

Introduction

Inverse problems Let x a random vector in \mathbb{R}^N , we observe

$$y = Ax + e, \quad \text{where } e \sim \mathcal{N}(0, \sigma^2 \text{Id}) \quad (1)$$

In practice, we observe a *realization* y of the random variable y :

$$y = A\bar{x} + e$$

- ▶ $A: \mathbb{R}^N \rightarrow \mathbb{R}^M$: forward linear operator (e.g., convolution, subsampling, tomography, partial Fourier (MRI), etc.)
- ▶ \bar{x} : true signal to recover (drawn from x) and e noise (drawn from e)

Learning-based methods: train a neural network ϕ_θ to map y to an estimate $\hat{x}(y) \stackrel{\text{def}}{=} \phi_\theta(By) \approx \bar{x}$, where B is a back-projection operator (e.g., A^T , A^+ or Id) and ϕ_θ is a neural network with parameters $\theta \in \Theta$.

→ **achieve SOTA performance but usually considered as black-boxes.**

Our contribution: a theoretical framework that relates CNNs to the MMSE estimators, deriving closed-form expression of the theoretical optimal solution, which is shown numerically to match the neural network outputs.

Supervised learning and MMSE estimator

Let $\mathcal{M} \stackrel{\text{def}}{=} \{\phi: \mathbb{R}^N \rightarrow \mathbb{R}^M \text{ measurable}\}$ is the set of all measurable functions.

$\mathcal{M}_\Theta \stackrel{\text{def}}{=} \{\phi_\theta: \theta \in \Theta\} \subseteq \mathcal{M}$ denote the range of a neural network architecture

$$\begin{aligned} \text{Learning: } x^* &= \phi_{\theta^*} \circ B, & \text{MMSE: } \hat{x}_{\text{MMSE}} &= \phi^* \circ B, \\ \theta^* &= \underset{\phi_\theta \in \mathcal{M}_\Theta}{\text{argmin}} \mathbb{E} [\|\phi_\theta(By) - x\|^2] & \phi^* &= \underset{\phi \in \mathcal{M}}{\text{argmin}} \mathbb{E} [\|\phi(By) - x\|^2] \end{aligned}$$

(Informal) Neural networks approximate the MMSE estimator.

Suppose that

- ▶ The network is **sufficiently expressive**: $\mathcal{M}_\Theta \approx \mathcal{M}$
- ▶ The optimization problem is solved to optimality: ϕ^* is a global minimizer of the ERM problem.

Then, ϕ^* approximates the MMSE estimator, that is $\phi^* \approx \hat{x}_{\text{MMSE}}$.

Empirical distribution

$$p_x = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \delta_x \quad \text{where } \mathcal{D} \text{ is a finite dataset}$$

Theorem 1 (Closed-form of MMSE). The MMSE estimator writes, for any $y \in \mathbb{R}^M$,

$$\hat{x}_{\text{MMSE}}(y) = \sum_{x \in \mathcal{D}} x \cdot w(x|y),$$

where $w(x|y) \propto \mathcal{N}(By; BAx, \sigma^2 BB^T)$.

Implication: a neural network with *sufficient capacity* can learn to implement the above formula.

- ▶ The formula is exact for the empirical MMSE estimator, which is the optimal solution of the ERM problem.
- ▶ In practice, we **could not reproduce** this behavior on real datasets, even with large MLPs or CNNs.
- ▶ What is **missing** from this formula? Our answer: **Inductive biases** of the neural network architecture!

CNNs and Inductive Biases

Two well-known inductive biases of CNNs

- ▶ **Translation equivariance:** a fully convolutional network is *equivariant* to translations.
- ▶ **Locality:** finite receptive fields.

