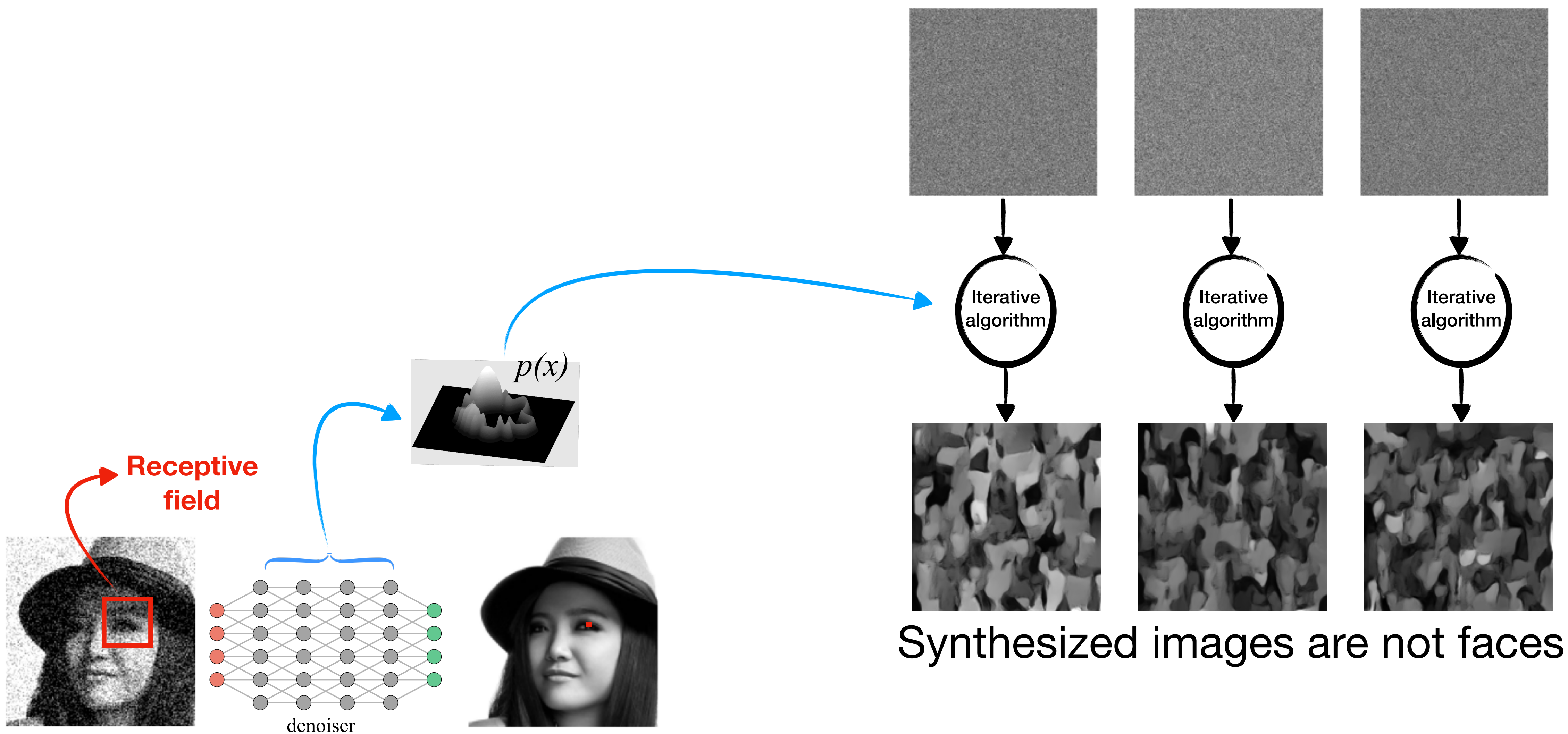
Hundreds of millions or
billions of parameters
with **global receptive fields**



Network receptive field



Synthesis fails without global receptive fields!

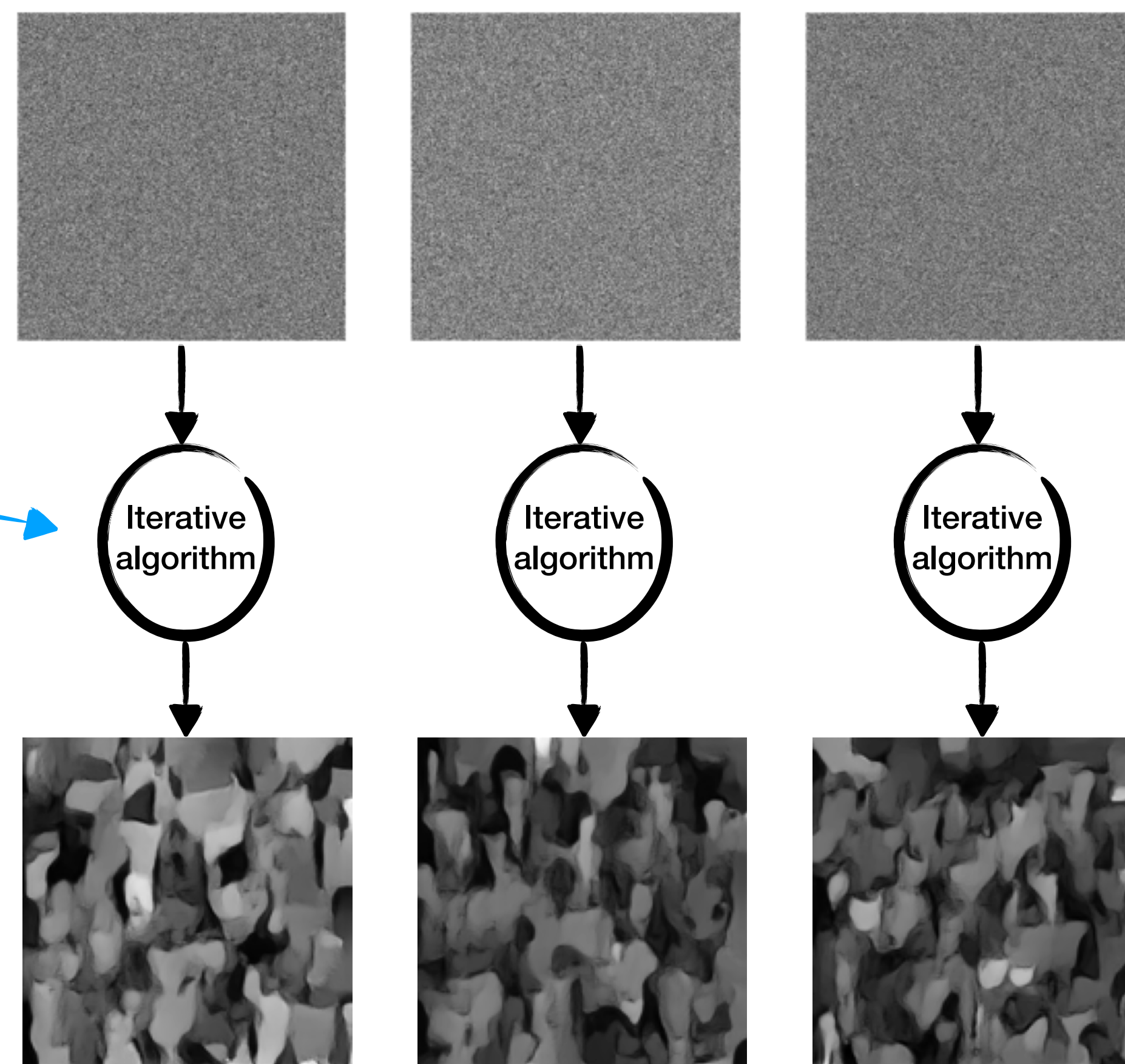
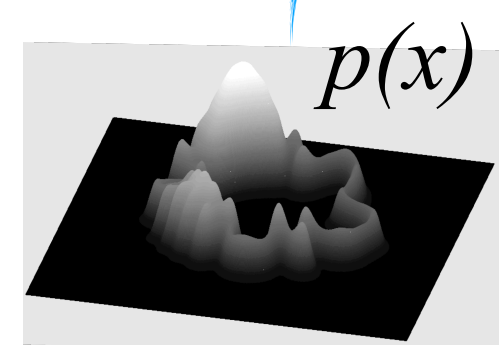
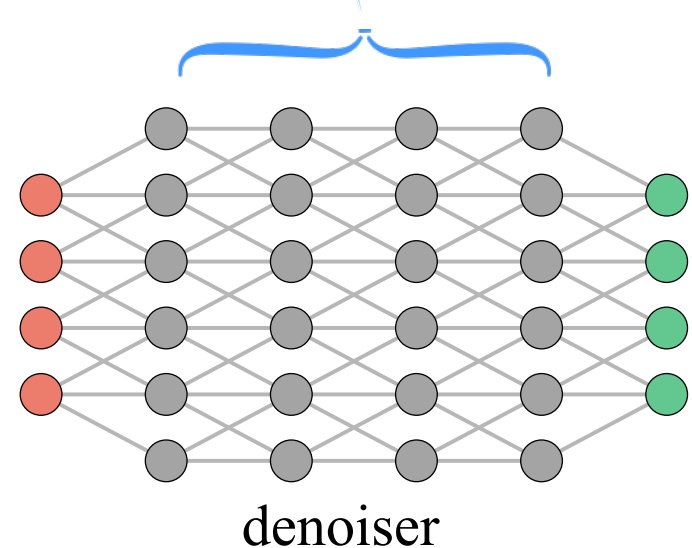


