

# Self Supervised Learning Methods for Imaging

*Escuela de Primavera de Deep Learning, Buenos Aires, Argentina*

*Julián Tachella, CNRS, École Normale Supérieure de Lyon*

# The Inverse problem

**Goal:** estimate signal  $x$  from  $y$

$$y = A(x) + \epsilon$$

measurements  
 $\in \mathbb{R}^m$

↑  
Physics

signal  $\in \mathbb{R}^n$

noise/error

We will focus on linear problems where the forward operator  $A$  is a matrix

# Examples

	$x$	$A$	$y$	reconstruction	
<b>Magnetic Resonance Imaging (MRI)</b> <small>A: undersampled Fourier models</small>					 Source: Brian Hargreaves
<b>Black Hole Imaging</b> <small>A: spatial-frequency e.g. Event Horizon Telescope (EHT)</small>					 M87* April 11, 2017 The Astrophysical Journal Letters, vol. 875, no. L1, 2019.
<b>Cryogenic electron microscopy (Cryo-EM)</b> <small>A: 2D projections of protein particles</small>					 Covid-19 virus' structure D. Wrapp et al. <i>Science</i> , vol. 367, no. 6483, 2020.

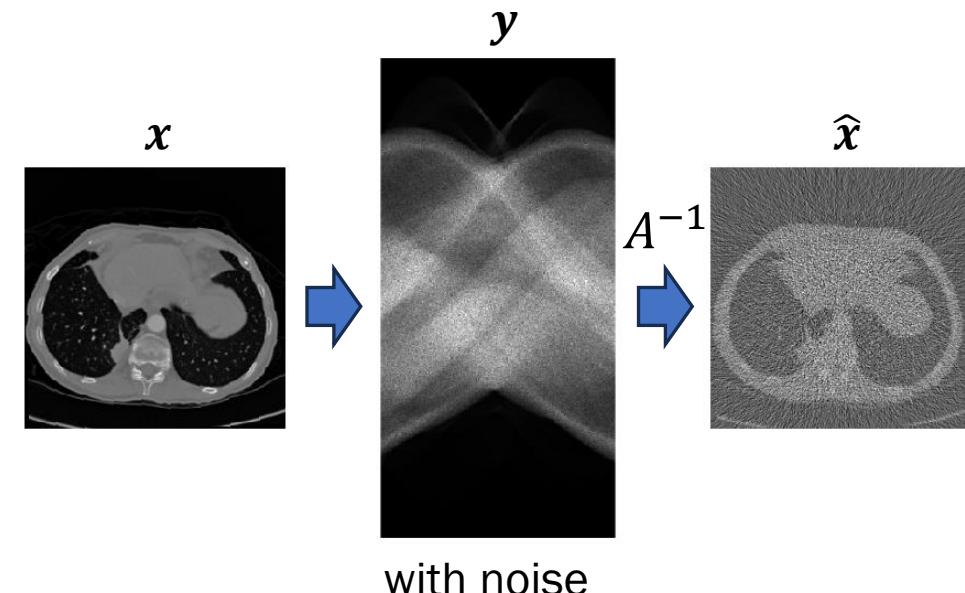
# Why it is hard to invert?

Measurements are usually corrupted by noise, e.g.

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

Can be additive, as above, or more complex, e.g. Poisson.

- Often, we do not know the exact noise distribution
- The forward operator may be poorly conditioned



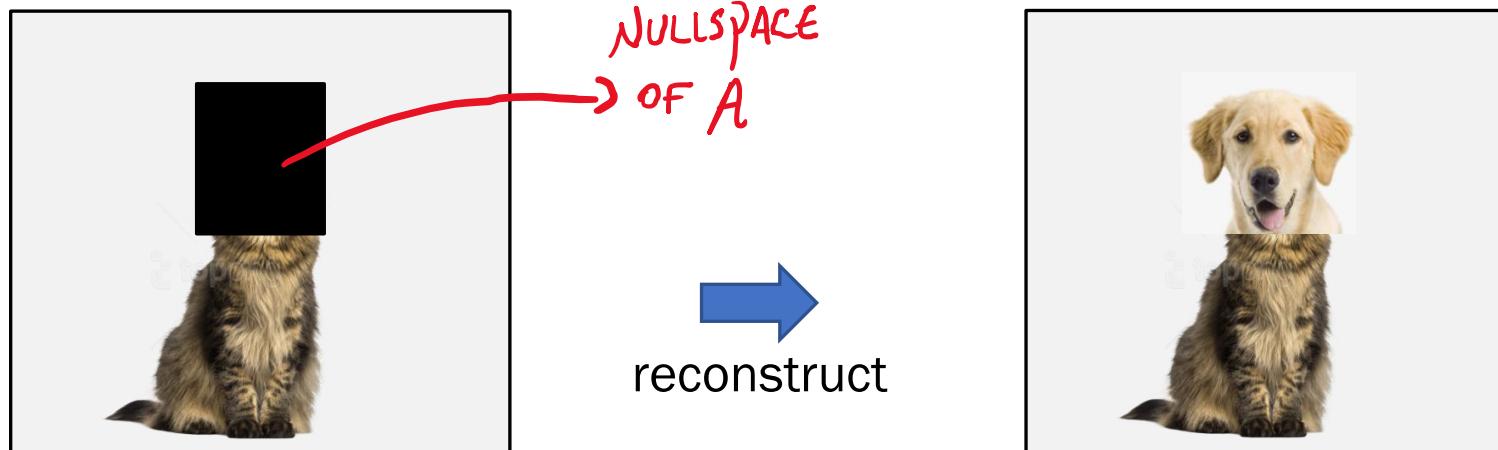
# Why it is hard to invert?

Even in the absence of noise,  $A$  may not be invertible, giving infinitely many  $\hat{x}$  consistent with  $y$ :

$$\hat{x} = A^\dagger y + v$$

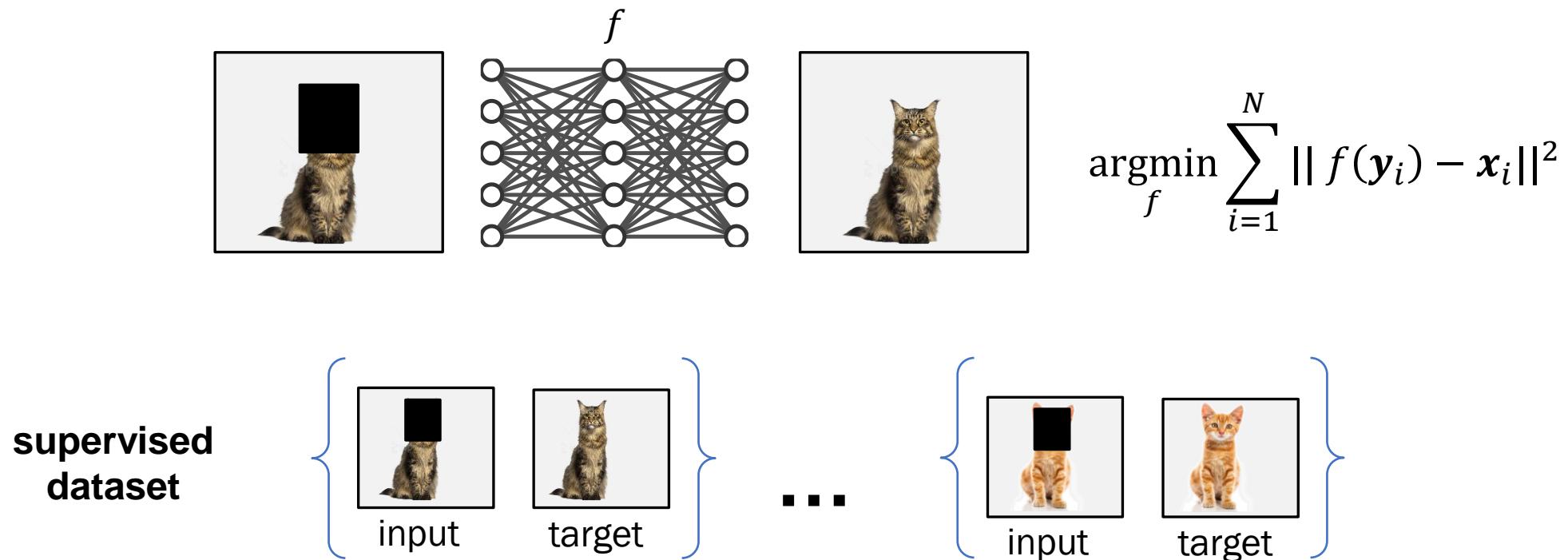
where  $A^\dagger$  is the pseudo-inverse of  $A$  and  $v$  is any vector in nullspace of  $A$

*Unique solution only possible if set of signals  $x$  is low-dimensional*



# 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,  $f_\theta$  is easy to evaluate

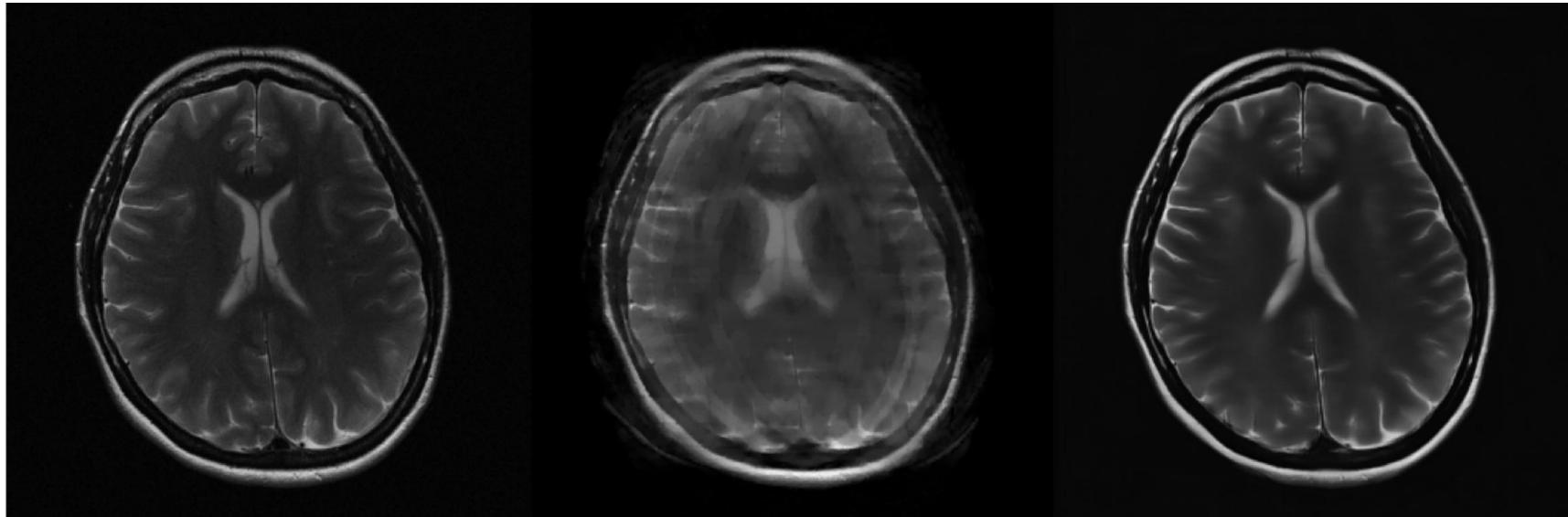
fastMRI

Accelerating MR Imaging with AI

Ground-truth

Total variation  
(28.2 dB)

