



UNSURE: Unknown Noise level Stein's Unbiased Risk Estimator

Julián Tachella, CNRS, École Normale Supérieure de Lyon

Joint work with Mike Davies (University of Edinburgh) and Laurent Jacques (UCLouvain)

Inverse Problems

Goal: recover signal x from y

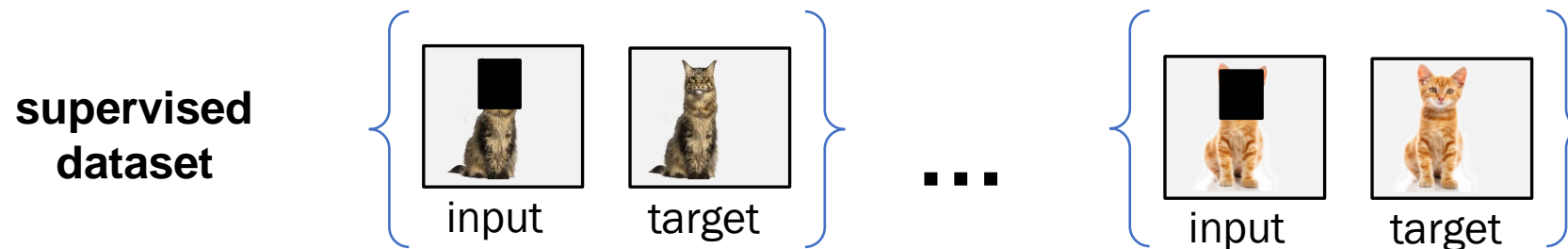
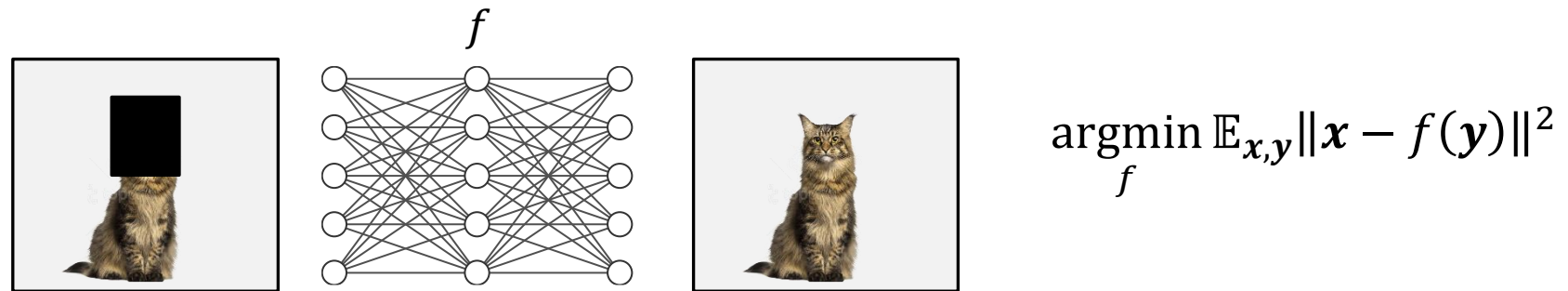
$$y = Ax + \epsilon$$

Diagram illustrating the components of the inverse problem equation $y = Ax + \epsilon$:

- y : MEASUREMENTS $\in \mathbb{R}^m$
- A : MEASUREMENT OPERATOR $\in \mathbb{R}^{m \times n}$
- x : SIGNAL $\in \mathbb{R}^n$
- ϵ : NOISE $\in \mathbb{R}^m$

Learning approach

Idea: use training pairs of signals and measurements to directly learn the inversion function



Learning approach

Advantages:

- State-of-the-art reconstructions
- Once trained, f_θ is easy to evaluate

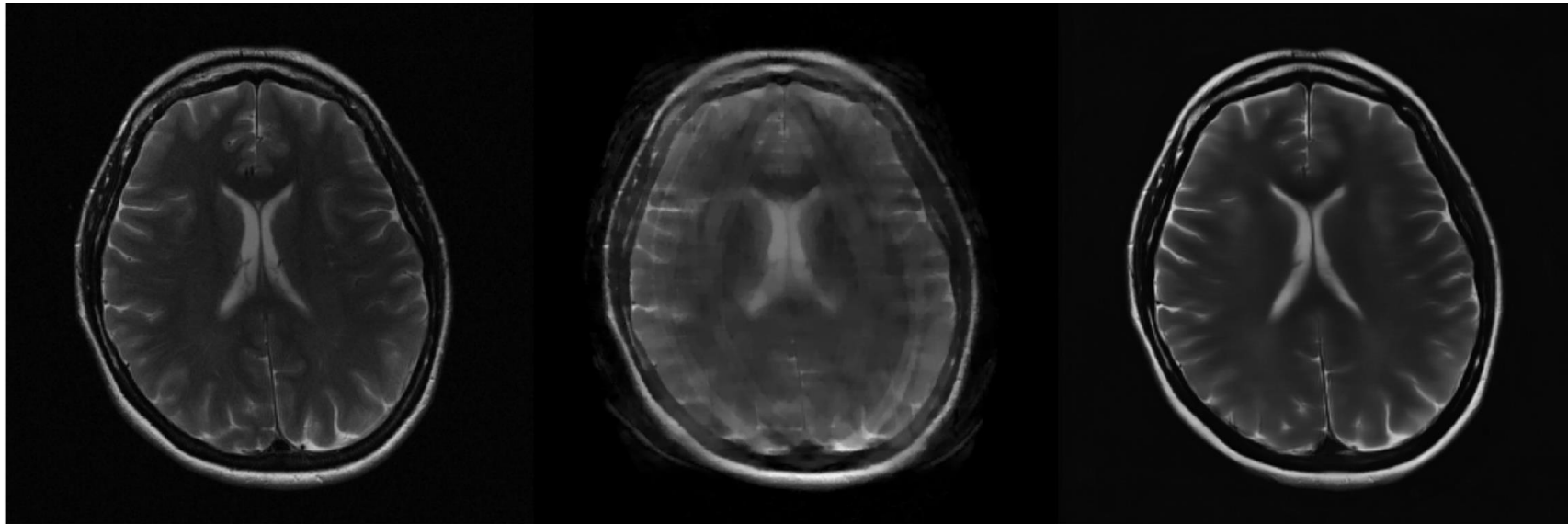
fastMRI

Accelerating MR Imaging with AI

Ground-truth

Total variation
(28.2 dB)

Deep network
(**34.5 dB**)