Synthesis fails without global receptive fields!

local and translation invariant RF \Rightarrow learn density is an MRF



Receptive field



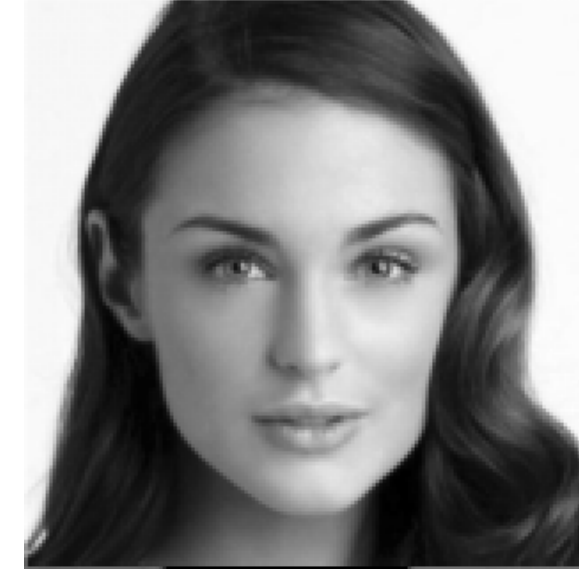
Synthesized images are not faces

**Do bigger images require
bigger models?**

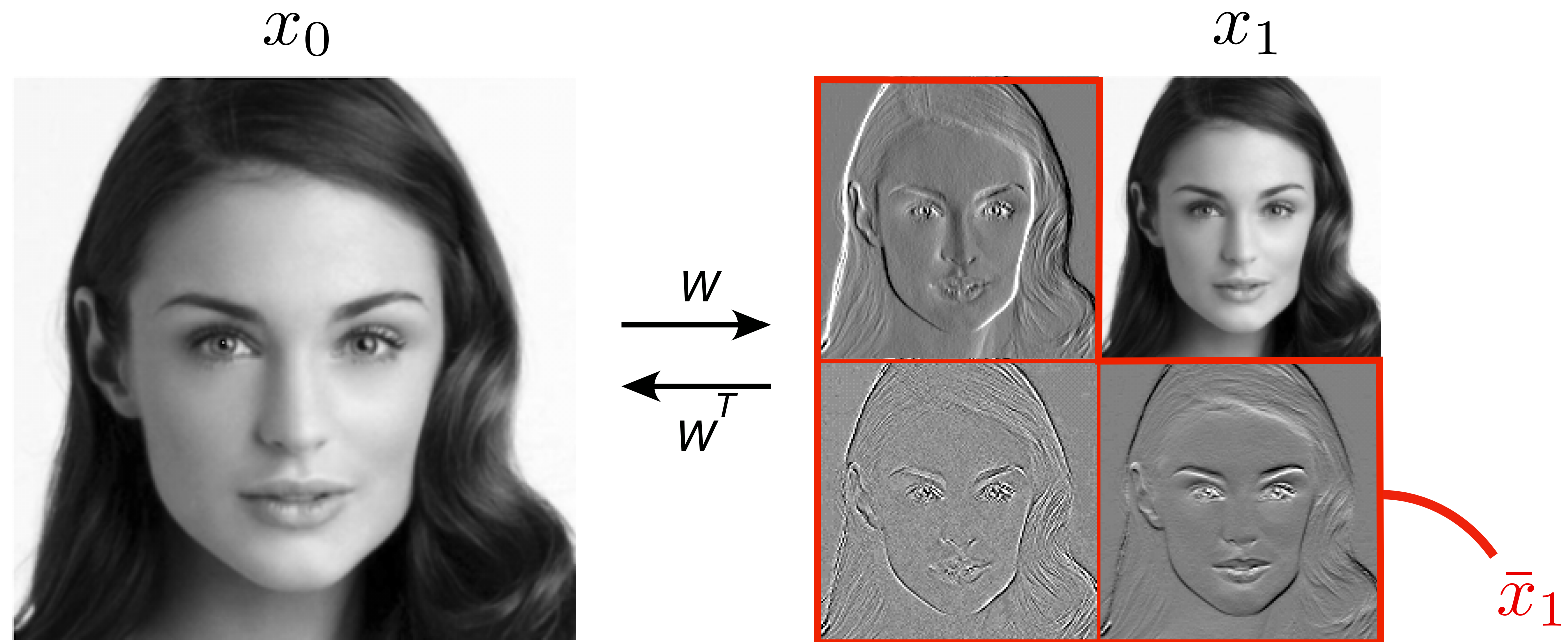
x_0



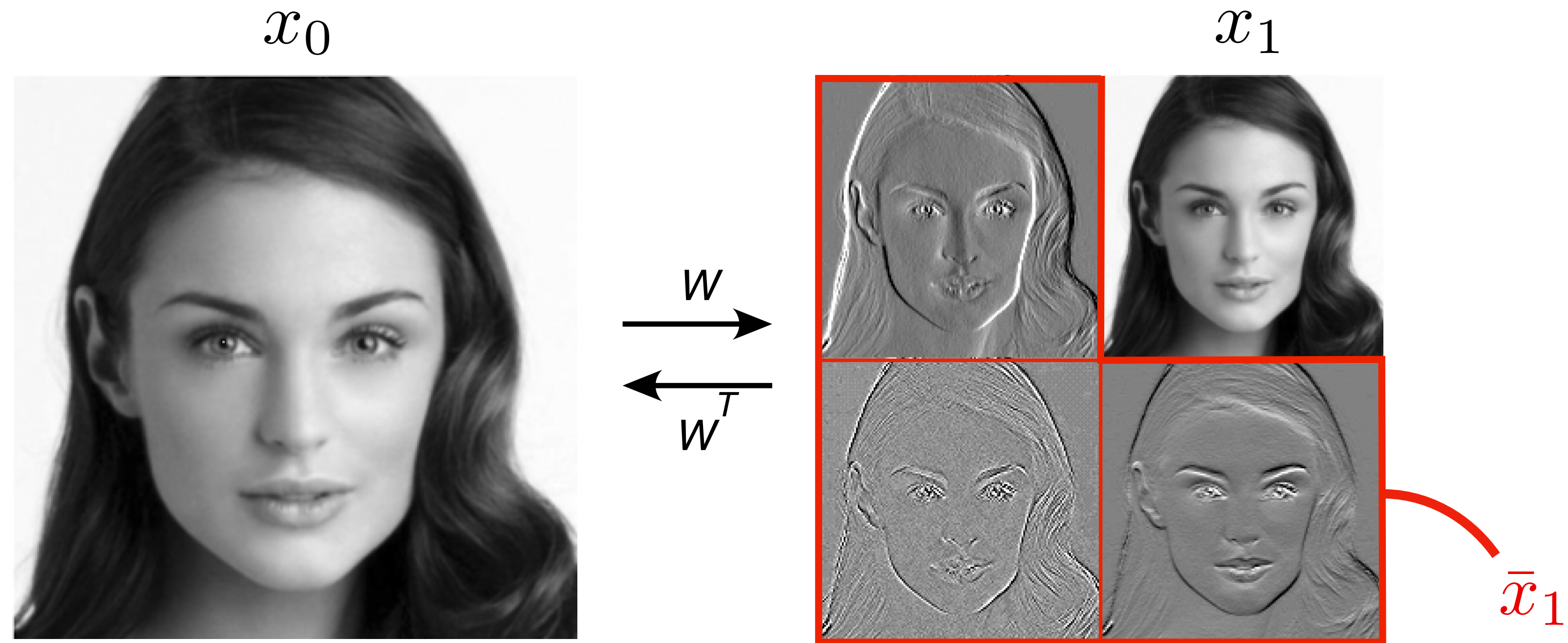
x_1



Wavelet decomposition

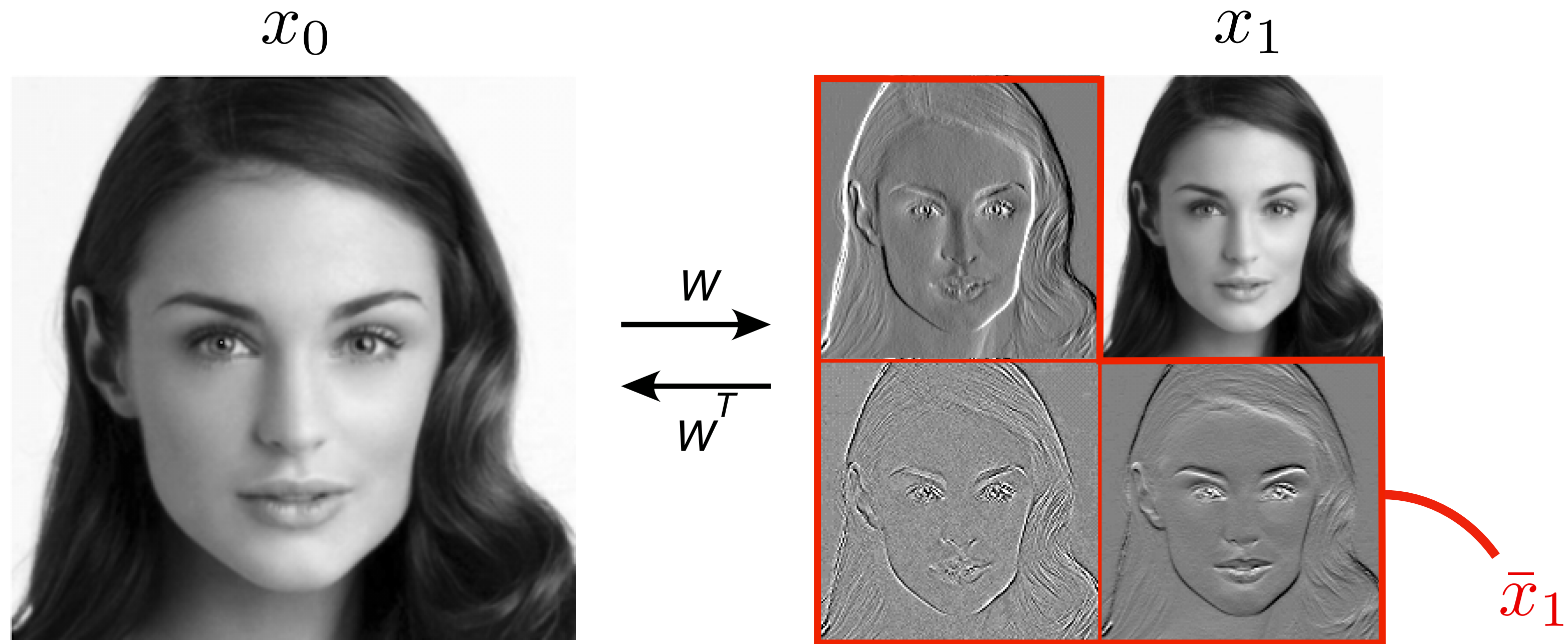


Wavelet decomposition