Reducing the search space of the MMSE estimator. Let \mathcal{M}_{LE} denote the set of all measurable function that are *local and translation equivariant*.

$$\mathcal{M}_\Theta \subseteq \mathcal{M}_{\text{LE}} \subseteq \mathcal{M}$$

Learning: $x^* = \phi_{\theta^*} \circ B$,

$$\theta^* = \underset{\phi_\theta \in \mathcal{M}_\Theta}{\text{argmin}} \mathbb{E} [\|\phi_\theta(By) - x\|^2]$$

Constrained MMSE: $\hat{x}_{\text{LE-MMSE}} = \phi^* \circ B$,

$$\phi^* = \underset{\phi \in \mathcal{M}_{\text{LE}}}{\text{argmin}} \mathbb{E} [\|\phi(By) - x\|^2]$$

Theorem 2 (Closed-form of LE-MMSE). Suppose that the N matrices $Q_n \stackrel{\text{def}}{=} \Pi_n B \in \mathbb{R}^{P \times M}$ have the *same rank* $r > 0$ for any n^a . The LE-MMSE estimator $\hat{x}_{\text{LE-MMSE}}$ admits, for any $y \in \mathbb{R}^M$ the following expression, defined pixel-wise for each pixel n' :

$$\hat{x}_{\text{LE-MMSE}}(y)[n'] = \sum_{x \in \mathcal{D}} \sum_{n=1}^N x[n] \cdot w_{n',n}(x|y),$$

where $w_{n',n}(x|y) \propto \mathcal{N}(Q_n y; Q_n A x, \sigma^2 Q_n Q_n^T)$.

^aThis assumption can be relaxed: a general formula by rank stratification is also available.

Key insights:

- ▶ The value at pixel n' of the LE-MMSE estimator is a weighted average of **all** the pixels in the training images. The reconstruction is a **patchwork** of training patches.
- ▶ Can produce images outside $\text{conv}(\mathcal{D})$.
- ▶ It is the **orthogonal projection** (in L^2 sense) of the empirical MMSE estimator onto \mathcal{M}_{LE} .