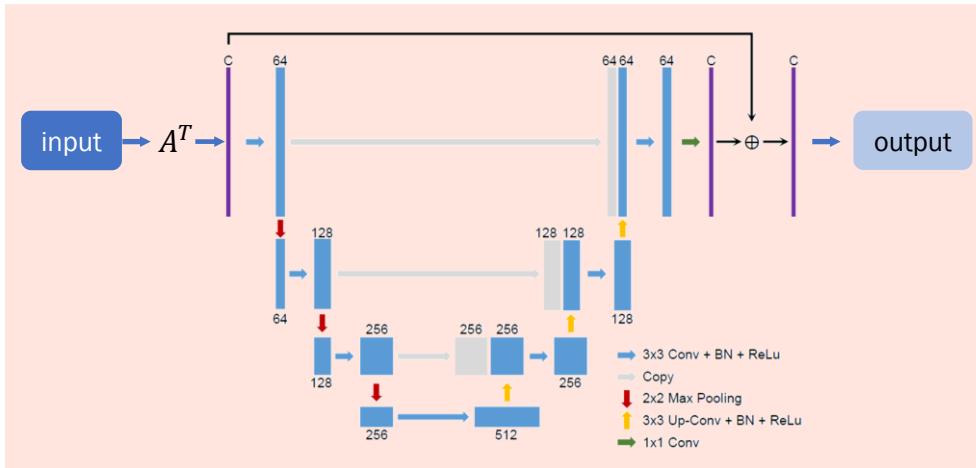
Deep network  
**(34.5 dB)**



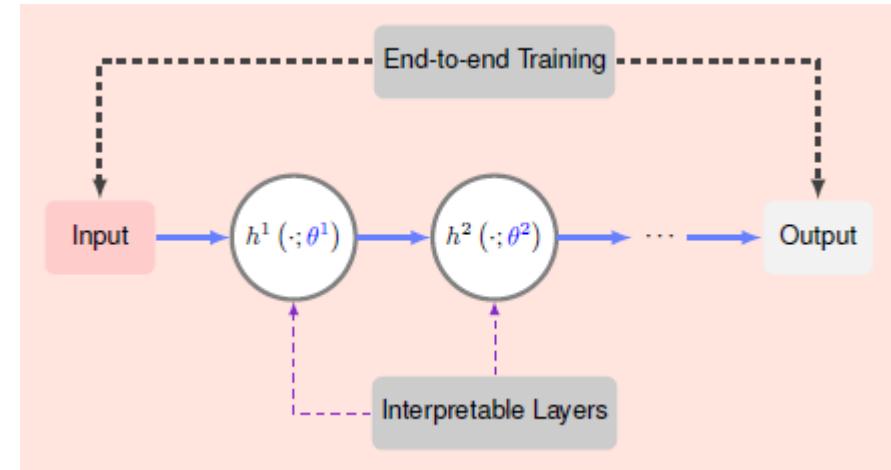
x8 accelerated MRI [Zbontar et al., 2019]

# Learning approach

Many architecture choices, e.g.



Back projected U-Net:  $\hat{x} = f(A^T y)$ , e.g. [Jin, 2017]



Unrolled networks:  $\hat{x} = f(y, A)$ , e.g. [Monga, 2020]

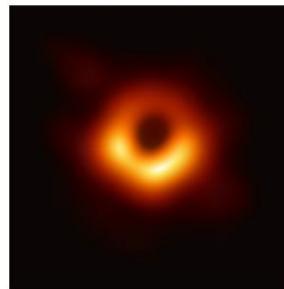
But also DnCNNs, DRUNet, SCUNet, DEQ, restormer, SwinIR, DiffPIR...

Here our focus will be on learning that is typically **architecture agnostic**

# Learning approach

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

- Medical and scientific imaging



- Distribution shift [Belthangady & Royer, 2019]

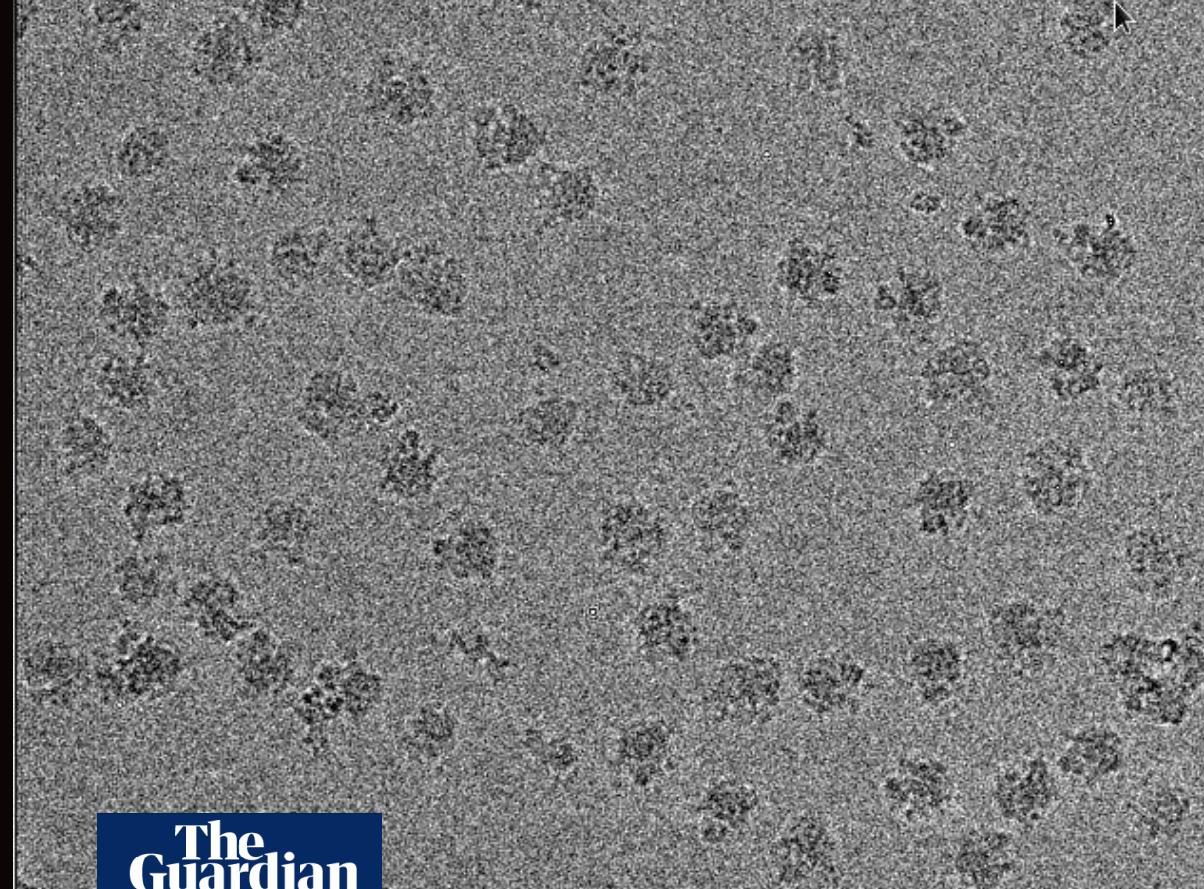
Training datasets	Witenagemot	Degraded image
abcdefg ...	witenagemot	
abc <del>defg</del> ...	Witenag <del>e</del> mot	
中文王国 ...	Witenag <del>e</del> mot	
a中b文c ...	Witenag <del>e</del> mot	
Old english word		Ground truth
		Witenagemot

# AI for Knowledge Discovery?

A black hole image showing a bright, circular, orange-yellow light source against a dark background.

The  
Guardian

Black hole picture captured for first time in space breakthrough

A grayscale image showing a complex, folded protein structure.

The  
Guardian

DeepMind uncovers structure of 200m proteins in scientific leap forward

# What this talk is **not** about

**Autoregressive models:** LLM pretraining on an autoregressive task.

**Self-supervised learning for feature learning:** Sim2Sim, masked autoencoders, DINOv2,3, etc. Focus on learning features for downstream tasks.