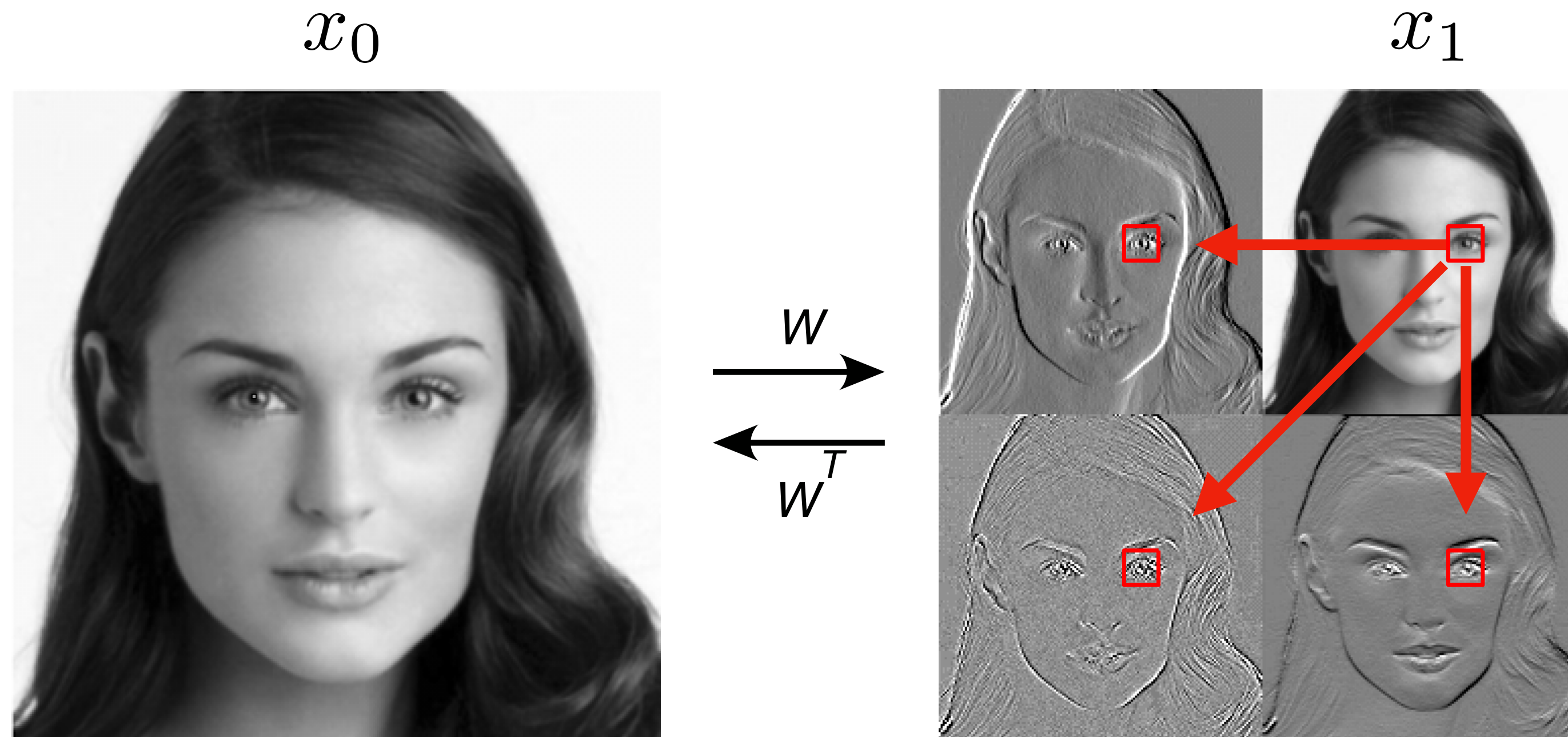
$$p(x_0) = p(x_1, \bar{x}_1)$$

Factorization of the prior



$$p(x_0) = p(x_1) p(\bar{x}_1|x_1)$$

Conditional densities are local



$$p(x_0) = p(x_1) p(\bar{x}_1 | x_1)$$

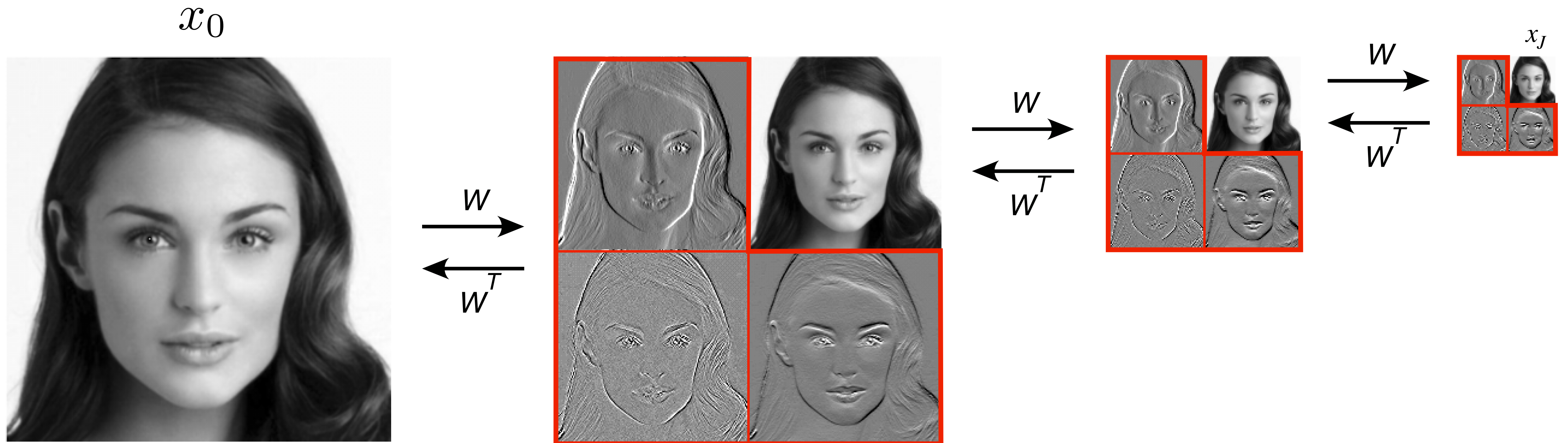


Global



Local

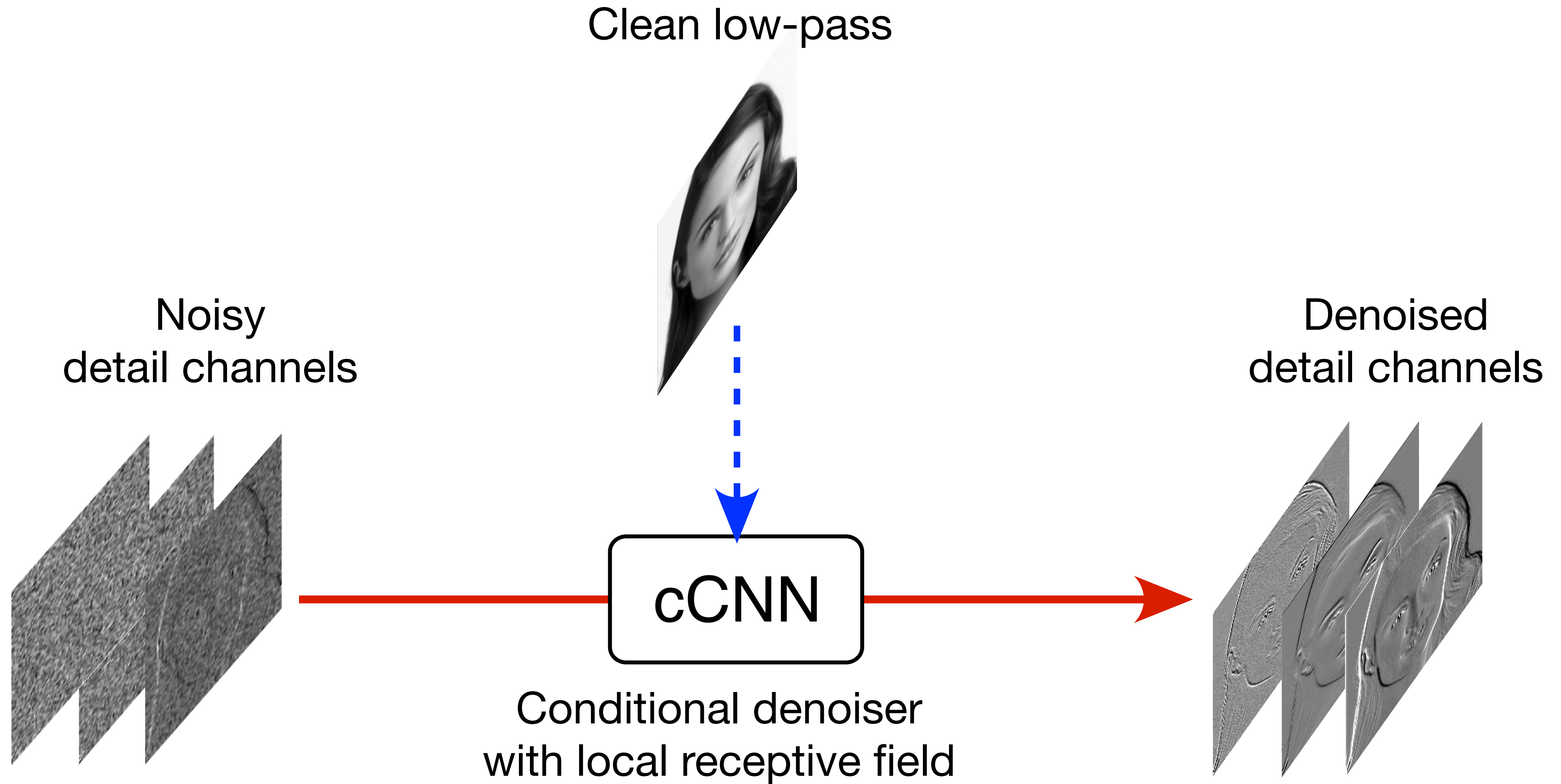
Multi-scale wavelet representation



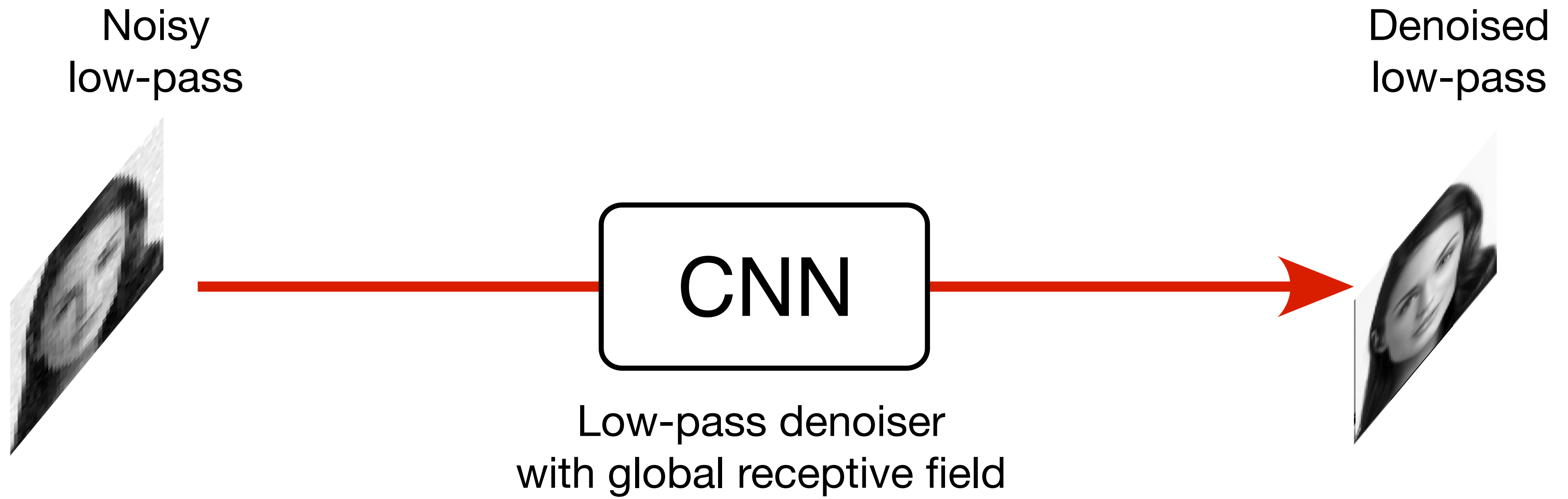
$$p(x_0) = p(x_J) \prod_{j=1}^J p(\bar{x}_j | x_j)$$