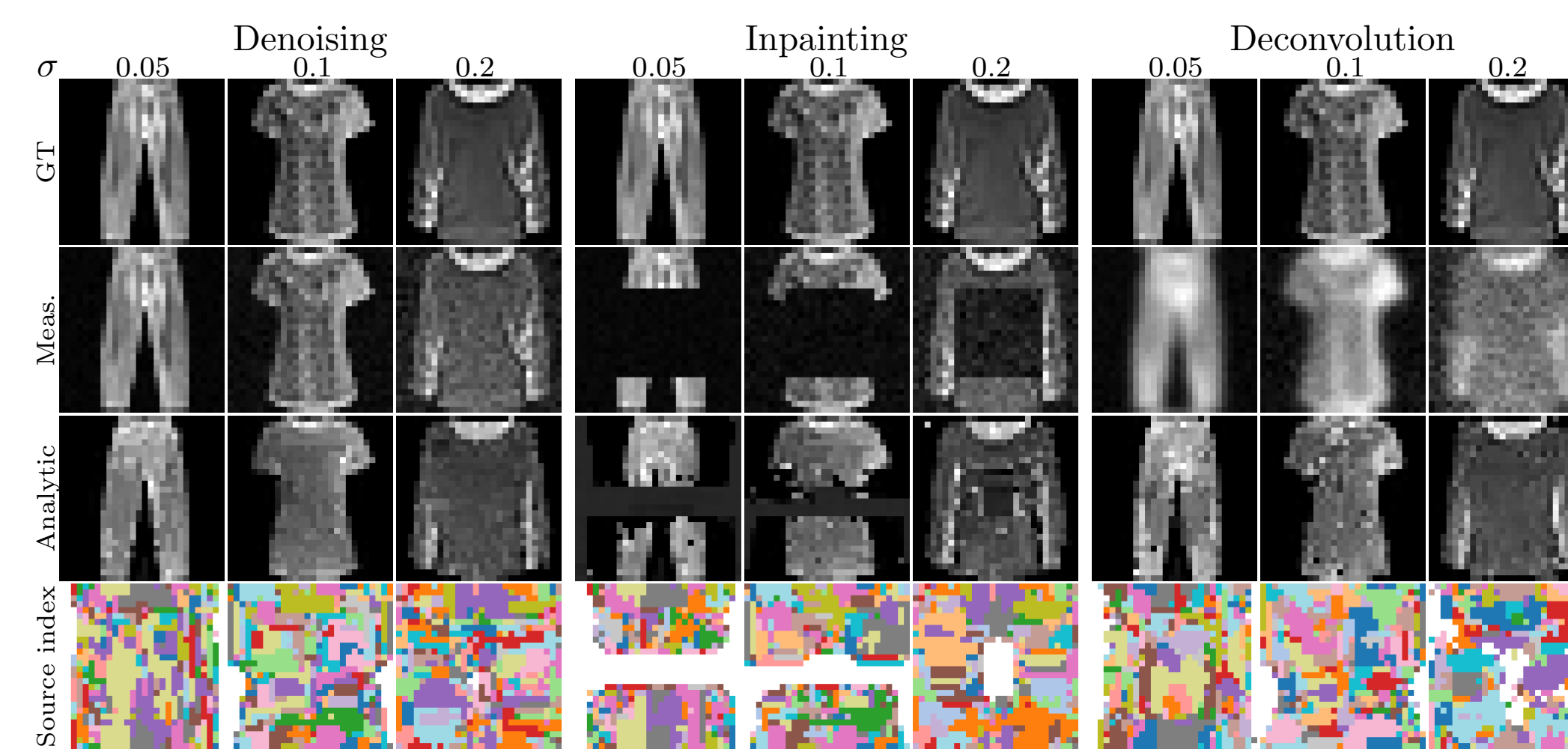


Figure: Pixel index of the reconstruction: contiguous regions coming from the same source image

The role of the back-projection operator B : estimator's variance

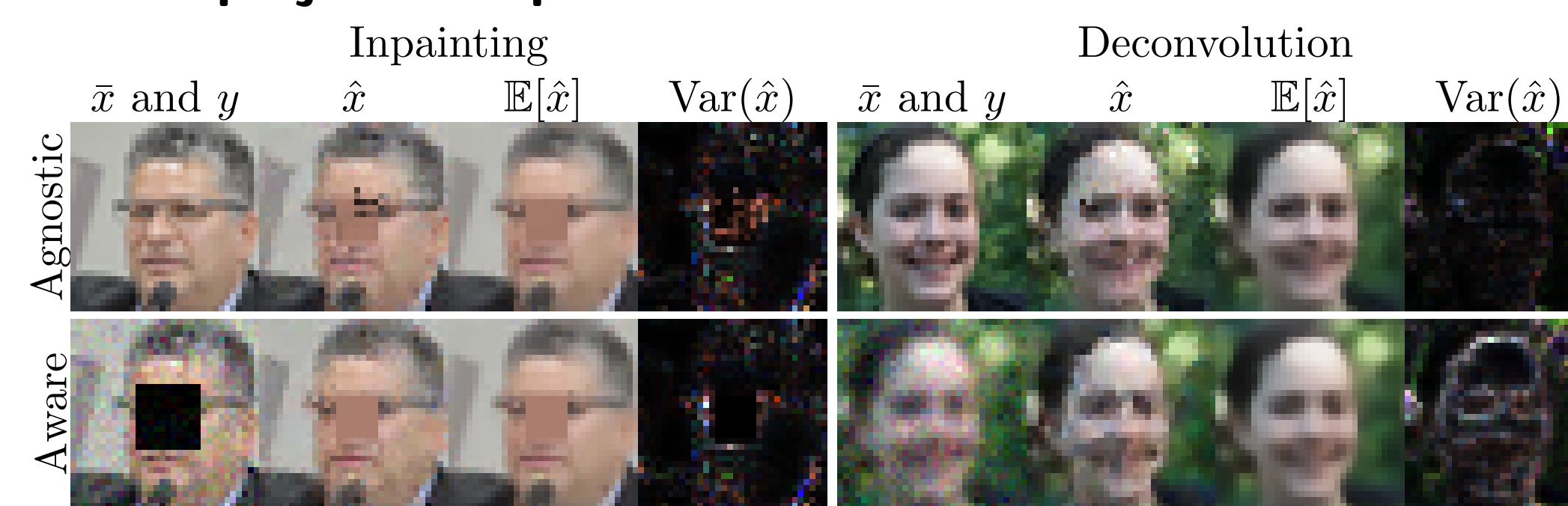


Figure: Physics-aware vs physics-agnostic solvers: pixel-wise variance.