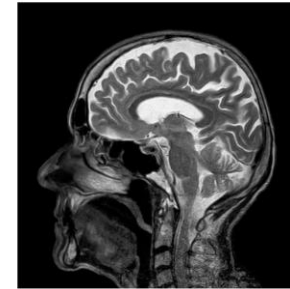
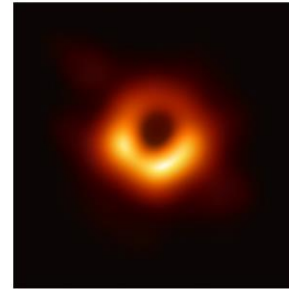


x8 accelerated MRI [Zbontar et al., 2019]

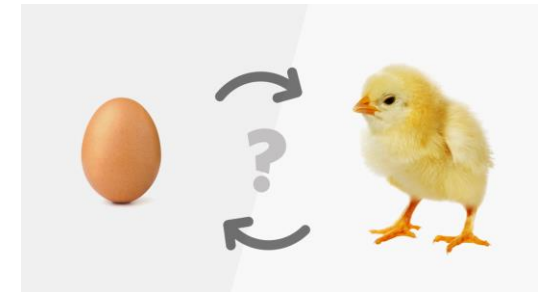
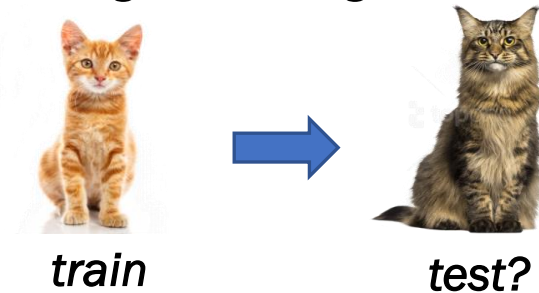
Learning approach

Main disadvantage: Obtaining training signals x_i can be expensive or impossible.

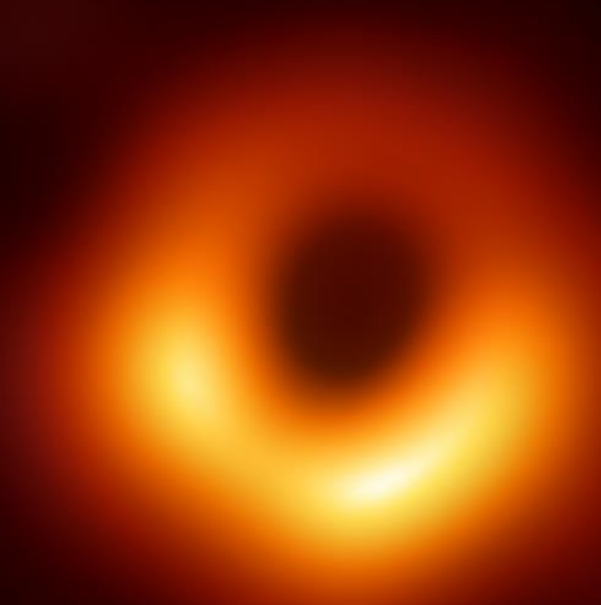
- Medical and scientific imaging



- Only solves inverse problems which we already know what to expect
- Risk of training with signals from a different distribution

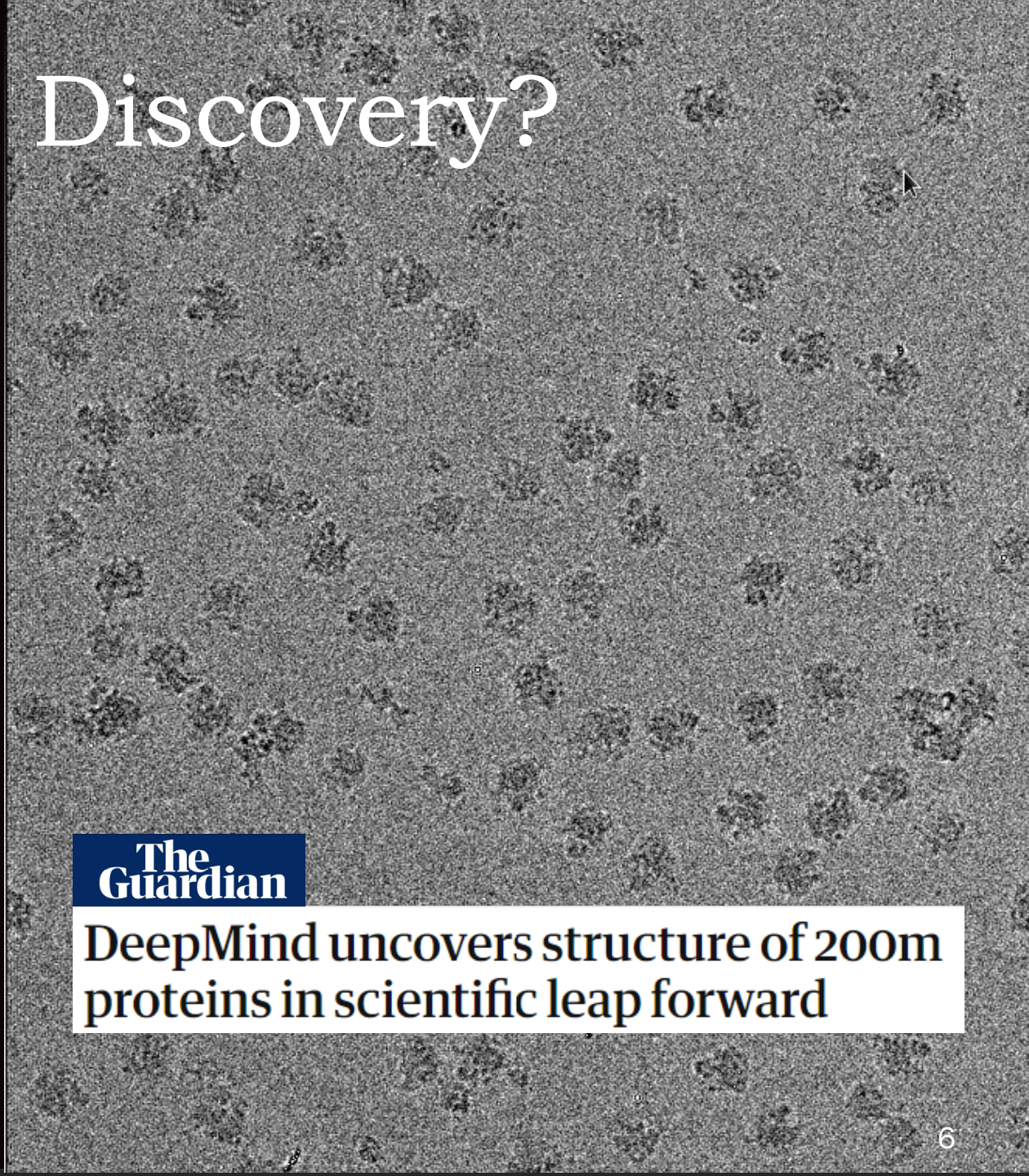


AI for Knowledge Discovery?