↑
↑
Global
Local

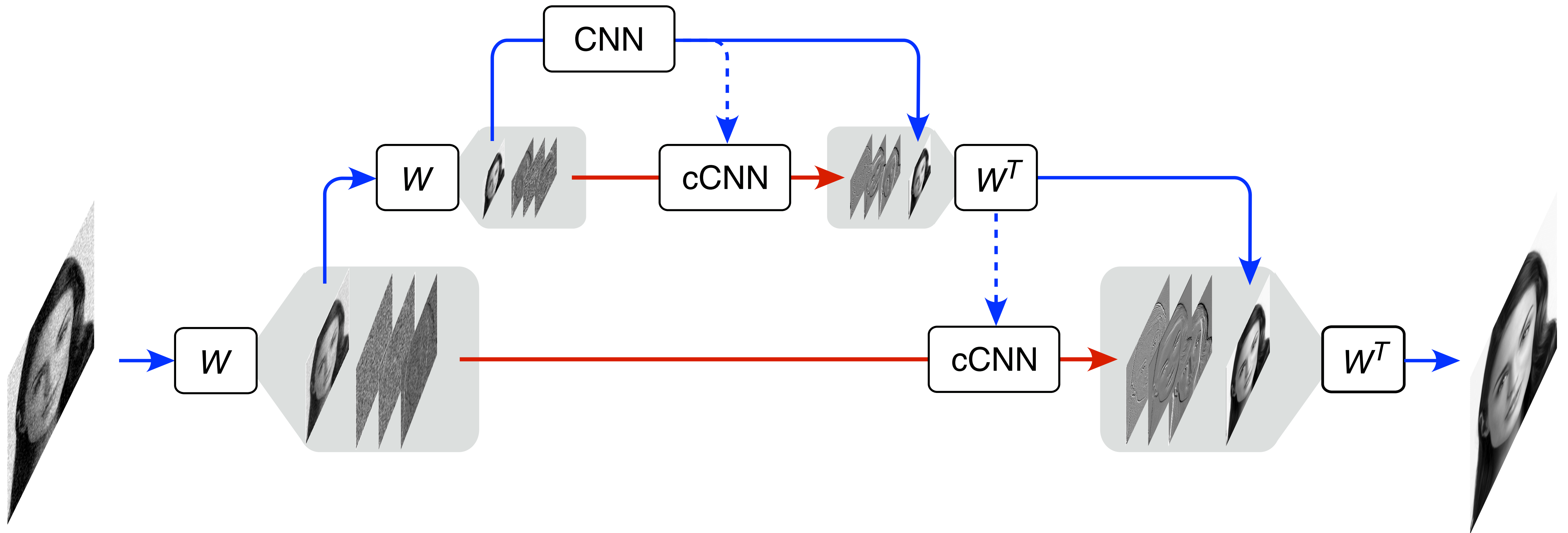
Conditional denoisers



Low pass denoiser

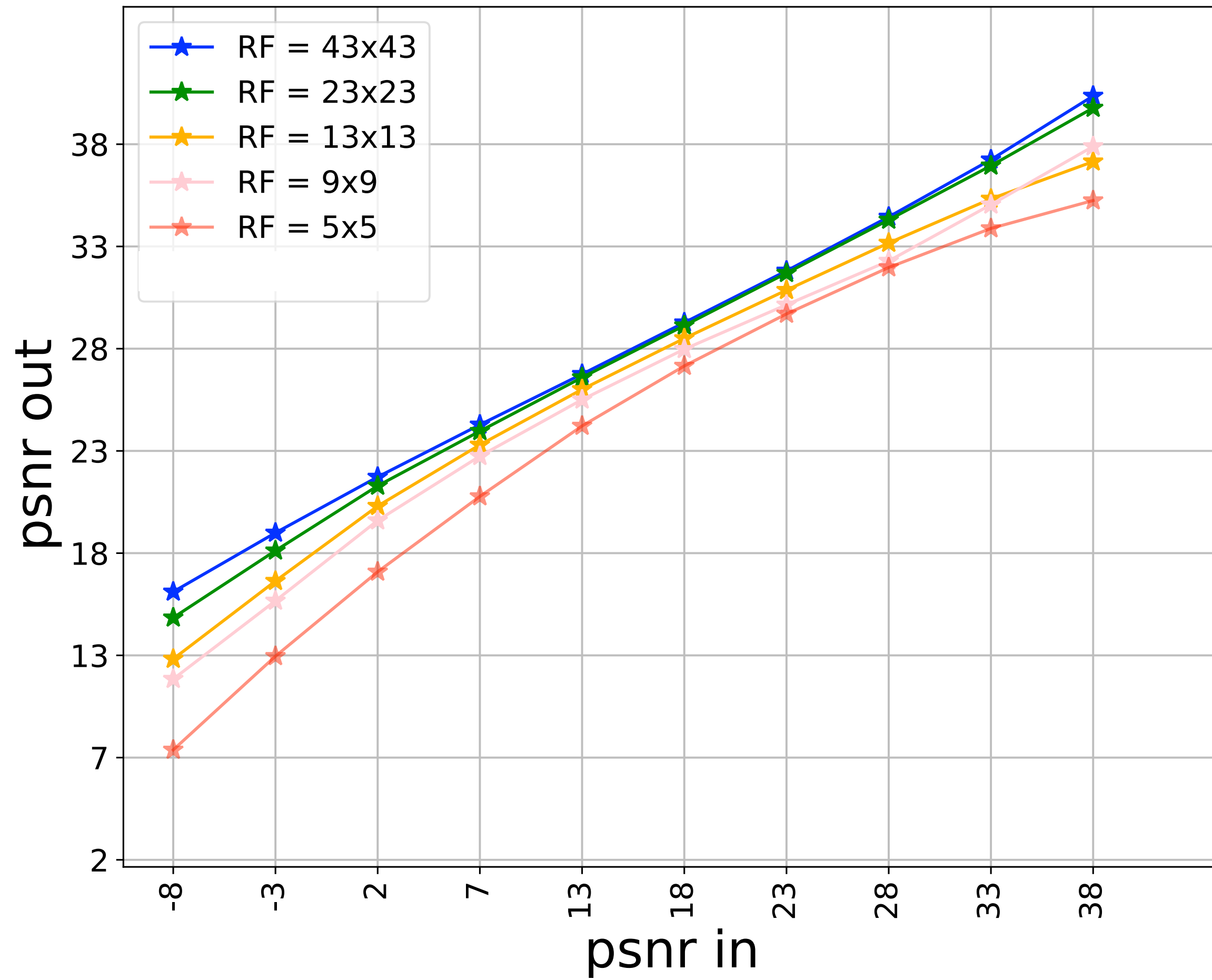