Physics-aware ($B = A^+$, bottom) has lower variance for inpainting (left), while physics-agnostic ($B = I$, top) has lower variance for deconvolution (right).

Numerical Results

Quantitative comparison: our closed-form matches the neural network outputs

	σ	Denoising			Inpainting			Deconvolution			Denoising			Inpainting			Deconvolution			PSNR (dB) - Neural vs Analytic
		0.05	0.20	0.80	0.05	0.20	0.80	0.05	0.20	0.80	0.05	0.20	0.80	0.05	0.20	0.80	0.05	0.20	0.80	
FFHQ-32	UNet2D	28.9	25.9	29.8	27.6	25.7	27.7	24.0	31.6	36.5	23.8	25.9	29.8	24.4	27.0	27.6	25.3	31.9	36.5	40.0 37.5 35.0 32.5 30.0 27.5 25.0 22.5
	ResNet	29.3	27.3	36.5	24.5	23.6	30.2	23.7	32.2	40.2	25.0	28.9	36.7	23.4	25.3	30.4	24.8	32.4	40.5	
	PatchMLP	28.1	27.5	32.1	23.2	23.6	24.6	23.6	30.6	33.8	24.9	29.0	32.3	22.9	25.4	24.3	24.6	30.8	33.9	
CIFAR10	UNet2D	30.8	31.0	34.1	22.5	24.1	25.0	26.1	34.2	36.7	25.9	32.0	33.9	21.8	24.2	24.9	26.0	34.6	36.8	
	ResNet	30.1	30.2	37.7	24.8	25.6	30.5	24.8	32.8	39.6	26.0	31.1	37.7	23.9	25.9	30.6	25.8	33.4	39.9	
	PatchMLP	30.2	29.9	31.8	24.6	26.1	25.4	25.4	32.5	32.7	25.5	30.7	32.0	23.4	26.6	26.0	25.4	32.7	32.9	
FashionMNIST	UNet2D	32.5	31.1	26.2	25.0	25.9	26.6	26.5	29.7	27.0	28.9	31.5	25.9	24.5	26.8	26.8	27.7	29.3	26.7	
	ResNet	31.4	28.9	36.3	21.7	24.3	30.7	24.7	30.9	41.6	28.0	29.6	36.1	21.2	24.8	30.6	25.9	30.9	41.3	
	PatchMLP	31.6	31.4	32.4	24.8	26.7	25.7	26.9	32.2	33.4	28.4	31.9	32.6	24.6	26.9	25.9	28.5	32.0	33.2	