The
Guardian

Black hole picture captured for first time in space breakthrough



The
Guardian

DeepMind uncovers structure of 200m proteins in scientific leap forward

Denoising problems

In this first part, we will focus on ‘denoising’ problems

$$\mathbf{y} = A\mathbf{x} + \boldsymbol{\epsilon}$$


where $A \in \mathbb{R}^{m \times n}$ is invertible (and thus $m \geq n$).

- We focus on $A = I$ for simplicity.
- All methods in this part can be extended to any invertible A .

MMSE estimators

We focus on ℓ_2 loss and minimum mean squared error estimators (MMSE)

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y}} \|\mathbf{x} - f(\mathbf{y})\|^2$$

 $f^*(\mathbf{y}) = \mathbb{E}\{\mathbf{x}|\mathbf{y}\}$

- Other estimators might be preferred, eg. perceptual [Blau and Michaeli, 2018]

Unsupervised Risk Estimators

Supervised loss

$$\mathcal{L}_{\text{SUP}}(\mathbf{x}, \mathbf{y}, f) = \|\mathbf{x} - f(\mathbf{y})\|^2 = \underbrace{\|\mathbf{y} - f(\mathbf{y})\|^2}_{\text{Measurement consistency}} + \underbrace{2f(\mathbf{y})^\top (\mathbf{y} - \mathbf{x})}_{\text{key term to approximate!} = f(\mathbf{y})^\top \boldsymbol{\epsilon}} + \text{const.}$$

Goal: build a self-supervised loss $\mathcal{L}_{\text{self}}$ such that

$$\mathbb{E}_{\mathbf{y}} \mathcal{L}_{\text{self}}(\mathbf{y}, f) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} \mathcal{L}_{\text{SUP}}(\mathbf{x}, \mathbf{y}, f) + \text{const.}$$

Stein's Unbiased Risk Estimator

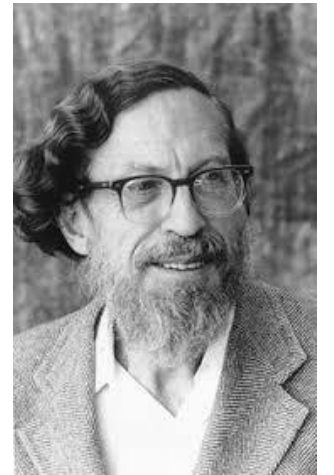
- **Stein's lemma** [Stein 1974] : Let $\mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, I\sigma^2)$, f be weakly differentiable, then

$$\mathbb{E}_{\mathbf{y}|\mathbf{x}} (\mathbf{y} - \mathbf{x})^\top f(\mathbf{y}) = \sigma^2 \mathbb{E}_{\mathbf{y}|\mathbf{x}} \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y})$$

$$\min_f \mathbb{E}_{\mathbf{y}} \|\mathbf{y} - f(\mathbf{y})\|^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y})$$

Measurement
consistency

Degrees of freedom [Efron, 2004]



- MMSE estimator $f^*(\mathbf{y}) = \mathbb{E}\{\mathbf{x}|\mathbf{y}\}$

Stein's Unbiased Risk Estimator

- **Hudson's lemma** [Hudson 1978] : Let $\mathbf{y}|\mathbf{x}$ exponential family, f be weakly differentiable, then

$$\mathbb{E}_{\mathbf{y}|\mathbf{x}} (\mathbf{y} - \mathbf{x})^\top f(\mathbf{y}) = \mathbb{E}_{\mathbf{y}|\mathbf{x}} \sum_i a(y_i) \frac{\delta f_i}{\delta y_i}(\mathbf{y})$$

- Gaussian noise: $a(y) = \sigma^2$
- Poisson noise: $a(y) \approx y$

Tweedie's Formula

The solution to SURE is **Tweedie's Formula**

$$\min_f \mathbb{E}_{\mathbf{y}} || \mathbf{y} - f(\mathbf{y}) ||^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y})$$

$$\min_f \mathbb{E}_{\mathbf{y}} || \mathbf{y} - f(\mathbf{y}) ||^2 - 2\sigma^2 \sum_i f_i(\mathbf{y}) \frac{\delta \log p_{\mathbf{y}}(\mathbf{y})}{\delta y_i}$$

$$\min_f \mathbb{E}_{\mathbf{y}} || f(\mathbf{y}) - \mathbf{y} - \sigma^2 \nabla \log p_{\mathbf{y}}(\mathbf{y}) ||^2$$



Integration by parts



Complete squares

$$\Rightarrow f(\mathbf{y}) = \mathbf{y} + \sigma^2 \nabla \log p_{\mathbf{y}}(\mathbf{y})$$

- **Noise2Score** [Kim and Ye, 2021] learns $\nabla \log p_{\mathbf{y}}(\mathbf{y})$ from noisy data + denoises with Tweedie.
- Key formula behind diffusion models, which can be trained self-supervised [Daras et al., 2024]

Cross-Validation Methods

What happens if we only know that $p(\mathbf{y}|\mathbf{x}) = \prod p(y_i|x_i)$?

$$\min_f ||\mathbf{y} - f(\mathbf{y})||^2 \text{ subject to } \frac{\delta f_i}{\delta y_i}(\mathbf{y}) = 0 \forall i, \forall \mathbf{y}$$

- SURE's perspective:

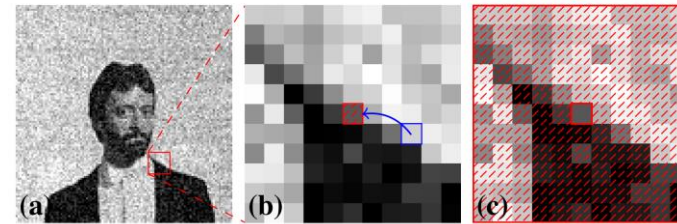
$$\min_f ||\mathbf{y} - f(\mathbf{y})||^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y})$$

- These methods are not MMSE optimal!
- f_i shouldn't depend on y_i : training or architecture

Cross-Validation Methods

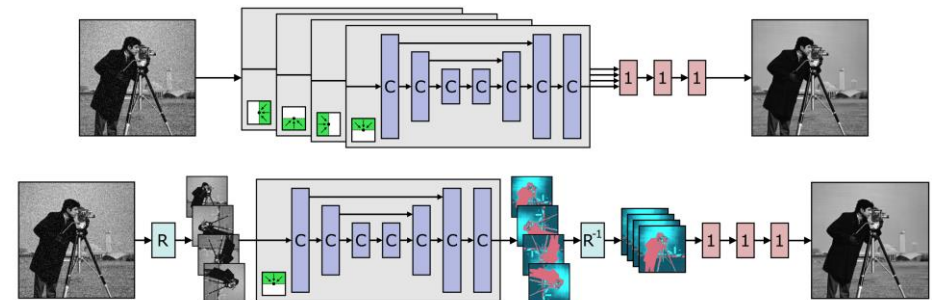
Noise2Void [Krull et al., 2019], **Noise2Self** [Batson, 2019], **Neighbor2Neighbor** [Huang, 2023]