**Diffusion models and PnP:** require pre-trained denoiser/scores with ground-truth data

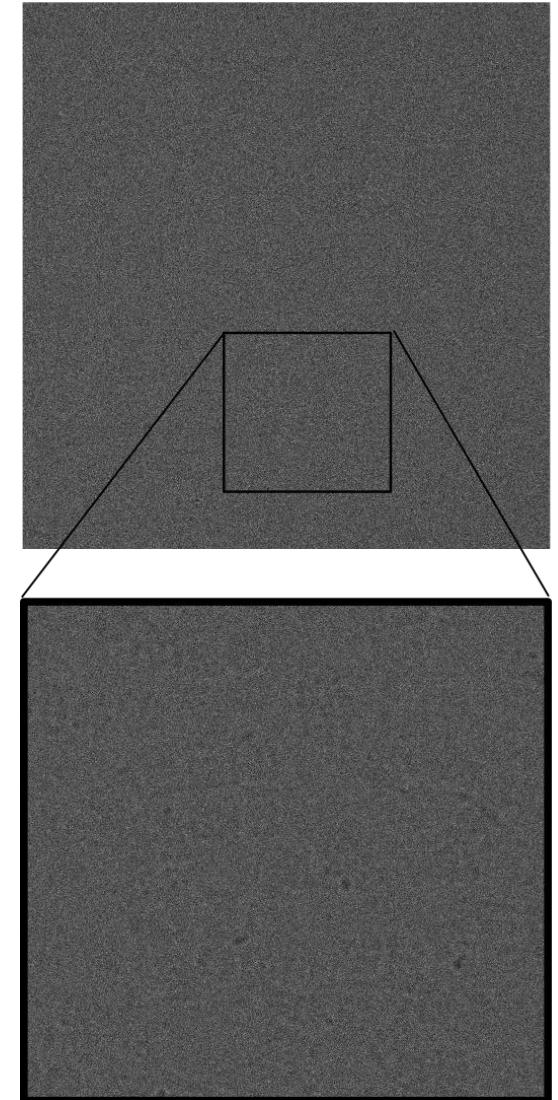
- The methods presented can be used to train denoisers without ground-truth!

# Purpose of this talk

How can we learn  $f$  from measurement  $\{\mathbf{y}_i\}_{i=1}^N$  data alone?

1. Noisy:  $\mathbf{y} = \mathbf{x} + \boldsymbol{\epsilon}$
2. Incomplete and noisy:  $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\epsilon}$

$$\operatorname{argmin}_f \sum_{i=1}^N \mathcal{L}(y_i, f)$$



# Best we can expect

We focus on  $\ell_2$  loss and minimum mean squared error estimators (MMSE)

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

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

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

# Self-supervised learning

**Approximating the supervised loss:**

1. Unbiased estimator

$$\mathbb{E}_y \mathcal{L}(y, f) = \mathbb{E}_{x,y} \|f(y) - x\|^2$$

2. Same minimizer

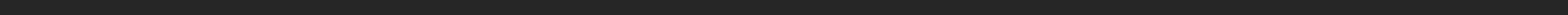
$$\operatorname{argmin}_f \mathbb{E}_y \mathcal{L}(y, f) = \operatorname{argmin}_f \mathbb{E}_{x,y} \|f(y) - x\|^2$$

3. Unbiased estimator under constraints

$$\mathbb{E}_y \mathcal{L}(y, f) = \mathbb{E}_{x,y} \|f(y) - x\|^2 \text{ for certain } f \neq \mathbb{E}\{x|y\}$$

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

# Part 2: Learning from noisy data



# Denoising problems

In this 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  $A = \mathbf{I}$  for simplicity.
- All methods in this part can be extended to any invertible  $A$ .

# Self-supervised risk estimators

## Supervised loss

$$\mathcal{L}_{\text{sup}}(x, y, f) = \|x - f(y)\|^2 = \|y - f(y)\|^2 + 2f(y)^T(y - x) + \text{const.}$$

Measurement consistency      key term to approximate!  
 $= f(y)^T \epsilon$

Naïve loss doesn't work!

$$\mathcal{L}_{\text{MC}}(y, f) = \|y - f(y)\|^2$$

$$\rightarrow f^*(y) = y$$

# Noise2Noise

**Mallows  $C_p$**  [Mallows, 1973], **Noise2Noise** [Lehtinen, 2018]

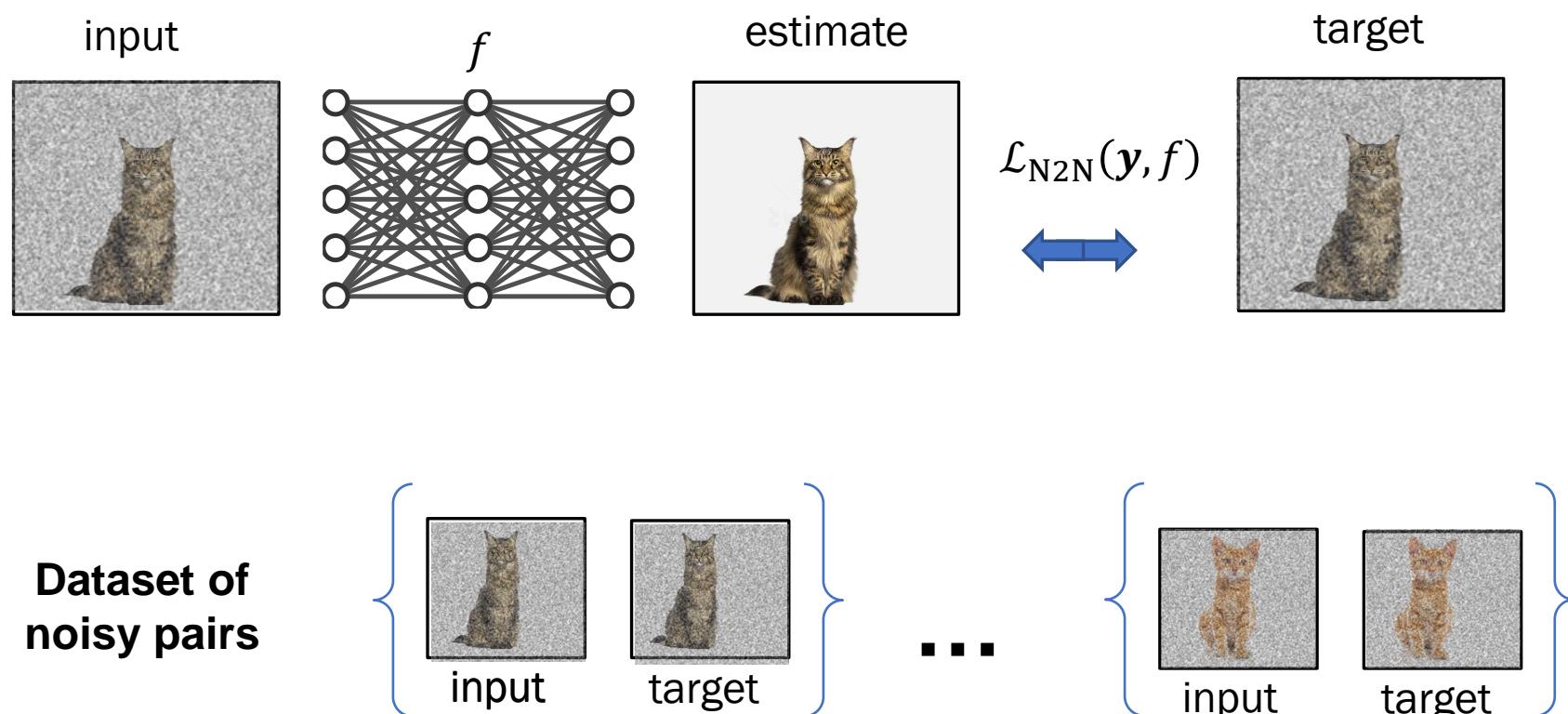
- **Independent** pairs  $\mathbf{y}_a = \mathbf{x} + \boldsymbol{\epsilon}_a$  and  $\mathbf{y}_b = \mathbf{x} + \boldsymbol{\epsilon}_b$  with  $\boldsymbol{\epsilon}_a, \boldsymbol{\epsilon}_b$  **independent**
- $\mathbb{E}_{\boldsymbol{\epsilon}_b|\mathbf{x}} \boldsymbol{\epsilon}_b = \mathbf{0}$

$$\mathcal{L}_{\text{Noise2Noise}}(\mathbf{y}_a, \mathbf{y}_b, f) = \|\mathbf{y}_b - f(\mathbf{y}_a)\|^2$$

$$\mathbb{E}_{\mathbf{y}_b|\mathbf{x}} f(\mathbf{y}_a)^\top (\mathbf{y}_b - \mathbf{x}) = f(\mathbf{x} + \boldsymbol{\epsilon}_a) \underbrace{\mathbb{E} \boldsymbol{\epsilon}_b}_{=0} = 0$$

- Also works for any noise distribution with  $\mathbb{E}_{\mathbf{y}_b|\mathbf{x}} \mathbf{y}_b = \mathbf{x}$

# Noise2Noise



**Not useful for the microscopy example!**

# Recorrupted2Recorrupted

**Recorrupted2Recorrupted** [Pang et al., 2021], **Coupled Bootstrap** [Oliveira et al., 2022], **Noisier2Noise** [Moran et al., 2020].

**Proposition:** Let  $\mathbf{y} \sim N(\mathbf{x}, I\sigma^2)$  and define

$$\begin{aligned}\mathbf{y}_a &= \mathbf{y} + \sqrt{\frac{1-\alpha}{\alpha}} \boldsymbol{\omega} \\ \mathbf{y}_b &= \frac{\mathbf{y}}{\alpha} - \frac{\mathbf{y}_a(1-\alpha)}{\alpha}\end{aligned}$$

where  $\boldsymbol{\omega} \sim N(\mathbf{0}, I\sigma^2)$  and  $\alpha \in \mathbb{R}$ , then  $\mathbf{y}_a$  and  $\mathbf{y}_b$  are **independent** random variables (fixed  $\mathbf{x}$ ).