Figure: PSNR between our closed-form formula and neural network outputs

Qualitative comparison: on various inverse problems and datasets

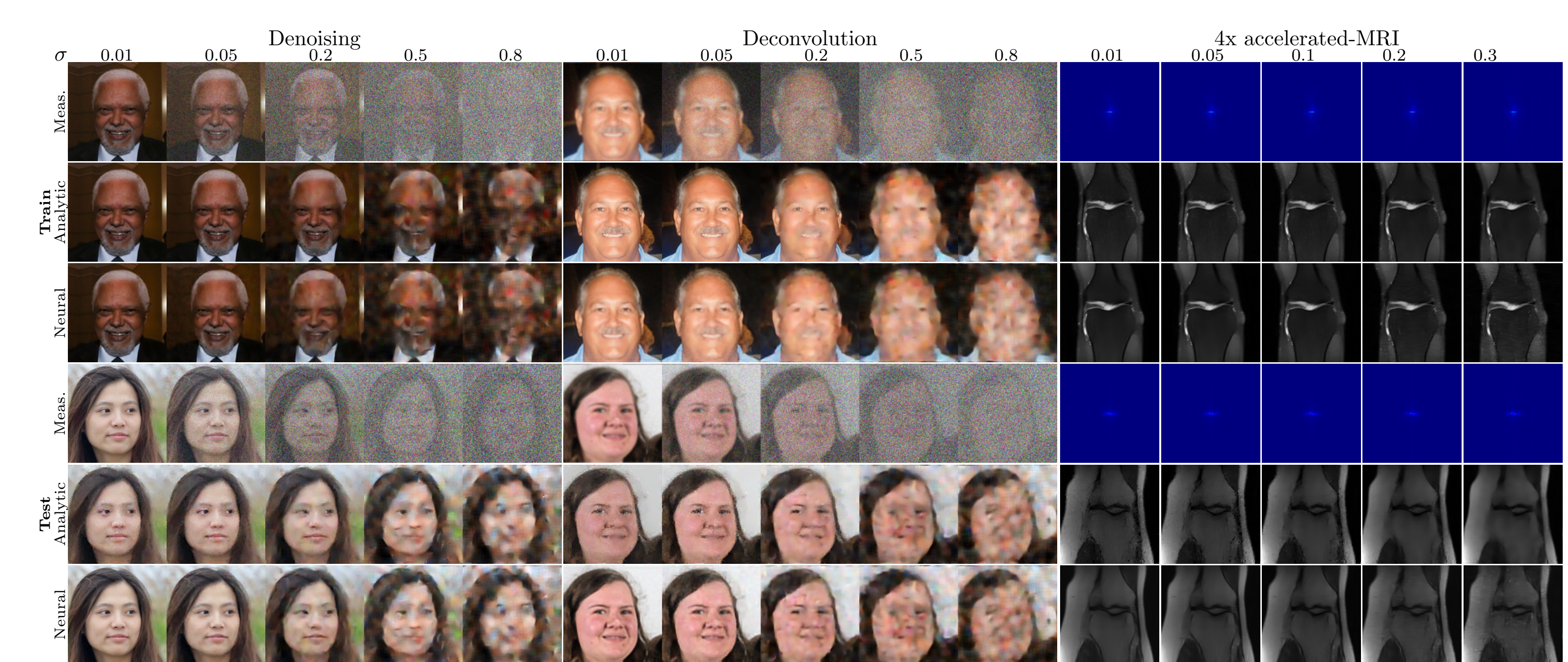
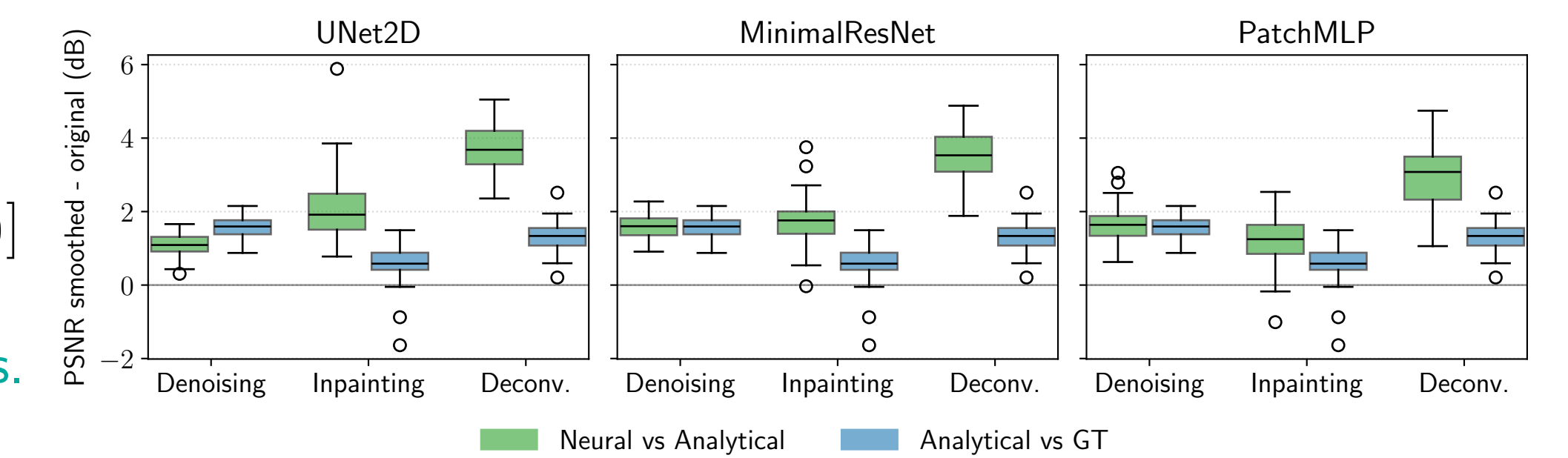


Figure: Qualitative comparison between our closed-form formula and neural network outputs.