- During training flip centre pixel
- Computes loss only on flipped pixels

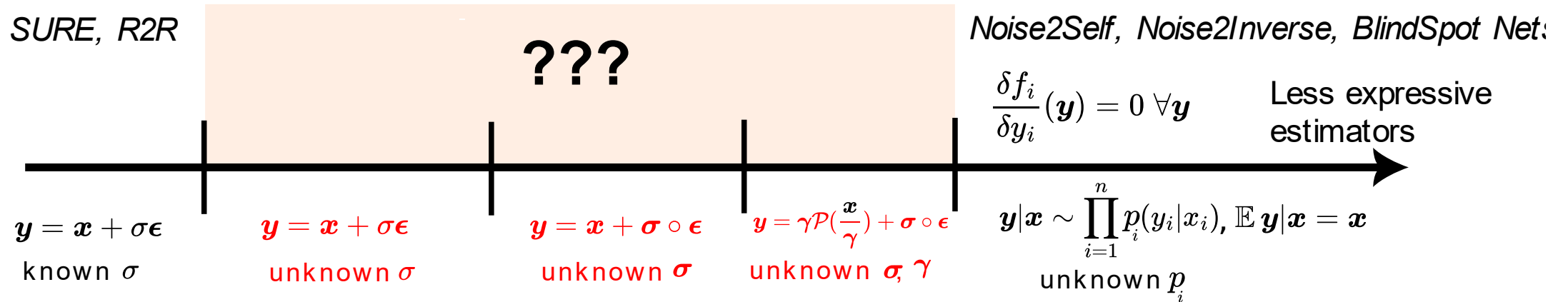


Blind spot networks [Laine et al., 2019], [Lee et al., 2022]

- Convolutional architecture that doesn't 'see' centre pixel by construction



Are We Missing Something?



Zero Expected Divergence

Assumption: Let $\mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, I\sigma^2)$.

$$\mathcal{L}_{\text{ZED}}(\mathbf{y}, f) = \|\mathbf{y} - f(\mathbf{y})\|^2 \text{ subject to } \mathbb{E}_{\mathbf{y}} \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y}) = 0$$

- SURE's perspective:

$$\mathcal{L}_{\text{ZED}}(\mathbf{y}, f) = \|\mathbf{y} - f(\mathbf{y})\|^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y})$$

- Not MMSE optimal (but almost)

ZED Denoisers

Assumption: Let $\mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, I\sigma^2)$.

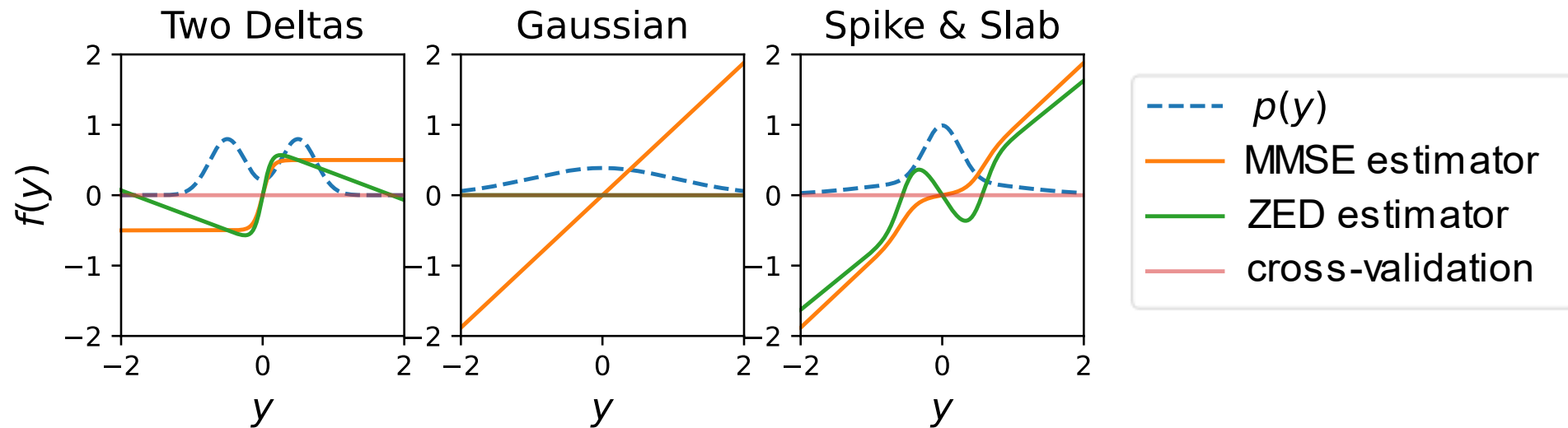
$$\min_f \mathbb{E}_{\mathbf{y}} \|\mathbf{y} - f(\mathbf{y})\|^2 \text{ subject to } \sum_i \mathbb{E}_{\mathbf{y}} \frac{\delta f_i}{\delta y_i}(\mathbf{y}) = 0$$

$$\max_{\eta} \min_f \mathbb{E}_{\mathbf{y}} \|\mathbf{y} - f(\mathbf{y})\|^2 + 2\eta \sum_i \mathbb{E}_{\mathbf{y}} \frac{\delta f_i}{\delta y_i}(\mathbf{y})$$

$$\Rightarrow f^{\text{ZED}}(\mathbf{y}) = \mathbf{y} + \hat{\eta} \nabla \log p_{\mathbf{y}}(\mathbf{y}) \quad \hat{\eta} = \left(\frac{1}{n} \mathbb{E}_{\mathbf{y}} \|\nabla \log p_{\mathbf{y}}(\mathbf{y})\|^2 \right)^{-1}$$

ZED Denoisers

Examples: separable prior $p(x) = \prod_i q(x_i)$



ZED Denoisers

How far is f^{ZED} from MMSE optimality?