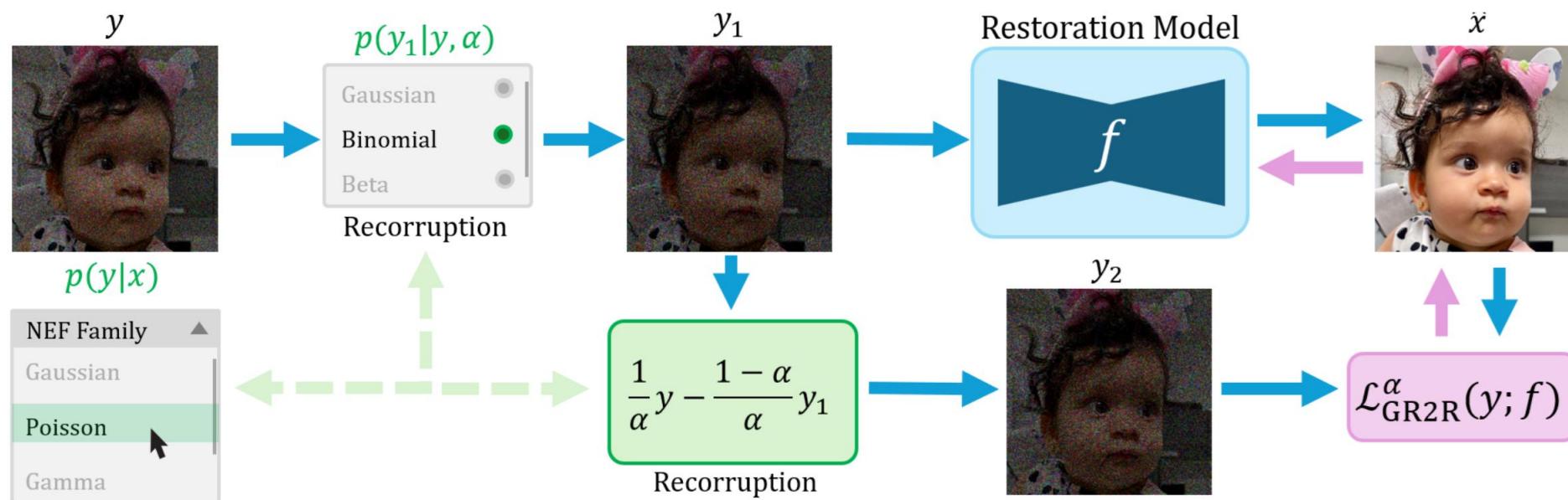
$$\mathcal{L}_{R2R}(\mathbf{y}, f) = \mathbb{E}_{\boldsymbol{\omega}} \| \mathbf{y}_b - f(\mathbf{y}_a) \|^2$$

- Price to pay:  $\text{SNR}(\mathbf{y}_a) < \text{SNR}(\mathbf{y})$
- At **test time**,  $f^{\text{test}}(\mathbf{y}) = \frac{1}{N} \sum_i f\left(\mathbf{y} + \sqrt{\frac{1-\alpha}{\alpha}} \boldsymbol{\omega}_i\right)$  with  $\boldsymbol{\omega}_i \sim \mathcal{N}(\mathbf{0}, I\sigma^2)$

# Recorrupted2Recorrupted

- Can be extended to other noise distributions [Monroy, Bacca and Tachella, CVPR 2025]

Model	Gaussian $y \sim \mathcal{N}(x, \Sigma)$	Poisson $z \sim \mathcal{P}(x/\gamma), y = \gamma z$	Gamma $y \sim \mathcal{G}(\ell, \ell/x)$
$y_1$	$y_1 = y + \sqrt{\frac{\alpha}{1-\alpha}}\omega, \quad \omega \sim \mathcal{N}(0, \Sigma)$	$y_1 = \frac{y - \gamma\omega}{1-\alpha}, \quad \omega \sim \text{Bin}(z, \alpha)$	$y_1 = y \circ (1 - \omega)/(1 - \alpha) \quad \omega \sim \text{Beta}(\ell\alpha, \ell(1 - \alpha))$



# 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

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

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

Measurement consistency      Degrees of freedom [Efron, 2004]

- **Hudson's lemma** [Hudson 1978] extends this result for the exponential family (eg. **Poisson Noise**)
- Beyond exponential family: **Poisson-Gaussian noise** [Le Montagner et al., 2014]  
[Raphan and Simoncelli, 2011]

# Stein's Unbiased Risk Estimator

**Monte Carlo SURE** [Efron 1975, Breiman 1992, Ramani et al., 2007]

SURE's divergence is generally approximated as

$$\sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y}) \approx \frac{\boldsymbol{\omega}^\top}{\alpha} (f(\mathbf{y}) - f(\mathbf{y} + \boldsymbol{\omega}\alpha))$$

where  $\alpha > 0$  small,  $\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}, I)$

- Noisier2Noise is equivalent to SURE when  $\alpha \rightarrow 0$  [Monroy Bacca and Tachella, CVPR 2025].

# Stein's Unbiased Risk Estimator

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

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

## Integration by parts

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

## Complete squares

$$\arg \min_f \mathbb{E}_y || f(y) - y - \sigma^2 \nabla \log p_y(y) ||^2$$

$$\rightarrow f(y) = y + \sigma^2 \nabla \log p_y(y)$$

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

# Summary So Far

	Train Eval	Test Eval	Single $y$	MMSE optimal	Unknown noise
Noise2Noise	1	1		✓	✓
R2R	1	>1	✓	✓	
SURE	2	1	✓	✓	

If we have a single  $y$  and don't know the noise distribution?

# UNSURE

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

**UNSURE** [Tachella et al., ICLR 2025]

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

- SURE's perspective:  $\mathcal{L}_{\text{SURE}}(\mathbf{y}, f) = \|\mathbf{y} - f(\mathbf{y})\|^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y})$
- Not MMSE optimal (but almost)

# UNSURE

- In practice, we use Lagrange multipliers

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

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

- Expected error

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

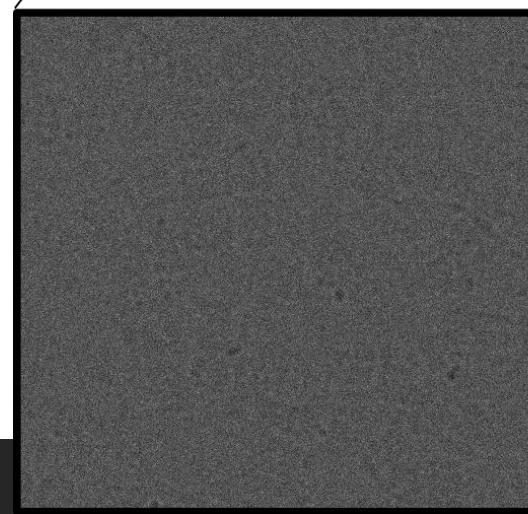
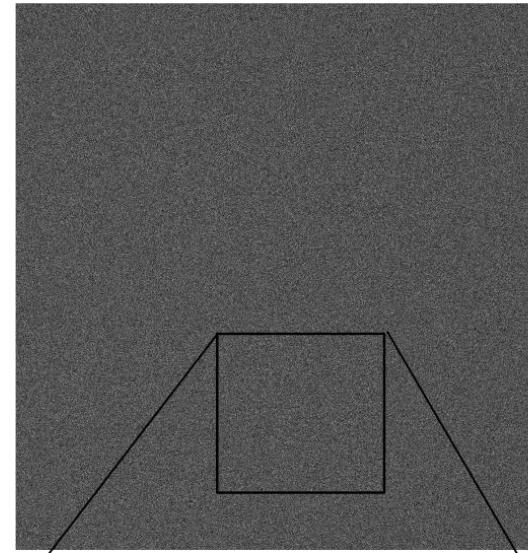
- UNSURE can be extended to **unknown noise covariance** and **Poisson Gaussian noise**

# Experiments

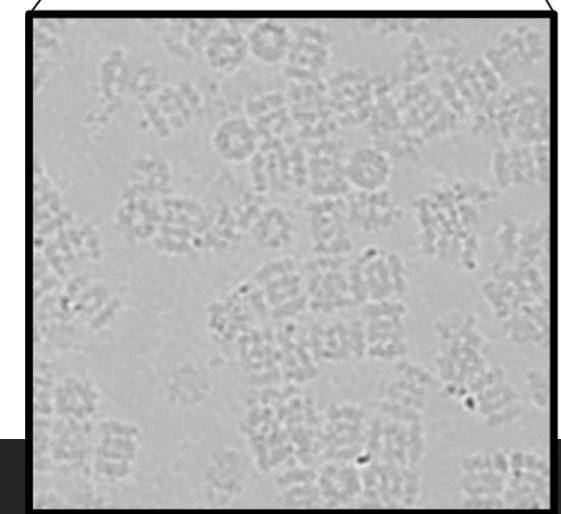
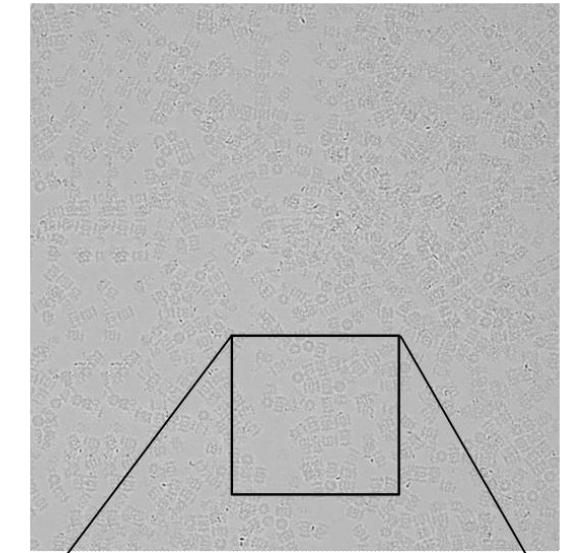
## Real data experiments

- Cryo electron microscopy images
- Extremely low SNR
- Approx. Poisson-Gaussian noise

Measurement



PG-UNSURE



# Cross-Validation Methods

**Assumption:**  $f_i$  does not depend on  $y_i$ , that is  $\frac{\delta f_i}{\delta y_i} = 0$ . Decomposable noise  $p(\mathbf{y}|\mathbf{x}) = \prod p(y_i|x_i)$

$$\mathbb{E}_{\mathbf{y}|\mathbf{x}} \sum_{i=1}^n f_i(\mathbf{y})(y_i - x_i) = \sum_{i=1}^n \mathbb{E}_{\mathbf{y}_{-i}|\mathbf{x}} f_i(\mathbf{y}_{-i}) \mathbb{E}_{y_i|x_i} (y_i - x_i) = 0$$

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

- SURE's perspective:  $\mathcal{L}_{SURE}(\mathbf{y}, f) = \|\mathbf{y} - f(\mathbf{y})\|^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\mathbf{y})$
- These methods are not MMSE optimal
- How to remove dependence on  $y_i$ : training or architecture

# Measurement Splitting

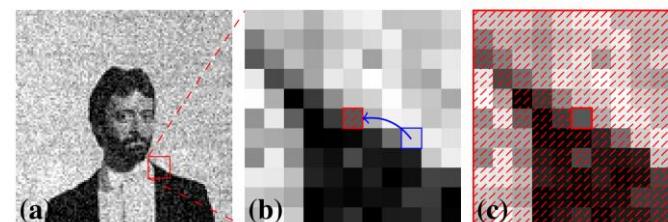
**Cross-validation** [Efron, 2004]: random split  $\mathbf{y} = \begin{bmatrix} \mathbf{y}_a \\ \mathbf{y}_b \end{bmatrix}$  at each iteration

$$\mathcal{L}_{N2V}(\mathbf{y}, f) = \mathbb{E}_{a,b} \|\mathbf{y}_b - \text{diag } \mathbf{m}_b f(\mathbf{y}_a)\|^2$$

where  $\mathbf{m}_b \in \{0,1\}^n$  masks out the pixels in  $\mathbf{y}_a$ .