Multi-scale wavelet conditional denoiser



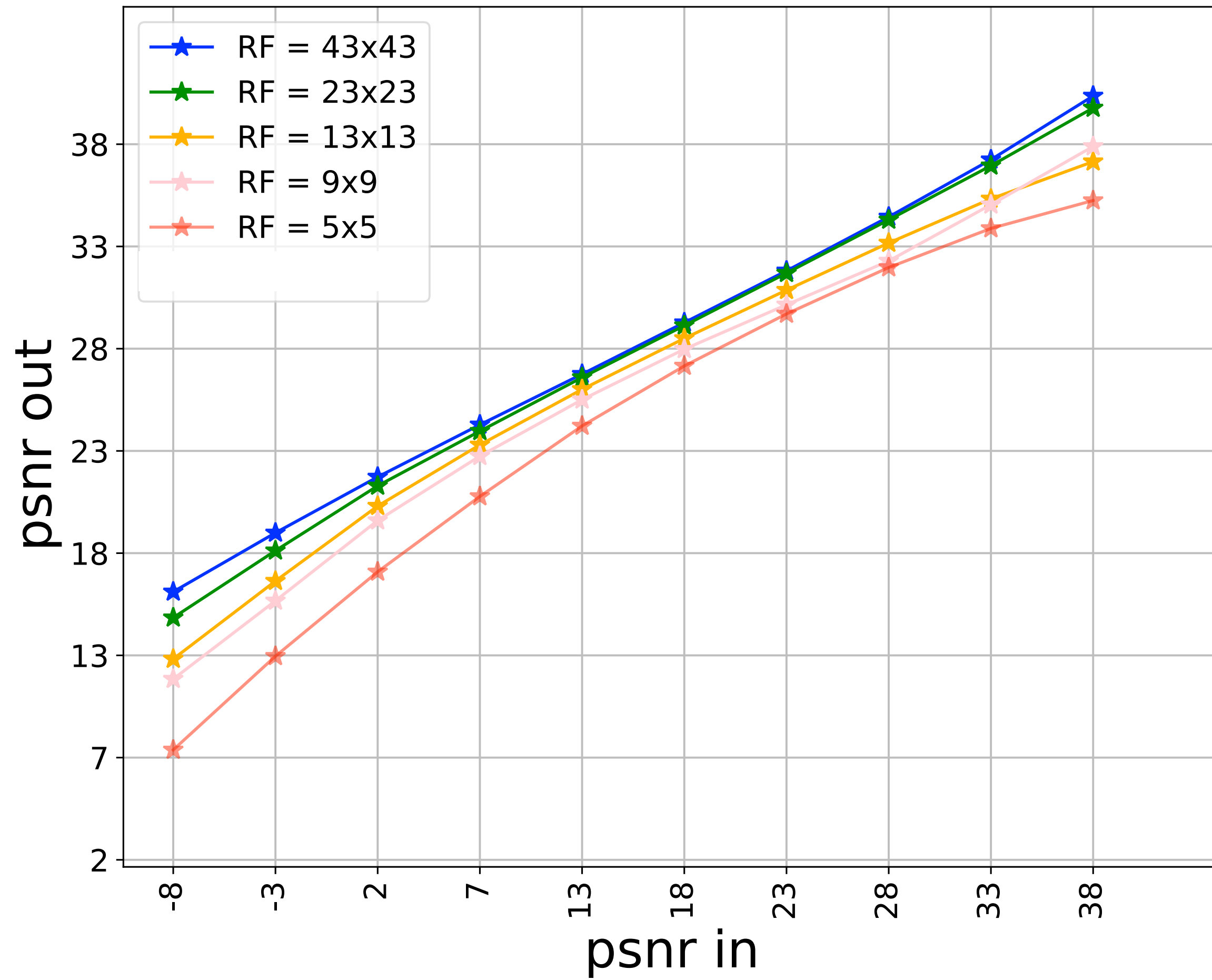
How local?

Pixel-domain denoiser

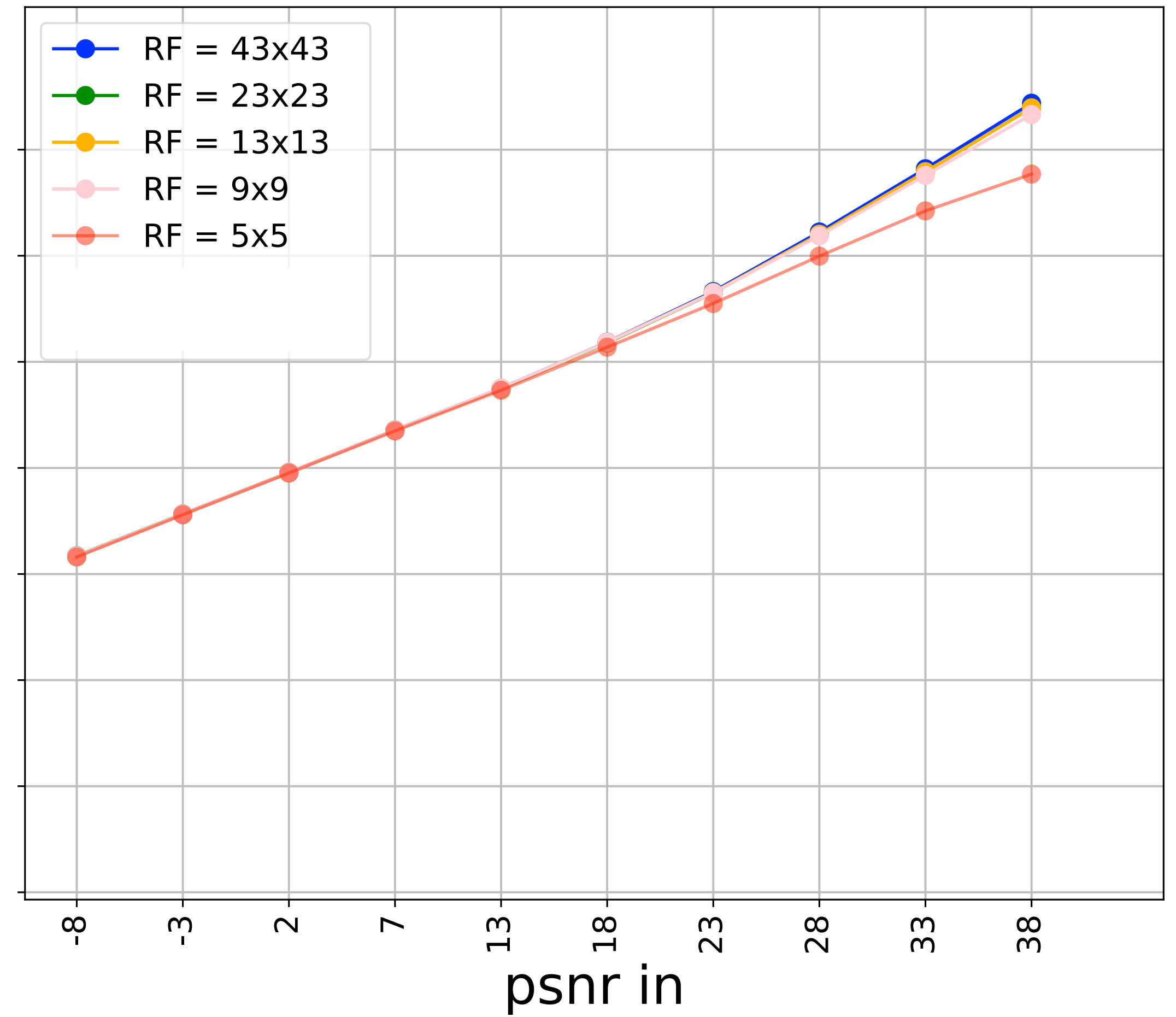


How local?