Theorem:

$$\frac{1}{n} \mathbb{E}_{\mathbf{x}, \mathbf{y}} || f^{\text{ZED}}(\mathbf{y}) - \mathbf{x} || = \sigma^2 \left(\frac{1}{1 - \frac{\text{MMSE}}{\sigma^2}} - 1 \right) \approx \text{MMSE} + \underbrace{\frac{\text{MMSE}^2}{\sigma^2}}$$

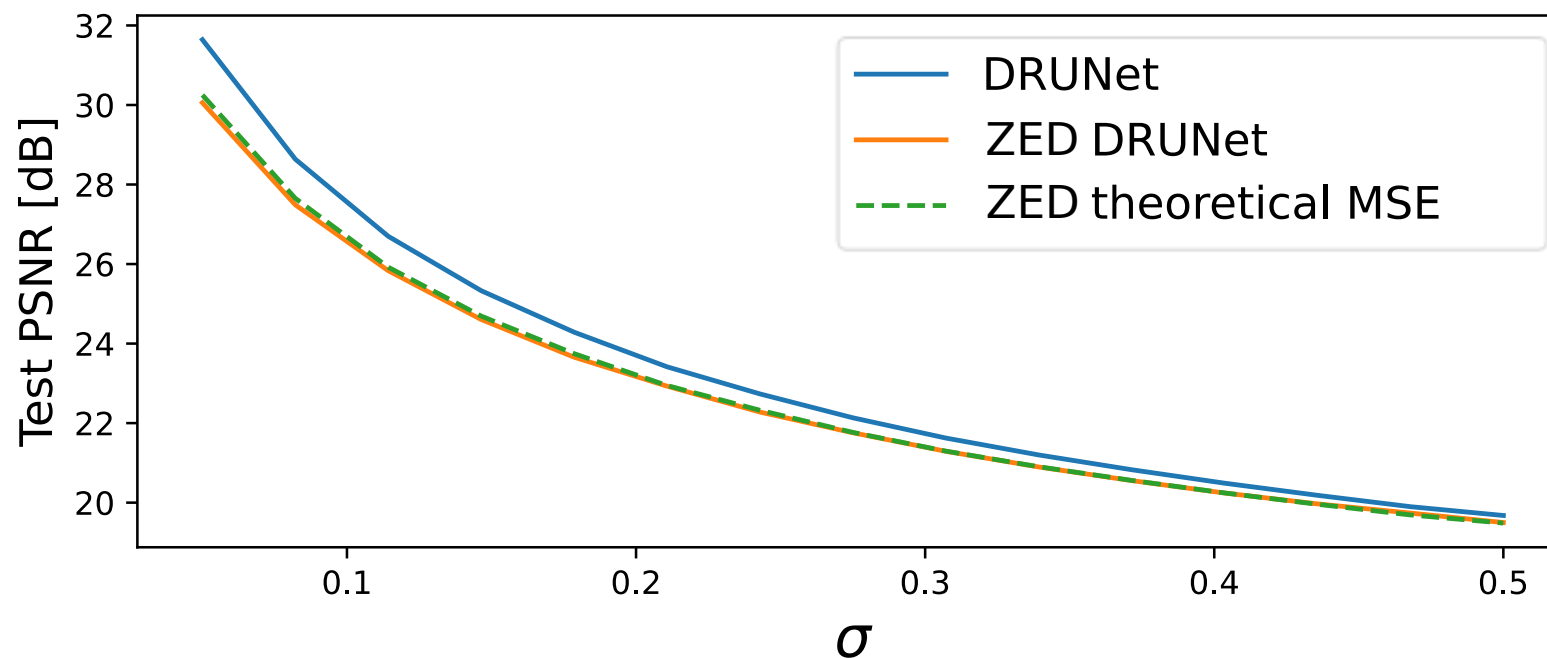
If supp. $p(\mathbf{x})$ is k -dimensional

$$\approx \left(\frac{k\sigma}{n} \right)^2$$

- $\hat{\eta} = \frac{\sigma^2}{1 - \frac{\text{MMSE}}{\sigma^2}}$ which serves as an estimator of the noise level

ZED Denoisers

Theorem is well verified by a pre-trained denoiser



Poisson-Gaussian Noise

- Noise model $\mathbf{y} = \gamma \mathbf{z} + \sigma \epsilon$, $\mathbf{z} \sim \mathcal{P}(\mathbf{x}/\gamma)$ unknown σ, γ

$$\min_f \mathbb{E}_{\mathbf{y}} \|\mathbf{y} - f(\mathbf{y})\|^2 \text{ subject to } \sum_i \mathbb{E}_{\mathbf{y}} \frac{\delta f_i}{\delta y_i}(\mathbf{y}) = 0 \text{ and } \sum_i \mathbb{E}_{\mathbf{y}} y_i \frac{\delta f_i}{\delta y_i}(\mathbf{y}) = 0$$

$$\max_{\eta, \gamma} \min_f \mathbb{E}_{\mathbf{y}} \|\mathbf{y} - f(\mathbf{y})\|^2 + 2\eta \operatorname{div} f(\mathbf{y}) + 2\gamma \nabla f(\mathbf{y})^\top \mathbf{y}$$

- Closed-form solution for f in the paper

Unknown Correlations

Non-isotropic noise $\mathbf{y}|\mathbf{x} \sim N(\mathbf{0}, \Sigma)$ with unknown Σ

Define a set of possible covariances $R = \{\Sigma_{\boldsymbol{\eta}} \in \mathcal{S}^{n \times n} : \boldsymbol{\eta} \in \mathbb{R}^p\}$

$$\max_{\boldsymbol{\eta}} \min_f \mathbb{E}_{\mathbf{y}} \|\mathbf{y} - f(\mathbf{y})\|^2 + 2 \operatorname{tr}(\Sigma_{\boldsymbol{\eta}} \frac{\delta f}{\delta \mathbf{y}})$$

$$\longrightarrow f(\mathbf{y}) = \mathbf{y} + \Sigma_{\hat{\boldsymbol{\eta}}} \nabla \log p_{\mathbf{y}}(\mathbf{y}) \quad \hat{\boldsymbol{\eta}} = \arg \min_{\boldsymbol{\eta}} \operatorname{tr}(\Sigma_{\boldsymbol{\eta}} \Sigma_{\boldsymbol{\eta}}^{\top} H) - 2 \operatorname{tr}(\Sigma_{\boldsymbol{\eta}})$$