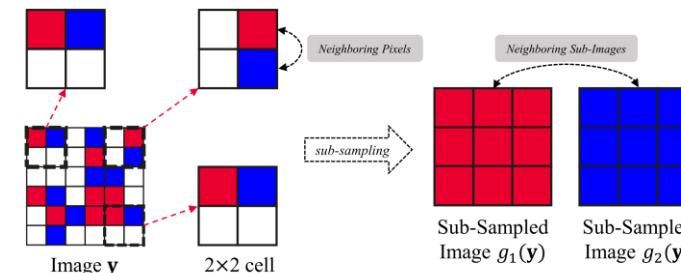
**Noise2Void** [Krull et al., 2019], **Noise2Self** [Batson, 2019]

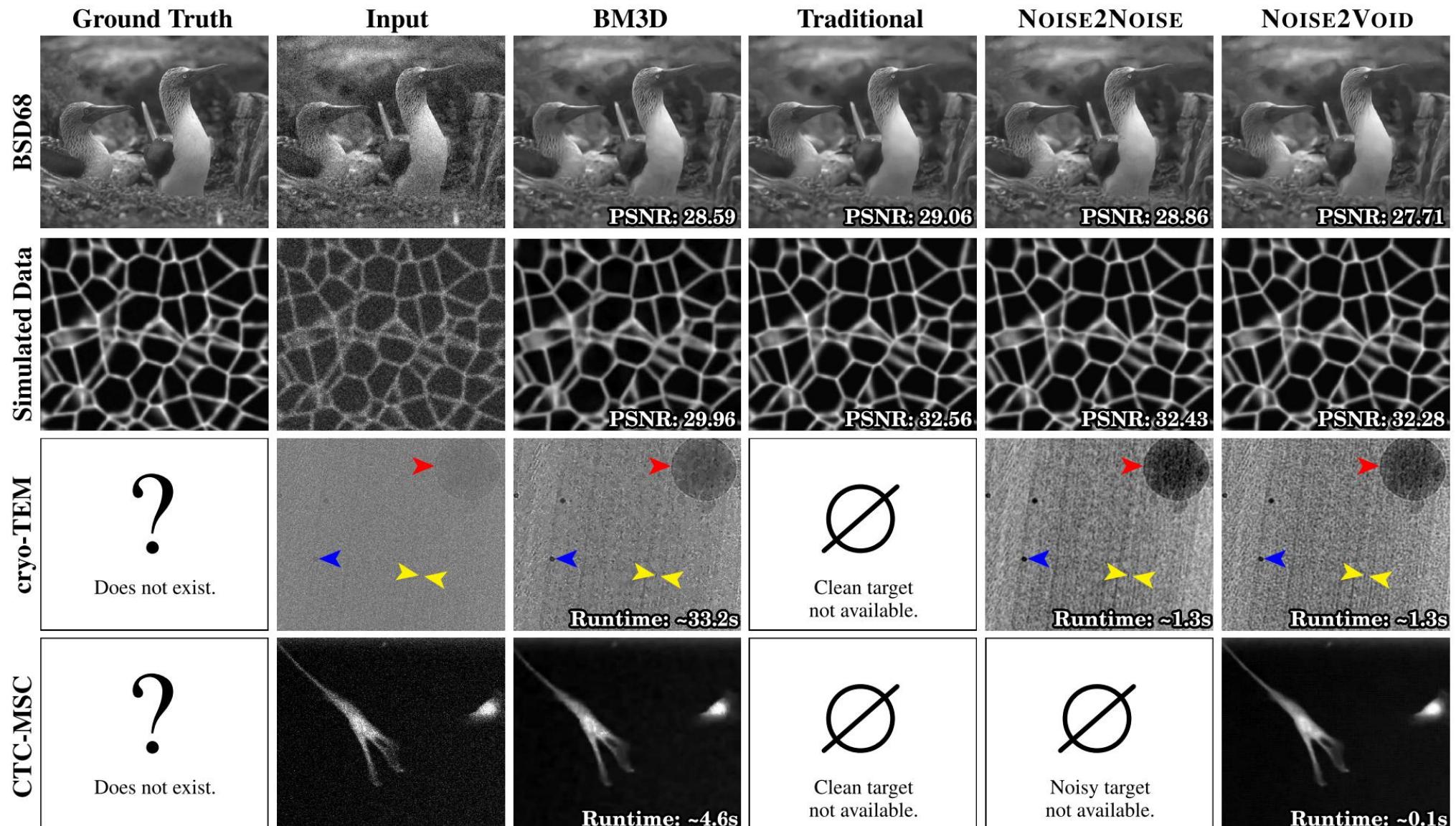
- During training flip centre pixel
- Computes loss only on flipped pixels



**Neighbor2Neighbor** [Huang, 2023]

- Use different subsampling as input and target
- Assumes scale invariance





# Measurement Splitting

At **test time**,  $f(\mathbf{y})$  is evaluated as

1. Test  $f$  as trained (expensive)

$$f^{\text{test}}(\mathbf{y}) = \frac{1}{N} \sum_i M f(\mathbf{y}_{\mathbf{a}_i}) \text{ with } \mathbf{y}_{\mathbf{a}_i} \sim p(\mathbf{y}_{\mathbf{a}}|\mathbf{y}) \text{ and } M = \left( \sum_i^N \text{diag}(\mathbf{m}_{b,i}) \right)^{-1}$$

2. Assume good generalization of  $f$  (cheap)

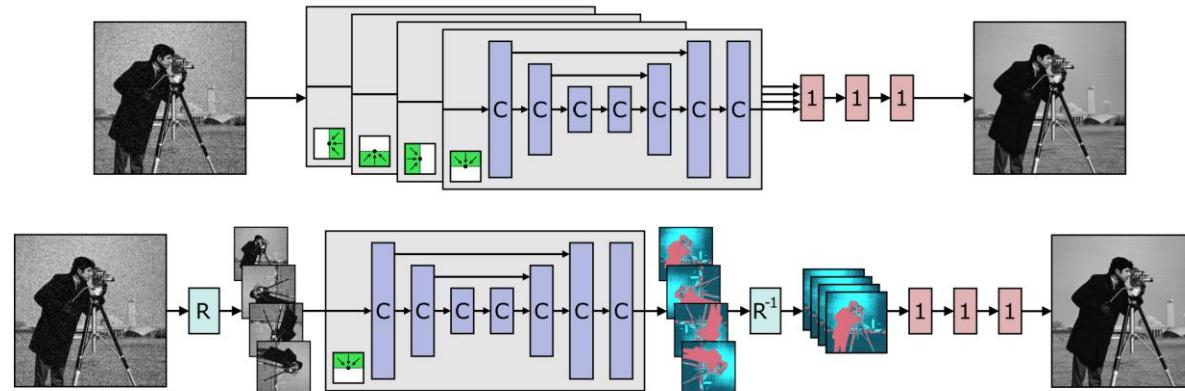
- $f^{\text{test}}(\mathbf{y}) = f(\mathbf{y}_{\mathbf{a}})$  with  $\mathbf{y}_{\mathbf{a}} \sim p(\mathbf{y}_{\mathbf{a}}|\mathbf{y})$
- $f^{\text{test}}(\mathbf{y}) = f(\mathbf{y})$

# Blind Spot Networks

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

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

$$\mathcal{L}_{\text{BS}}(\mathbf{y}, f_{\text{BS}}) = \|\mathbf{y} - f_{\text{BS}}(\mathbf{y})\|^2$$

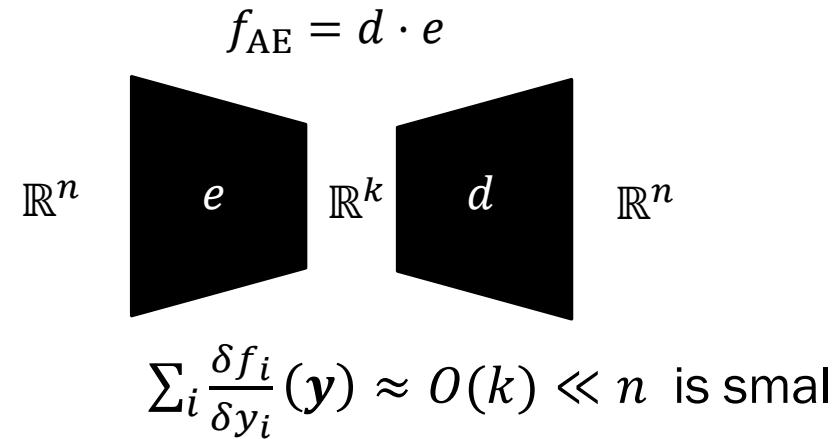


# Autoencoders

## Autoencoders

**Assume**

- $f$  has a strong bottleneck



$$\mathcal{L}_{\text{AE}}(\mathbf{y}, f) = \|\mathbf{y} - f_{\text{AE}}(\mathbf{y})\|^2$$

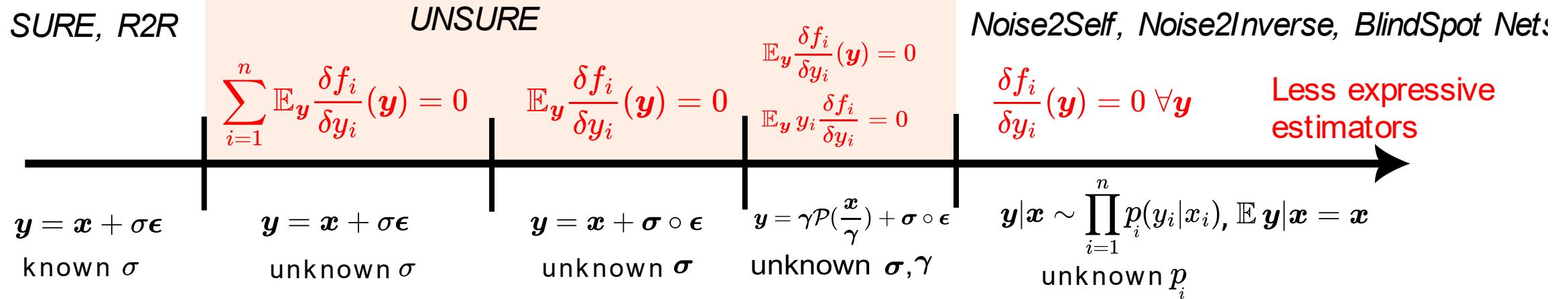
- Noise distribution is ‘high-dimensional’ whereas signal distribution is ‘low-dimensional’