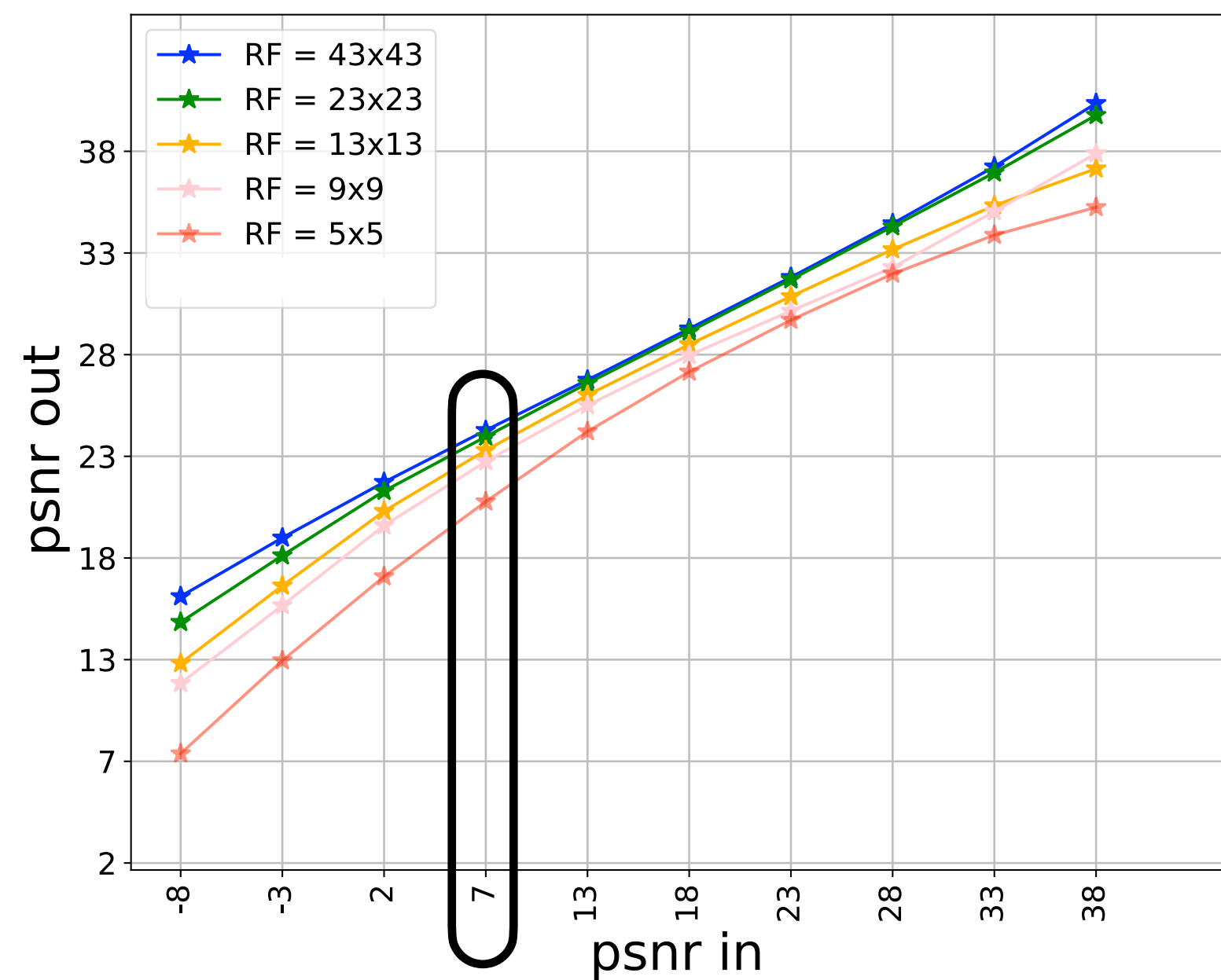
Pixel-domain denoiser



Wavelet-domain denoiser



Pixel-domain denoiser

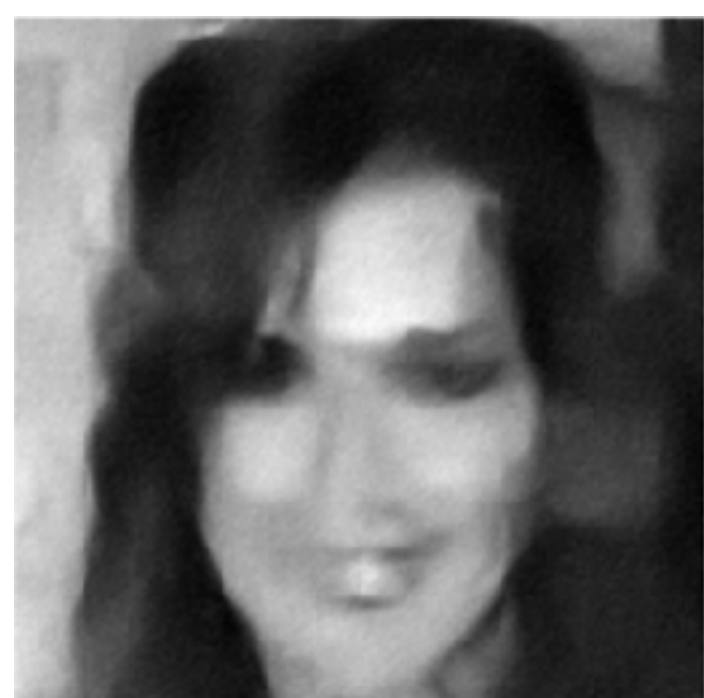
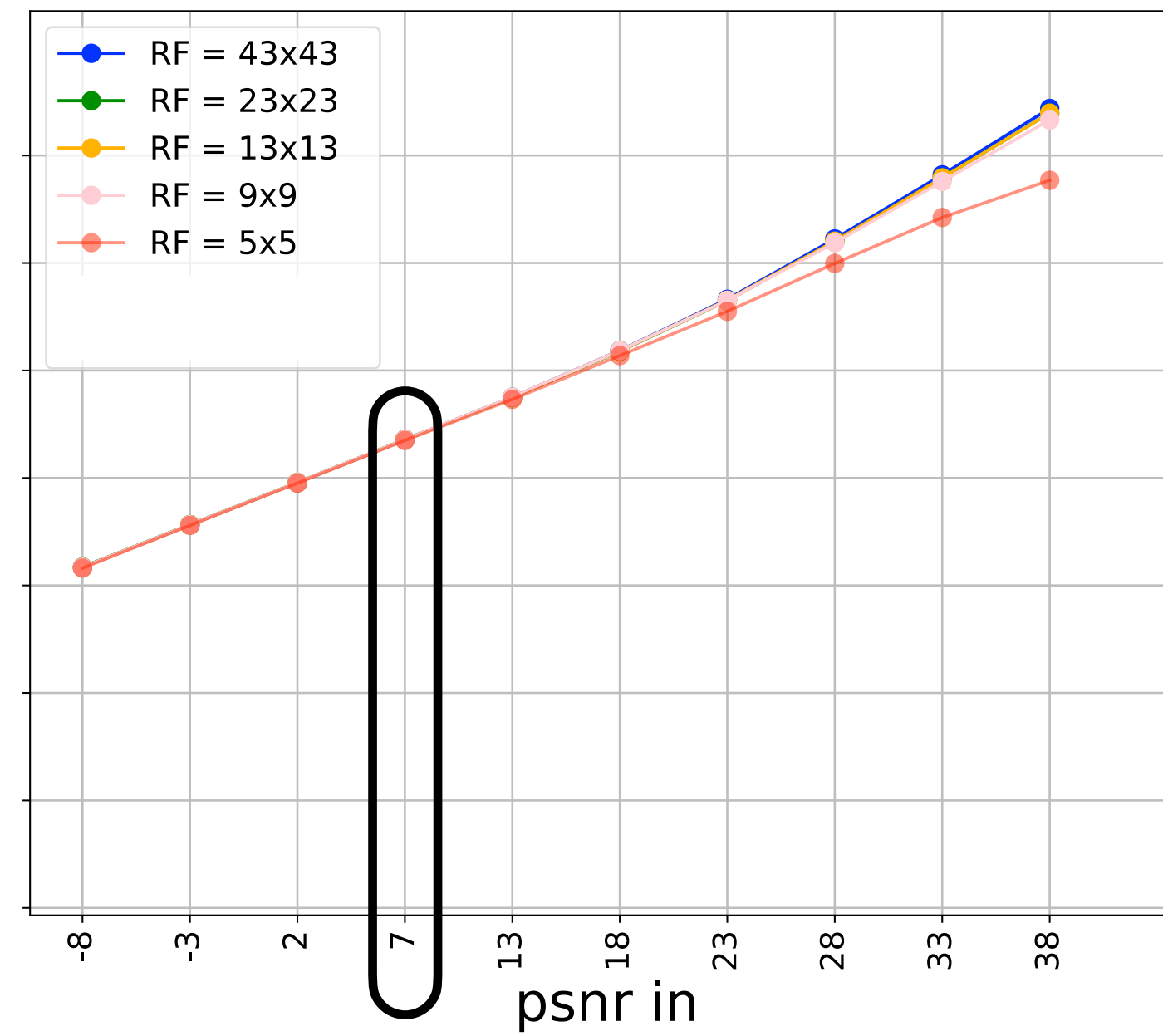