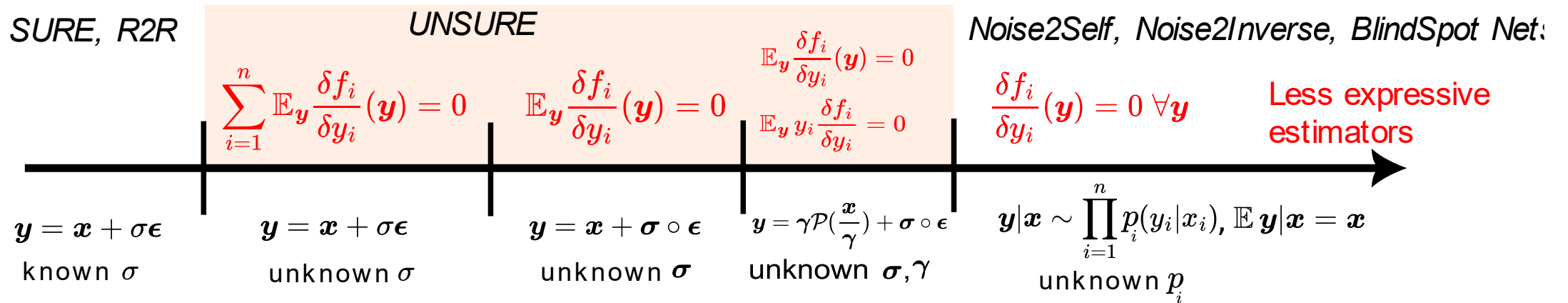
$$\text{with } H = \mathbb{E}_{\mathbf{y}} \nabla \log p_{\mathbf{y}}(\mathbf{y}) \nabla \log p_{\mathbf{y}}(\mathbf{y})^{\top}$$

Unknown Correlations

Practical examples:

- Unknown per-pixel noise level $[\Sigma_{\hat{\eta}}]_{i,i} = \frac{1}{\mathbb{E}_{\mathbf{y}} \frac{\delta \log p_{\mathbf{y}}(\mathbf{y})}{\delta y_i}}$
- Unknown spatial correlation $\Sigma_{\hat{\eta}} = \text{circ } \hat{\eta}$ with $\hat{\eta} = F^{-1}(\frac{1}{F\mathbf{h}})$ and \mathbf{h} the autocorrelation of the score up to $\pm r$ pixels

Summary



Implementation

We parameterize f as a deep neural network, small $\alpha > 0$, $\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}, I)$

UNSURE: Solve Lagrangian problem, approximating divergence as [Ramani et al., 2007]

$$\text{tr}(\Sigma \frac{\delta f}{\delta \mathbf{y}}) \approx \frac{(\Sigma \boldsymbol{\omega})^\top}{\alpha} (f(\mathbf{y}) - f(\mathbf{y} + \boldsymbol{\omega} \alpha))$$

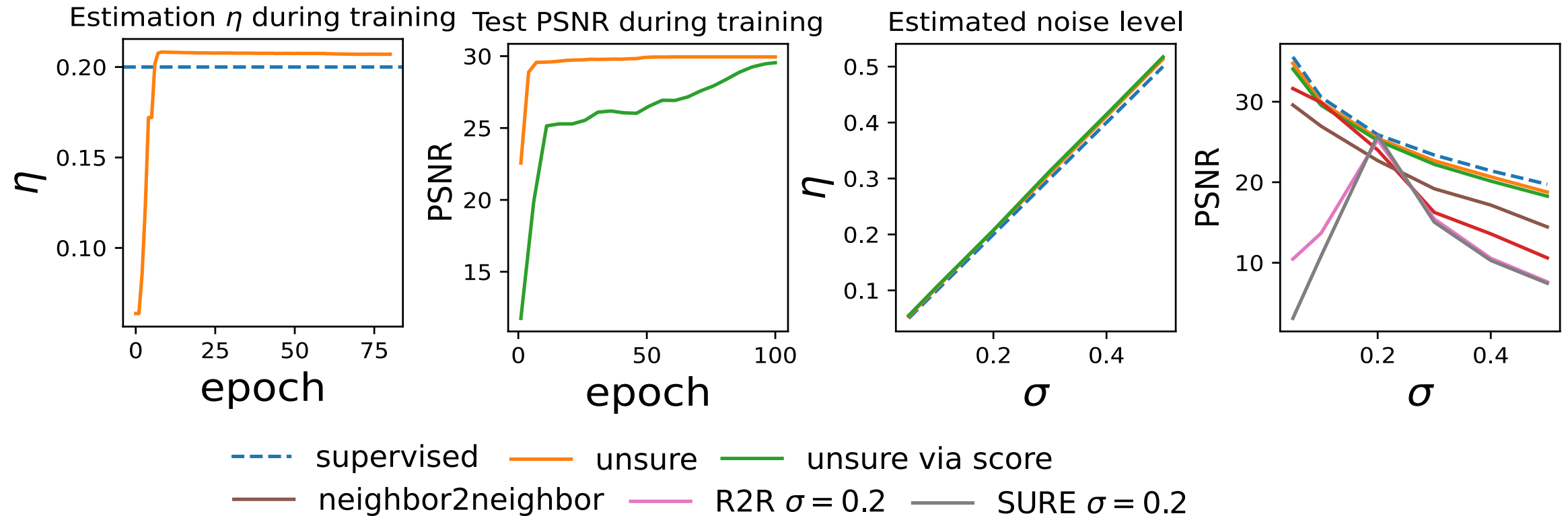
UNSURE via score: Learn score $s(\mathbf{y}) \approx \nabla \log p_{\mathbf{y}}(\mathbf{y})$ via

$$\arg \min_s \mathbb{E}_{\mathbf{y}, \boldsymbol{\omega}} ||\boldsymbol{\omega} - \alpha s(\mathbf{y} + \boldsymbol{\omega} \alpha)||^2$$

and use $f(\mathbf{y}) = \mathbf{y} + \Sigma_{\hat{\boldsymbol{\eta}}} s(\mathbf{y})$ to denoise at test time.

Experiments

- MNIST dataset
- Isotropic Gaussian noise



Experiments

- DIV2K dataset (320x320 RGB)
- Spatially correlated Gaussian noise

Method	Noise2Void	Neighbor2Neighbor	UNSURE (unknown Σ)	SURE (known Σ)	Supervised
PSNR [dB]	19.09 ± 1.79	23.61 ± 0.13	28.72 ± 1.03	29.77 ± 1.22	29.91 ± 1.26

Kernel size η	1×1	3×3	5×5
PSNR [dB]	23.62	28.72 ± 1.03	27.38 ± 0.88



Beyond Denoising

For $A \neq I$, most estimators can be adapted to approximate

$$\mathbb{E}_{\mathbf{x}, \mathbf{y}} ||A^\dagger A(\mathbf{x} - f(\mathbf{y}))||^2$$

where A^\dagger is the pseudoinverse of A .