# Summary

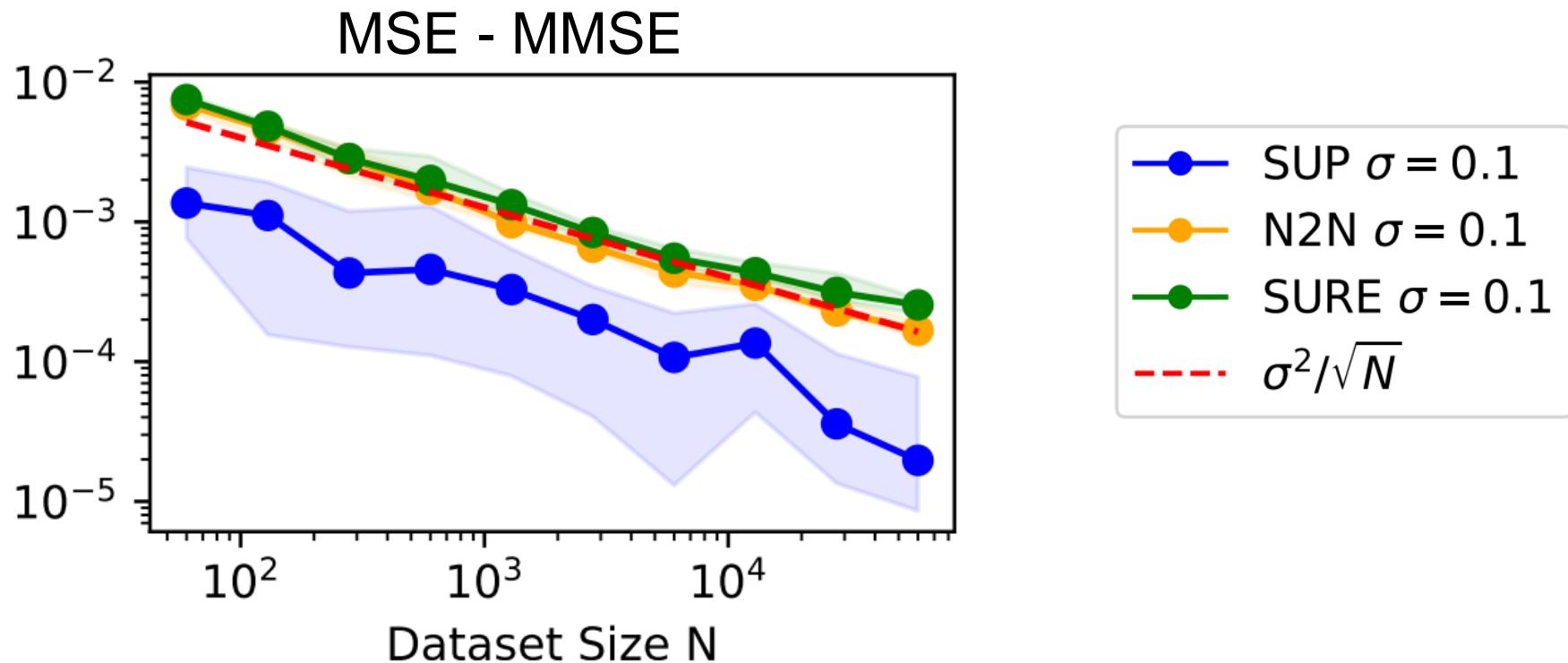
	Train Eval	Test Eval	Single $y$	MMSE optimal	Unknown separable noise	Unknown coloured noise
Noise2Noise	1	1		✓	✓	✓
R2R	1	>1	✓	✓		
SURE	2	1	✓	✓		
UNSURE	2	1	✓	✓	✓	✓
Noise2Void	1	1	✓		✓	
Blind Spot	1	>1	✓		✓	
Autoencoders	1	1			✓	✓

**No free lunch:** less assumptions about noise = less optimal estimator

# Summary

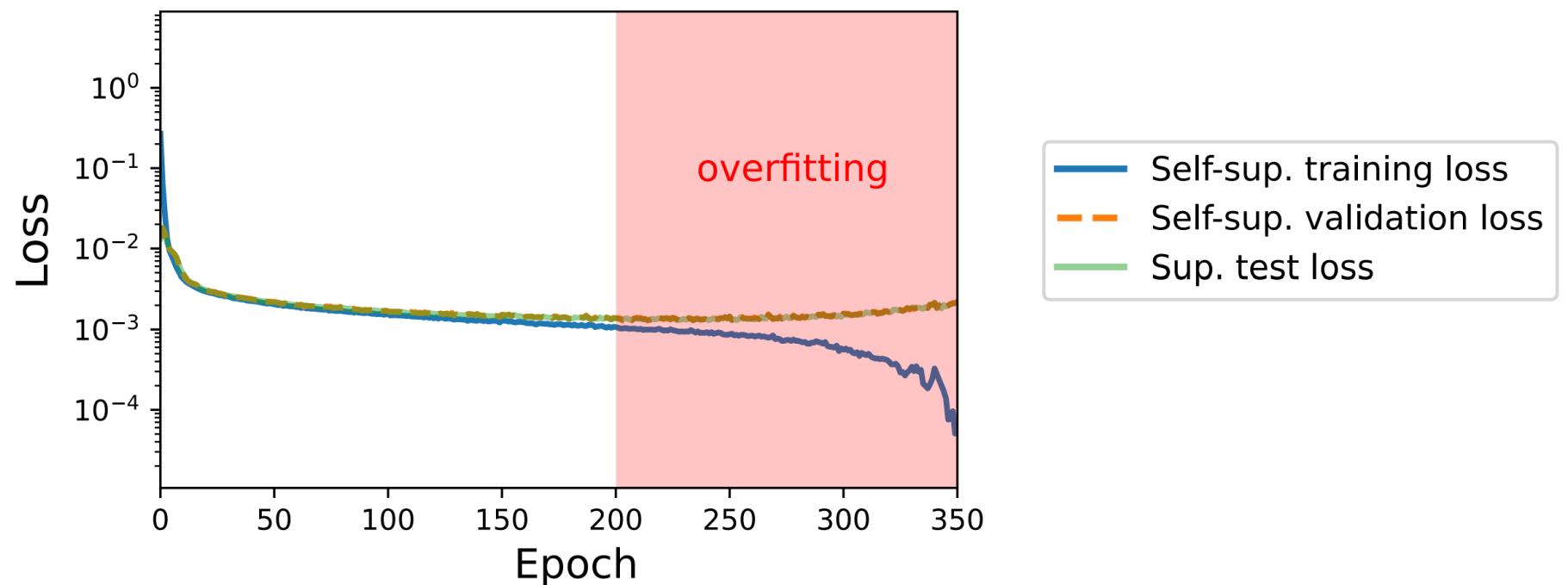


# Sample Complexity



# Self-supervised validation

- Follow the same validation practices in supervised learning

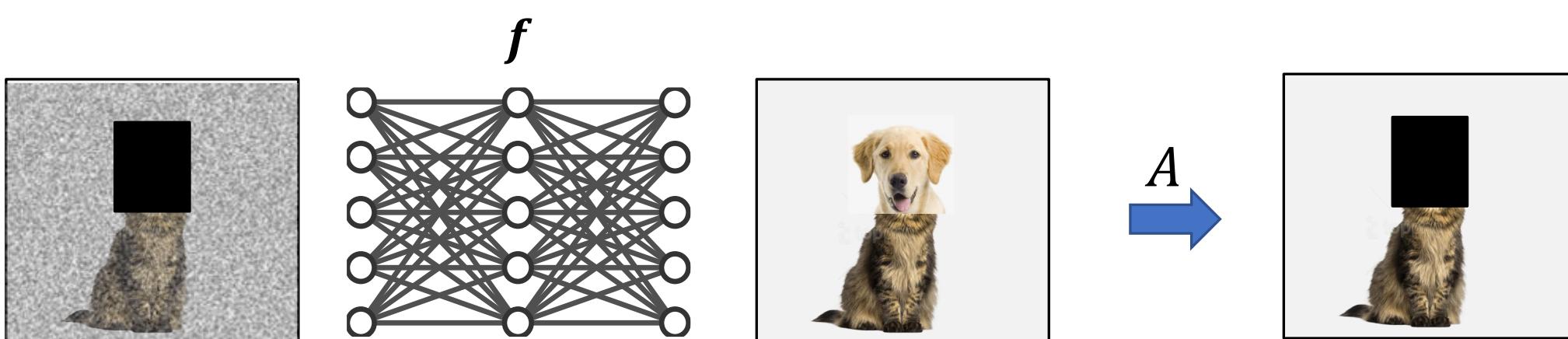


# Incomplete Measurements?

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

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

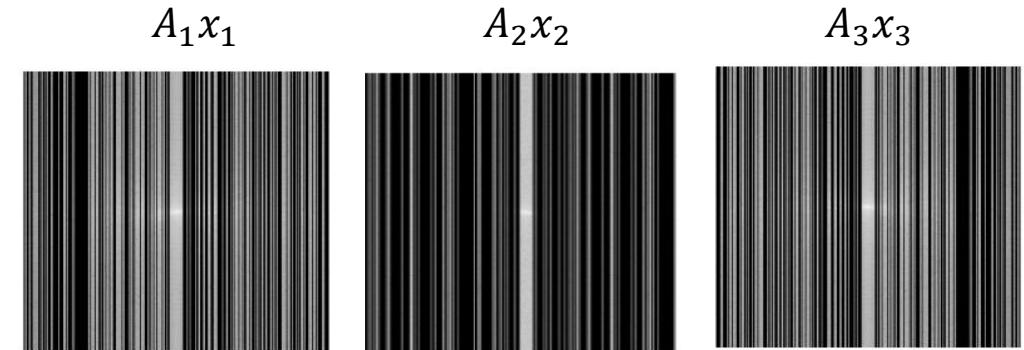
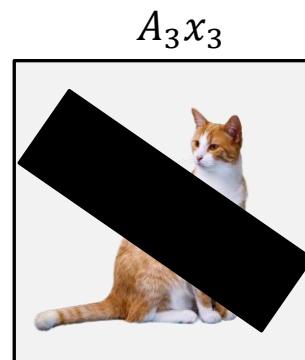
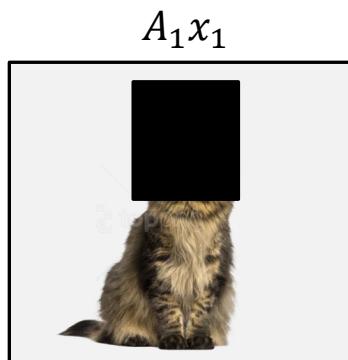
In this case, the risk does not penalise  $f(\mathbf{y})$  in the **nullspace** of  $A$ !