psnr in = 7



Image size: 320x320

Wavelet-domain denoiser



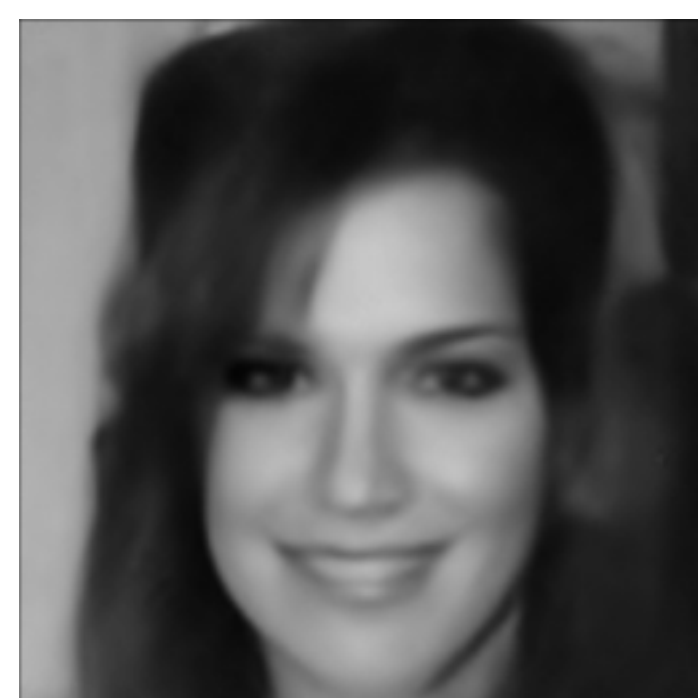
RF = 43x43



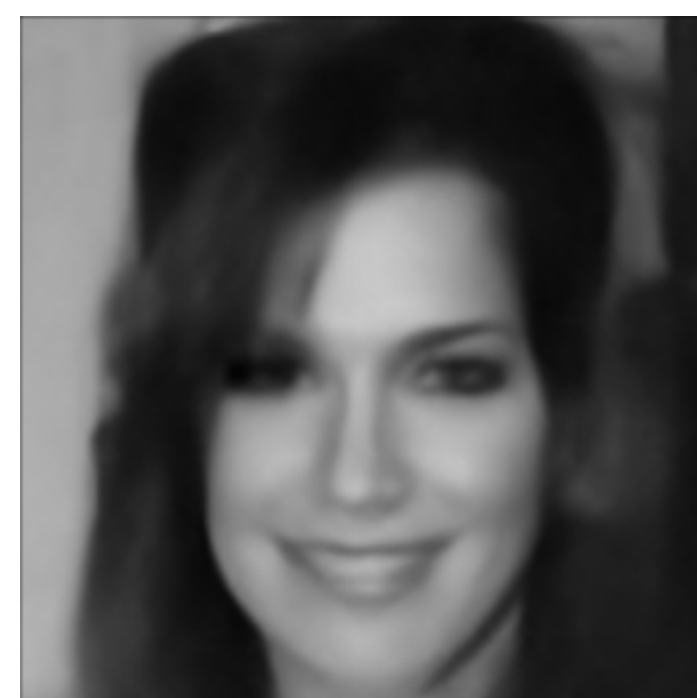
RF = 23x23



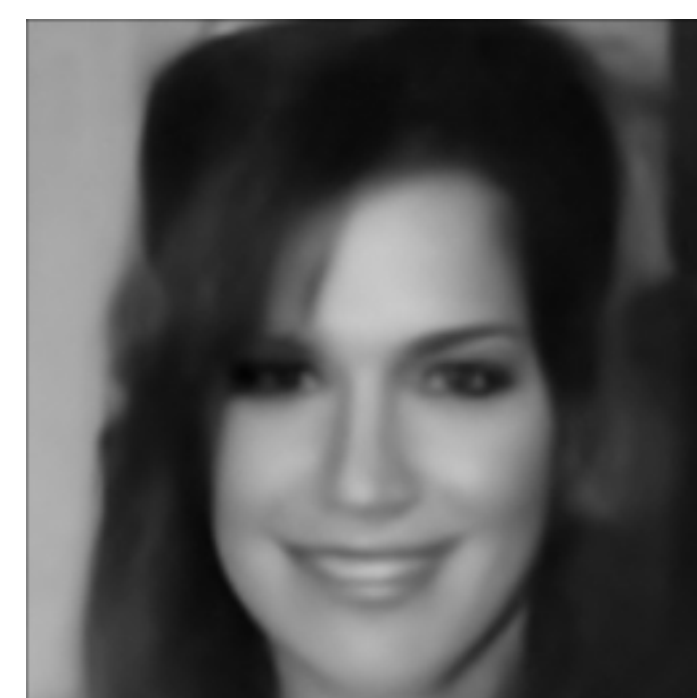
RF = 9x9



RF = 43x43

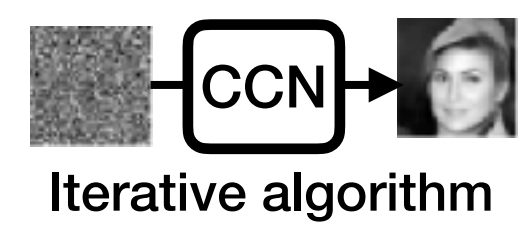


RF = 23x23

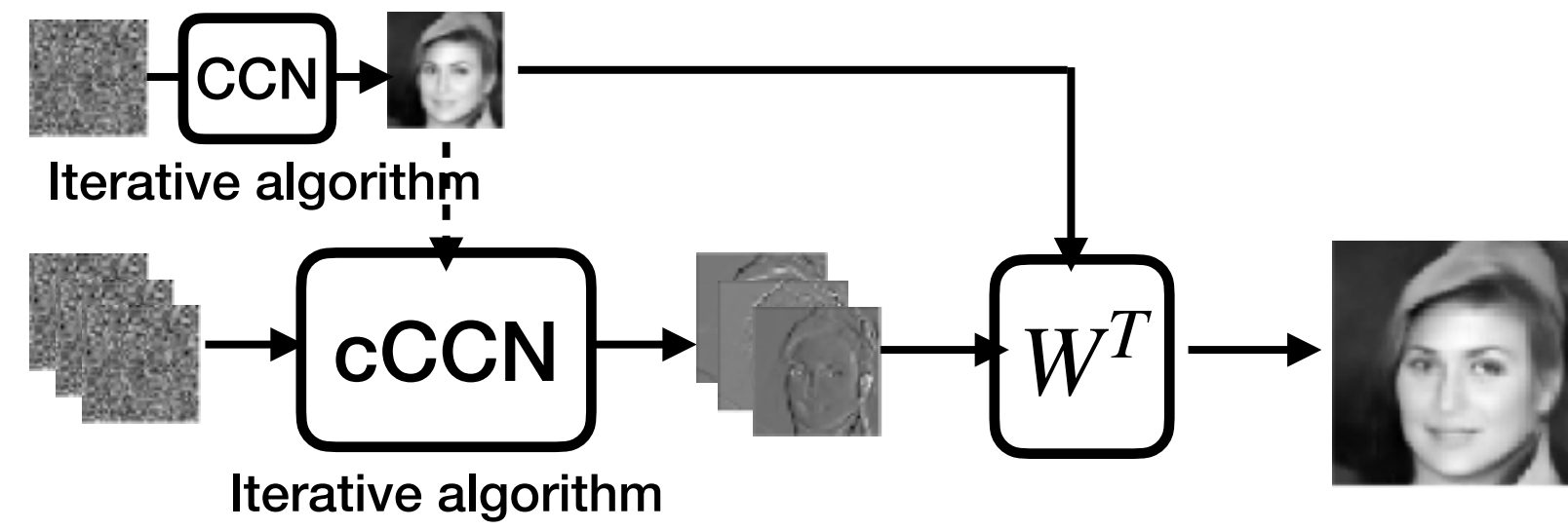


RF = 9x9

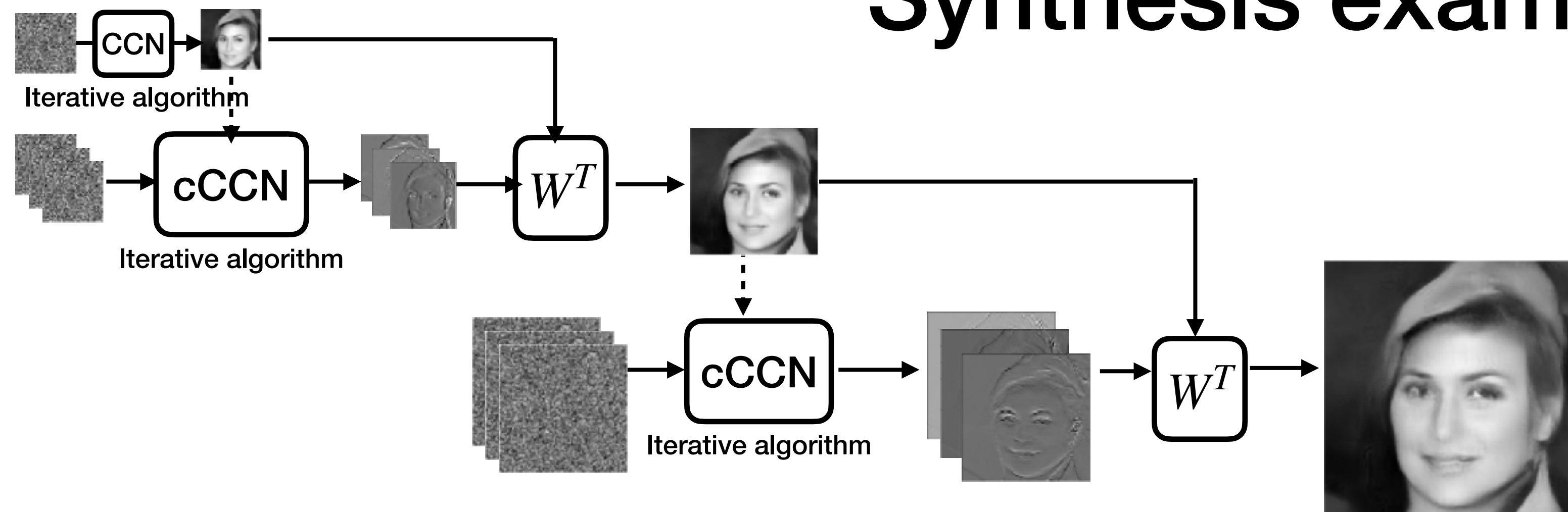
Synthesis example



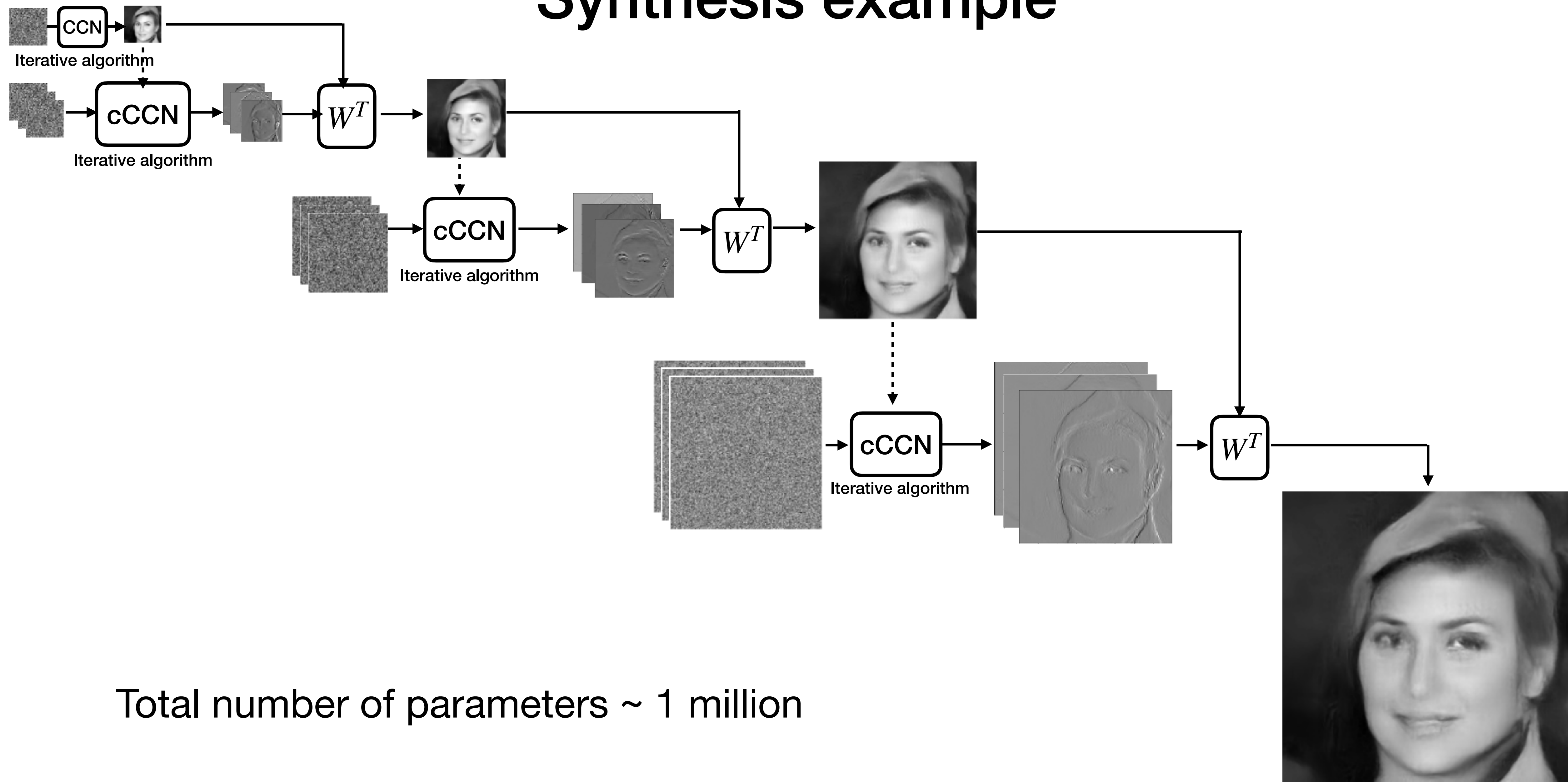
Synthesis example



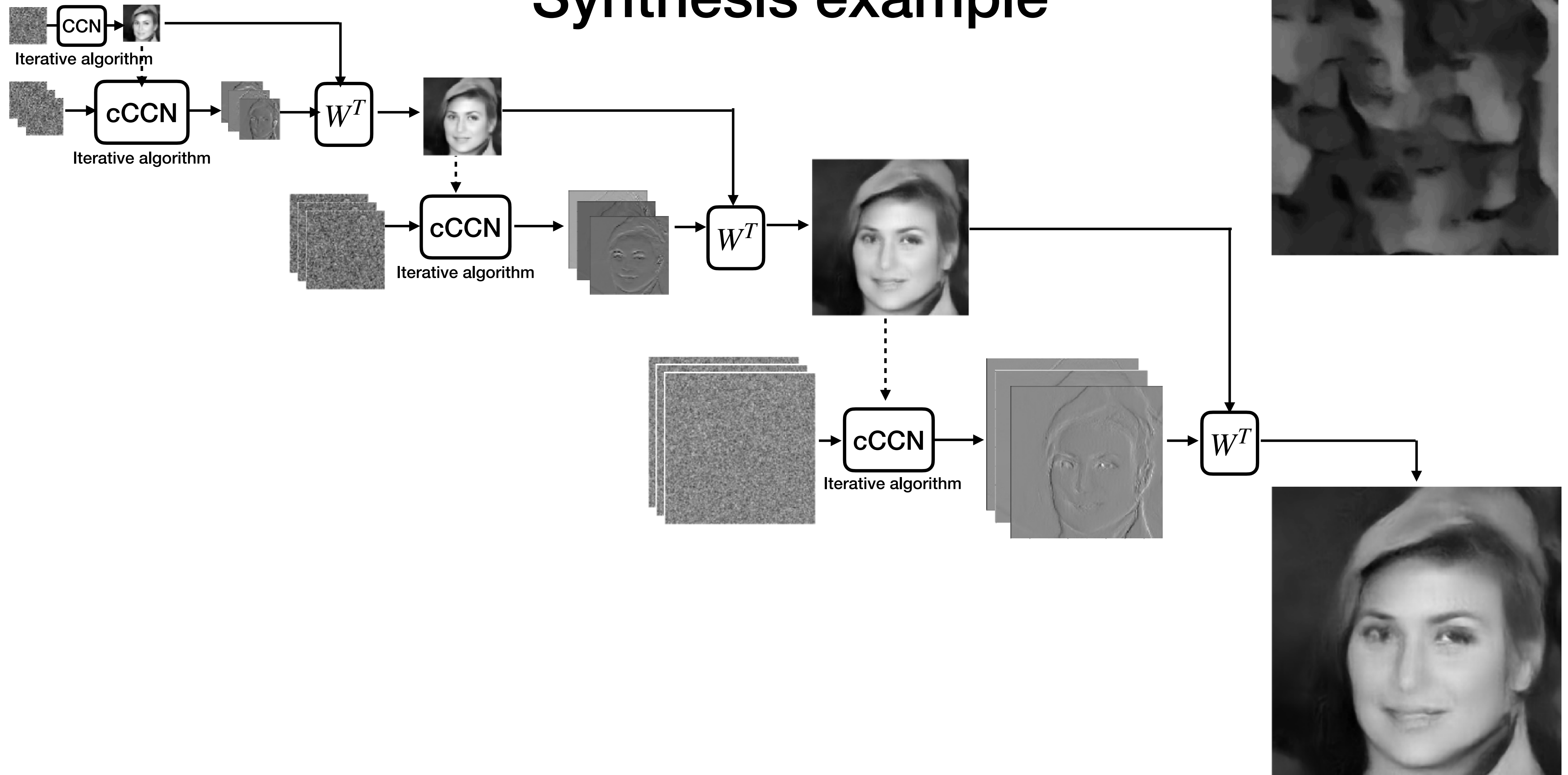
Synthesis example



Synthesis example



Synthesis example



To sum up:

- We can model probability of large images with small networks.
- The global structure is captured by a global prior over a small low-pass image.
- Details can be modeled using local (Markov) conditional probability distributions in the wavelet domain.

Thank you!

Kadkhodaie, Guth, Simoncelli, Mallat, “*Generalization in diffusion models arises from geometry-adaptive harmonic representation*”. ICLR 2024

Kadkhodaie, Guth, Mallat, & Simoncelli, “*Learning multi-scale local conditional probability models of images.*” ICLR 2023