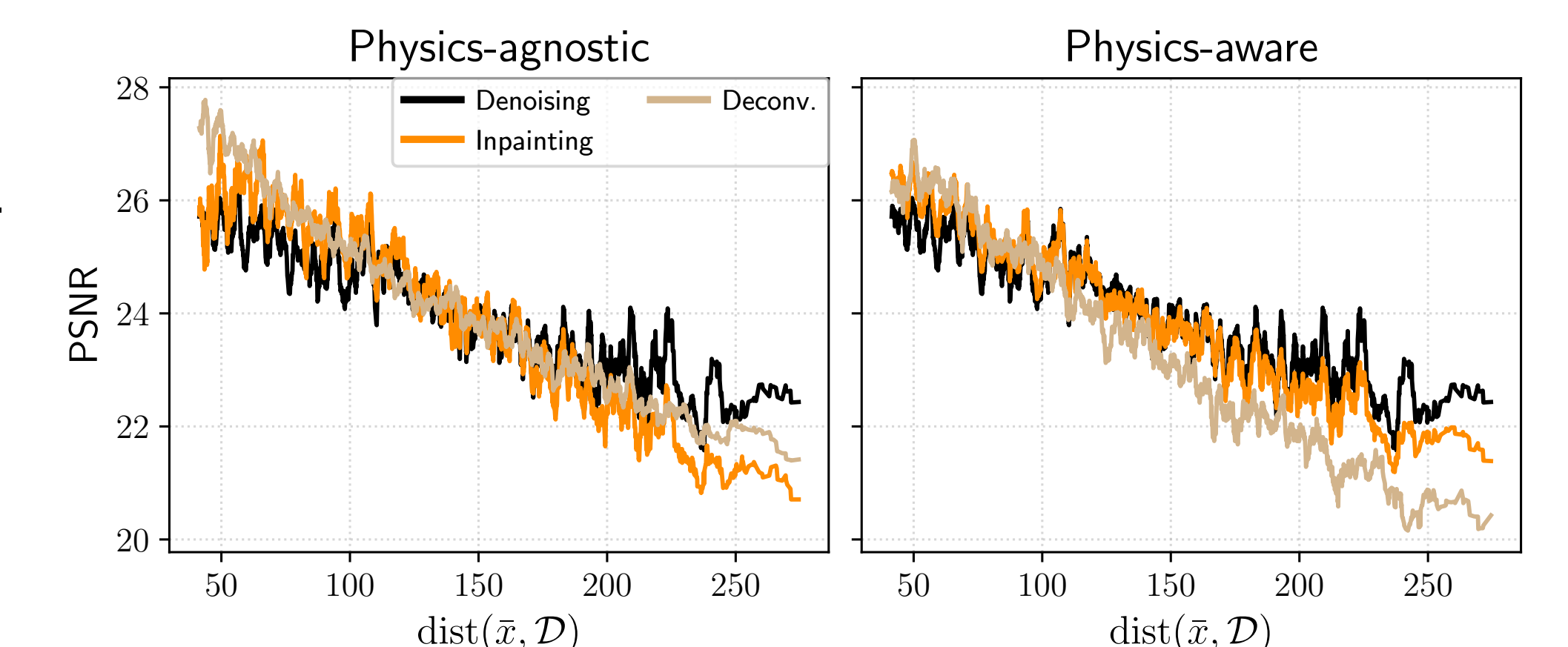
Discussion

Spectral bias: CNNs exhibit a bias towards low frequencies, smoother images. Setting

$\hat{x}_{\text{LE-MMSE}}^{\text{smooth}}(y) \stackrel{\text{def}}{=} \mathbb{E} [\hat{x}_{\text{LE-MMSE}}(y + \varepsilon \cdot z)]$ improves the both reconstruction and the match with neural networks.



Limitation: The alignment **degrades** in low-density regions (eg. out-of-distribution, low-noise levels or large patch sizes).



Next extensions

- ▶ Toward recovery guarantee: stability and reconstruction guarantees from the closed-form formula.
- ▶ Beyond Gaussian noise and MSE loss and toward more inductive biases.