# Learning from Measurements

**How to learn from only  $y$ ?**

- Access multiple operators  $y_i = A_{g_i}x_i$  with  $g \in \{1, \dots, G\}$
- Each  $A_g$  with different nullspace
- Offers the possibility for learning using multiple measurement operators



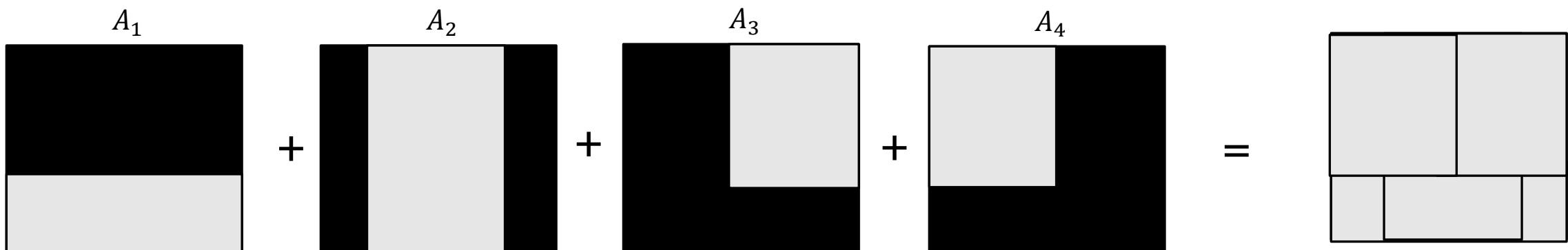
# Necessary Condition

**Intuition:** we need that the operators  $A_1, A_2, \dots, A_G$  cover the whole ambient space [T., 2022].

**Proposition:** Learning reconstruction mapping  $f$  from observed measurements possible only if

$$\text{rank}(\mathbb{E}_g A_g^\top A_g) = n$$

and thus, if  $m \geq n/G$ .



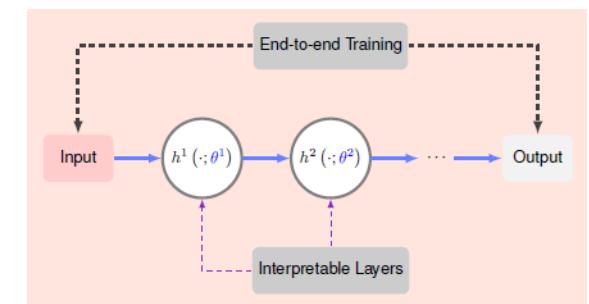
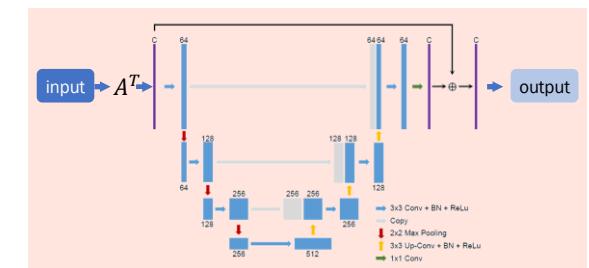
# Learning Approach

We will consider networks  $\hat{x} = f(\mathbf{y}, A)$ , where  $f$  is also a function of measurement operator e.g.,

- Filtered back projection  $f(\mathbf{y}, A) = f(A^\dagger \mathbf{y})$
- Unrolled networks...
- Naïve **measurement consistency** loss:

$$\mathcal{L}_{MC}(f) = \mathbb{E}_{\mathbf{y}, g} \|\mathbf{y} - A_g f(\mathbf{y}, A_g)\|^2$$

Without noise, a minimizer is the trivial solution  $f(\mathbf{y}, A_g) = A_g^\dagger \mathbf{y}, \forall g$

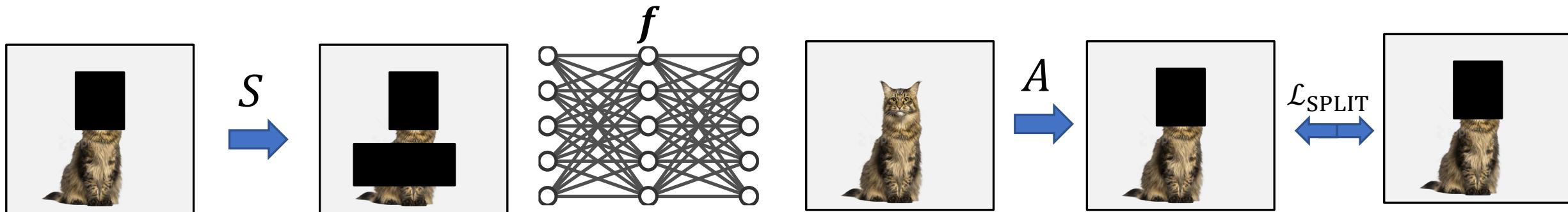


# Measurement Splitting Revisited

**Self-supervised learning via data undersampling (SSDU) [Yaman et al., 2019]**

Assume clean measurements, sample additional mask  $S \sim p(S | A_g)$  and mask inputs

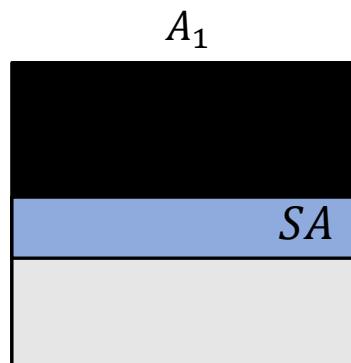
$$\mathcal{L}_{\text{SPLIT}}(\mathbf{y}, f) = \mathbb{E}_S \| \mathbf{y} - A f(S\mathbf{y}, SA) \|^2$$



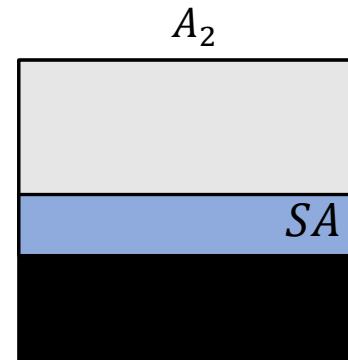
# Measurement Splitting Revisited

**Theorem [Daras et al. 2024].** If  $\mathbb{E}\{A^\top A | SA\}$  has full rank, it has same minimizer as supervised loss,  
 $f^*(S\mathbf{y}, SA) = \mathbb{E}\{\mathbf{x} | S\mathbf{y}, SA\}$

Inpainting example with  $G = 2$  operators



$$f_1^*(\mathbf{y}) = \mathbb{E}\{A_1 \mathbf{x} | S\mathbf{y}, SA\}$$



$$f_2^*(\mathbf{y}) = \mathbb{E}\{A_2 \mathbf{x} | S\mathbf{y}, SA\}$$

$$f^*(\mathbf{y}) = f_1^*(\mathbf{y}) + f_2^*(\mathbf{y}) = \mathbb{E}\{\mathbf{x} | S\mathbf{y}, SA\}$$

# Measurement Splitting Revisited

What happens if we have **noisy data**?

$$\begin{aligned}\mathcal{L}_{\text{SPLIT}}(\mathbf{y}, f) &= \mathbb{E}_S \|\mathbf{y} - A f(S\mathbf{y}, SA)\|^2 \\ &= \mathbb{E}_S \|\underbrace{S\mathbf{y} - SA f(S\mathbf{y}, SA)}_{} \|^2 + \|\underbrace{(I - S)\mathbf{y} - (I - S)A f(S\mathbf{y}, SA)}_{} \|^2\end{aligned}$$

Can be replaced by SURE,  
R2R, etc.

Not affected by separable noise

# Using All Measurements

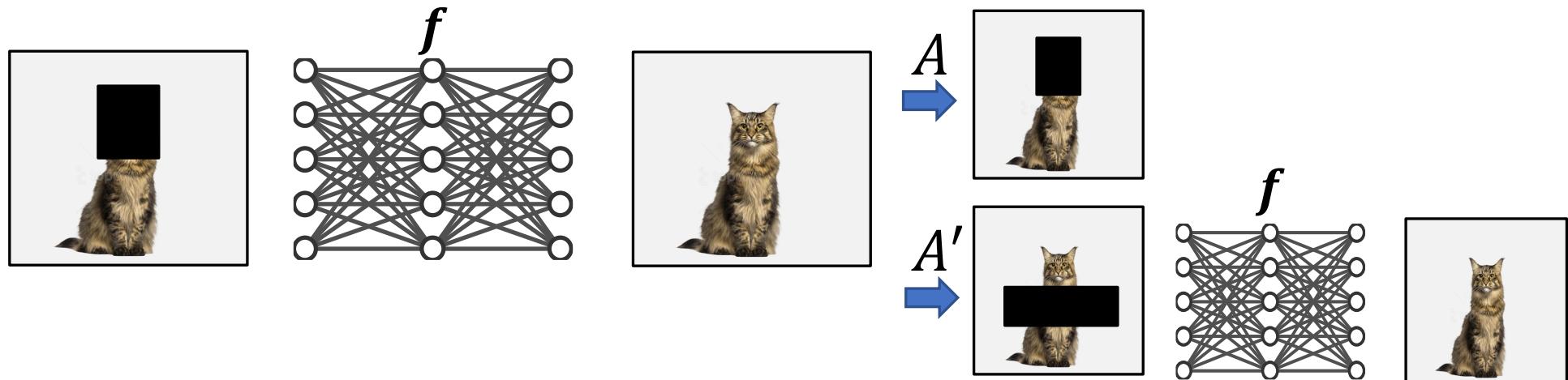
Can we use all measurements?