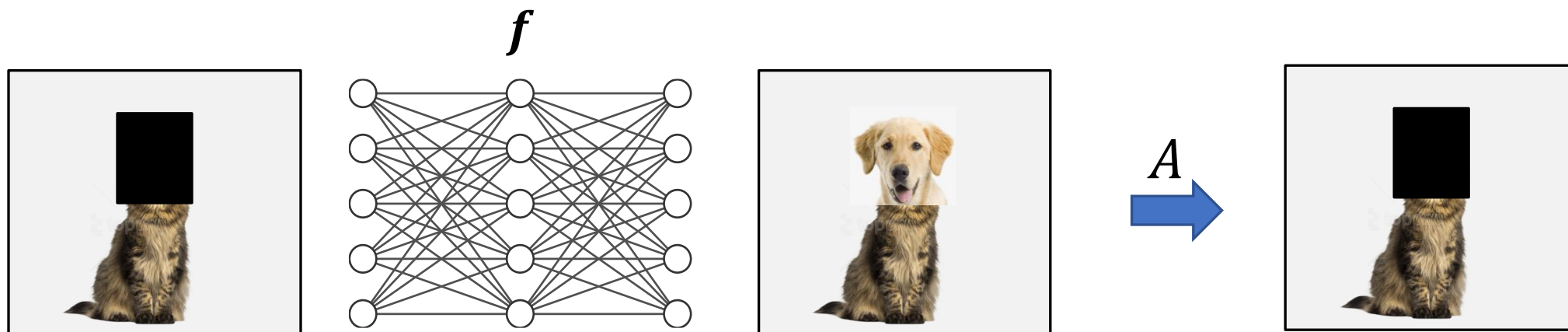
For example, **GSURE** [Eldar, 2008] writes for Gaussian noise

$$\mathcal{L}_{GSURE}(\mathbf{y}, f) = ||A^\dagger \mathbf{y} - A^\dagger A f(\mathbf{y})||^2 + 2\sigma^2 \sum_i \frac{\delta[A^\dagger A \cdot f]_i}{\delta y_i}(\mathbf{y})$$

Incomplete Measurements?

1. If A is invertible, we have $A^\dagger A = I$
2. If A is not invertible, $\mathbb{E}_{x,y} ||A^\dagger A(\mathbf{x} - f(\mathbf{y}))||^2 \neq \mathbb{E}_{x,y} ||\mathbf{x} - f(\mathbf{y})||^2$

In this case, the risk does not penalise $f(\mathbf{y})$ in the **nullspace** of A !



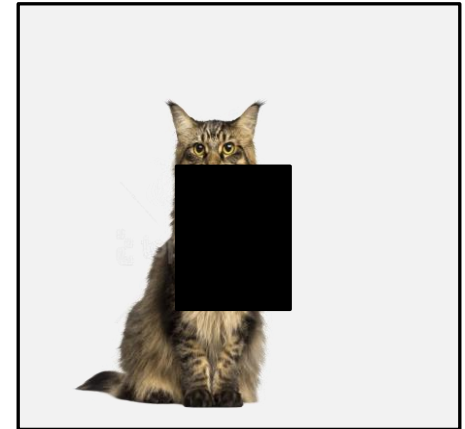
Symmetry Prior

Equivariant Imaging [Chen et al., 2021]

For all $g \in G$ we have

$$\mathbf{y} = A\mathbf{x} = \underbrace{AT_g}_{A_g} \overbrace{T_g^{-1}\mathbf{x}}^{\mathbf{x}'} = A_g\mathbf{x}'$$

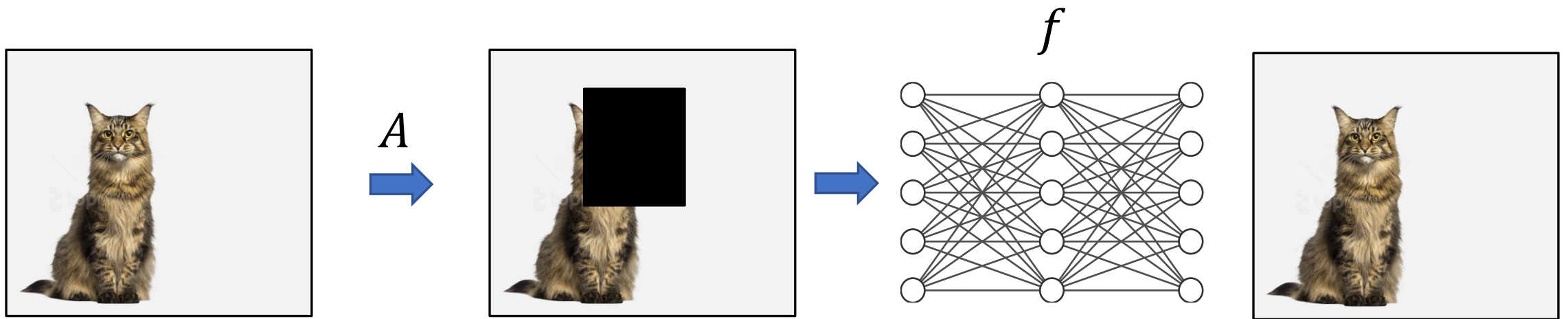
- We get multiple virtual operators $\{A_g\}_{g \in G}$ 'for free'!
- Each AT_g might have a different nullspace



Equivariant Imaging

How can we enforce equivariance in practice?

Idea: we should have $f(AT_g x) = T_g f(Ax)$, i.e. $f \circ A$ should be G -equivariant



Equivariant Imaging

We can leverage invariance of X to transformations T_g to learn in the nullspace [Chen, 2021]

$$\mathcal{L}_{EI}(\mathbf{y}, f) = \mathbb{E}_g || T_g \hat{\mathbf{x}} - f(AT_g \hat{\mathbf{x}}) ||^2$$

where $\hat{\mathbf{x}} = f(\mathbf{y})$ is used as reference

Robust Equivariant Imaging++

enforces equivariance of $f \circ A$

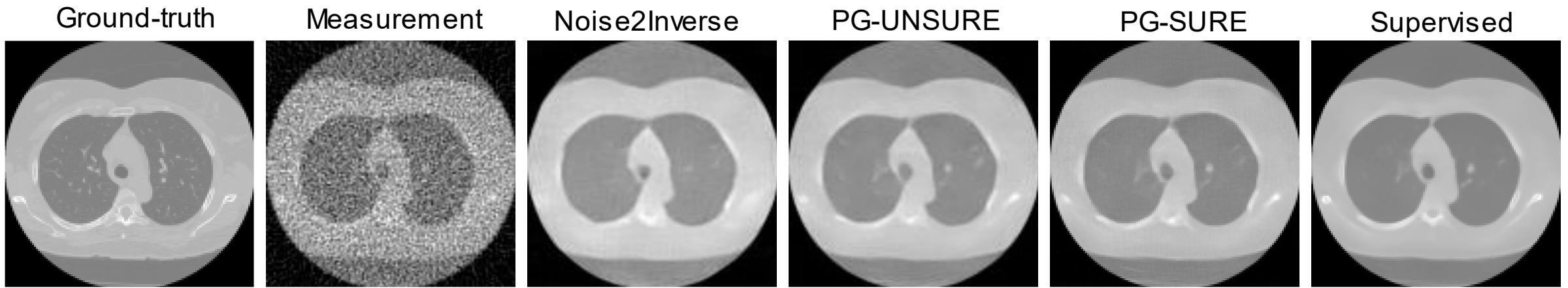
$$\mathcal{L}_{REI}(\mathbf{y}, f) = \mathcal{L}_{\text{UNSURE}}(\mathbf{y}, f) + \mathcal{L}_{EI}(\mathbf{y}, f)$$

Handles noisy measurements of unknown noise level

Experiments

- LIDC-IDRI dataset
- Poisson-Gaussian noise
- Uses equivariant imaging with rotations

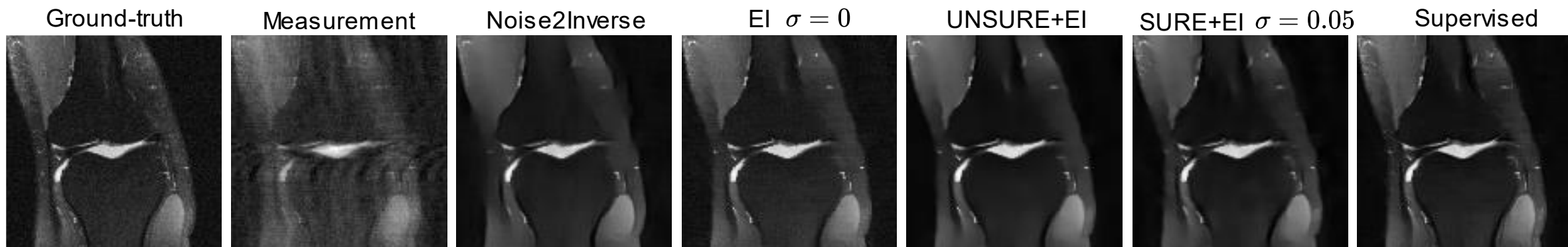
Method	Noise2Inverse	PG-UNSURE (unknown σ, γ)	PG-SURE (known σ, γ)	Supervised
PSNR [dB]	32.54 ± 0.71	33.31 ± 0.57	33.76 ± 0.61	34.67 ± 0.68



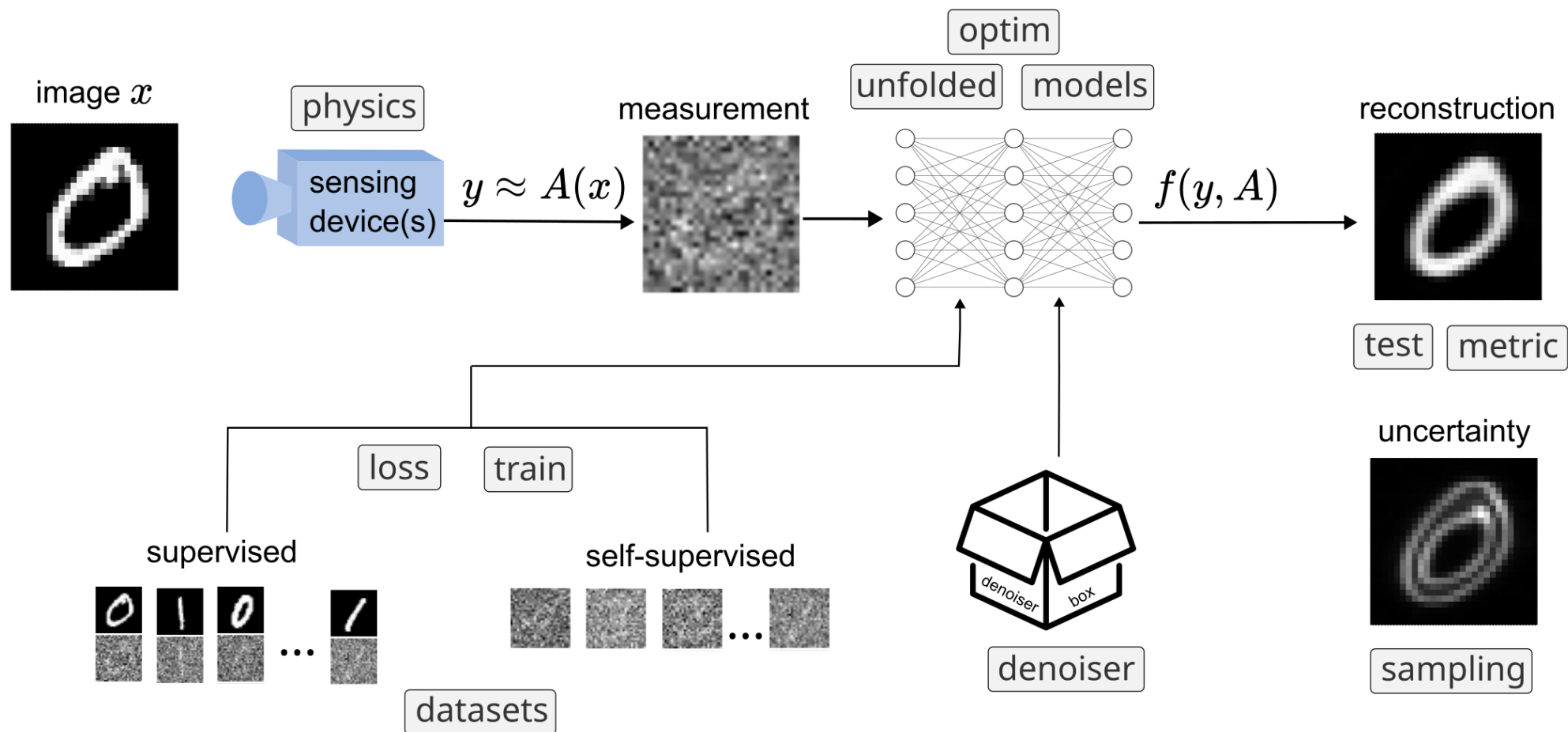
Experiments

- FastMRI dataset
- Gaussian noise

Method	CV + EI	EI (assumes $\sigma = 0$)	UNSURE + EI (unknown σ)	SURE + EI (assumes $\sigma = 0.05$)	Supervised
PSNR [dB]	33.25 ± 1.14	34.32 ± 0.91	35.73 ± 1.45	28.05 ± 4.73	36.63 ± 1.38



Deep Inverse



Uncertainty Quantification

Can we measure the uncertainty of the reconstructions?

Self-supervised losses can also be used for uncertainty quantification!

- SURE can be used to assess reconstruction error in denoising
- SURE4SURE [Bellec et al., 2021] gives error variance estimates.
- EI loss can be seen as a bootstrapping technique [T. & Pereyra, 2024] with well calibrated uncertainty estimates



References

Slides and codes of a recent 3-hour tutorial can be found here:

<https://tachella.github.io/projects/selfsuptutorial/>