**Multi Operator Imaging (MOI)** [Tachella et al., NeurIPS 2022]

$$\mathcal{L}_{\text{MOI}}(\mathbf{y}, f) = \underbrace{\|\mathbf{y} - Af(\mathbf{y}, A)\|^2}_{\text{Replaced by SURE, R2R, etc.}} + \underbrace{\mathbb{E}_{A'} \|f(A'\hat{\mathbf{x}}, A') - \hat{\mathbf{x}}\|^2}_{\text{Enforces } f(A\mathbf{x}, A) \approx f(A'\mathbf{x}, A')}$$

Replaced by SURE, R2R, etc.

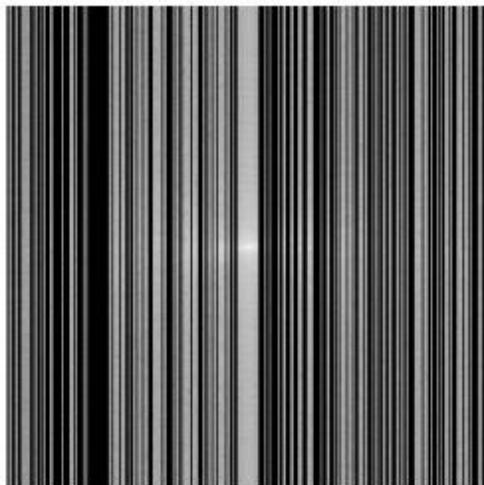
Enforces  $f(A\mathbf{x}, A) \approx f(A'\mathbf{x}, A')$



# Magnetic Resonance Imaging

- Unrolled network
- FastMRI dataset (single coil)
- $A_g$  are subsets of Fourier measurements (x4 downsampling)

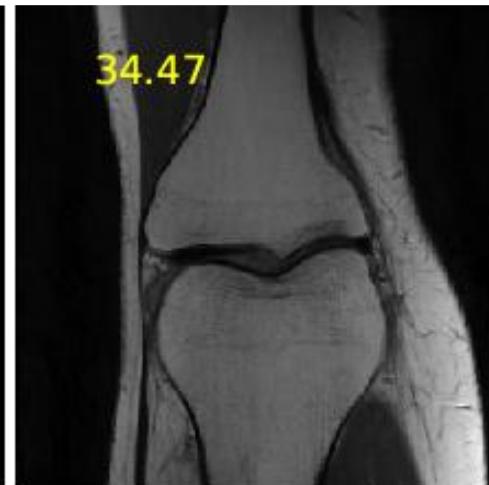
Measurements  $y$



Signal  $x$



Supervised



Measurement Splitting

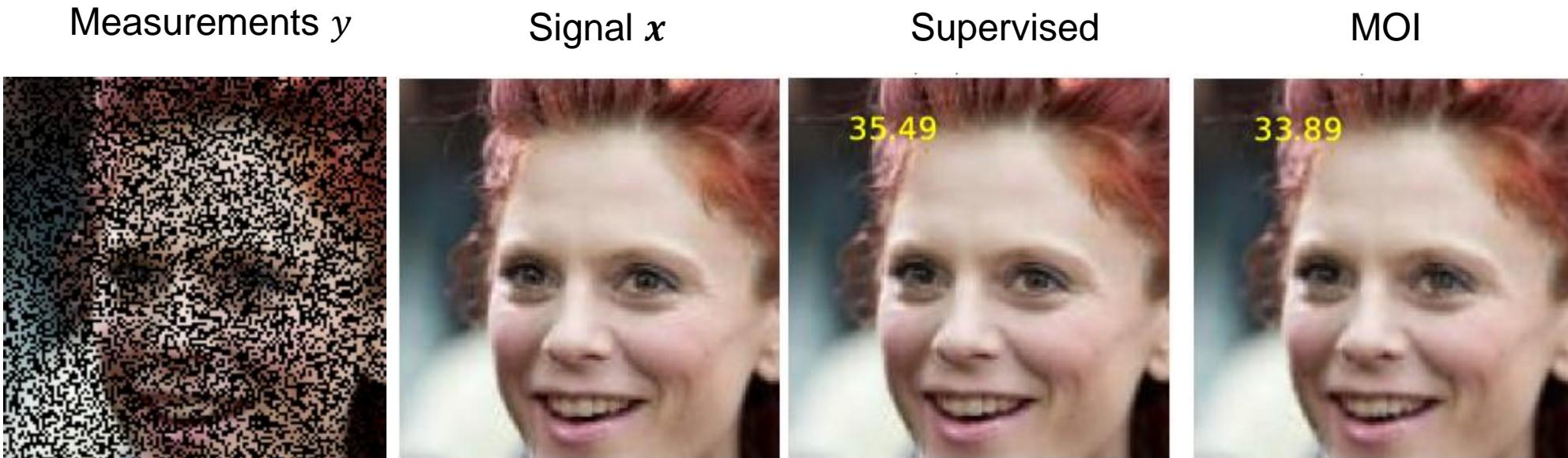


MOI



# Inpainting

- U-Net network
- CelebA dataset
- $A_g$  are inpainting masks



# Part 4: Learning with equivariance

# Symmetry Prior

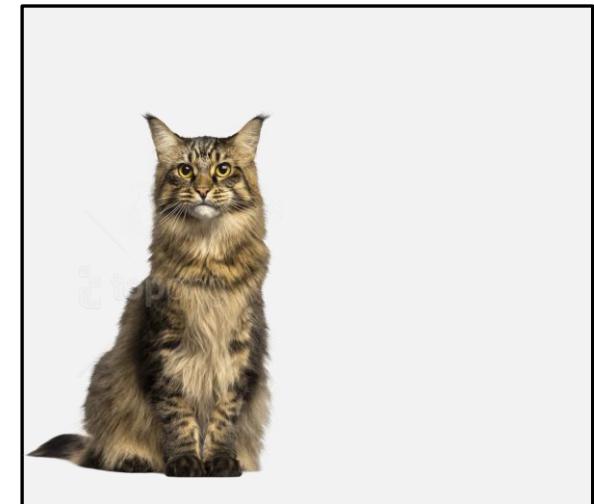
**Idea:** Most natural signals sets  $\mathcal{X}$  are invariant to groups of transformations.

*Example:* natural images are translation invariant

- Mathematically, a set  $\mathcal{X}$  is invariant to  $\{T_g \in \mathbb{R}^{n \times n}\}_{g \in G}$  if

$$\forall x \in \mathcal{X}, \forall g \in G, T_g x \in \mathcal{X}$$

**Other symmetries:** rotations, permutation, amplitude



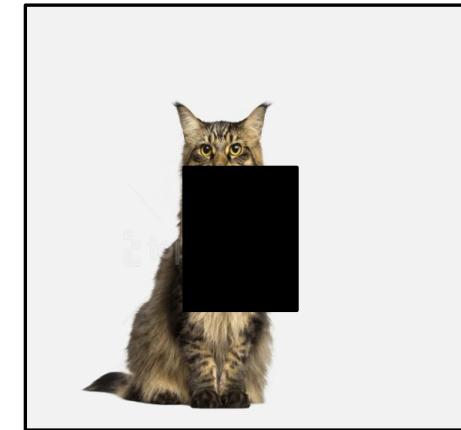
# Symmetry Prior

**Equivariant Imaging** [Chen, Davies and Tachella, ICCV 2021]

For all  $g \in G$  we have

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

$\overbrace{\quad\quad\quad}^{\mathbf{x}'}$   
 $\underbrace{\quad\quad\quad}_{A_g}$



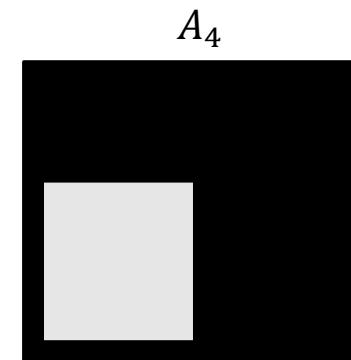
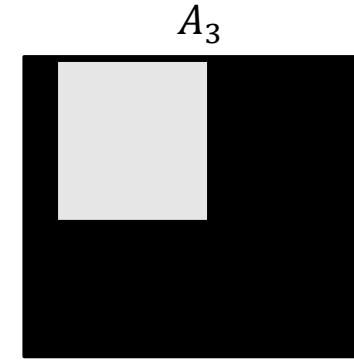
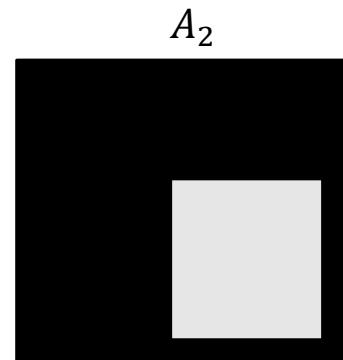
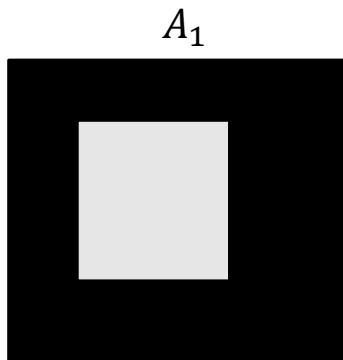
- We get multiple virtual operators  $\{A_g\}_{g \in G}$  ‘for free’!
- Each  $AT_g$  might have a different nullspace

# Necessary condition

**Proposition [T. et al., 2023]:** Learning reconstruction mapping  $f$  from observed measurements possible only if

$$\text{rank}(\mathbb{E}_g T_g^\top A^\top A T_g) = n,$$

and thus if  $m \geq \max_j \frac{c_j}{s_j} \geq \frac{n}{|G|}$  where  $s_j$  and  $c_j$  are dimension and multiplicity of irreps.



# (Non)-Equivariant Operators

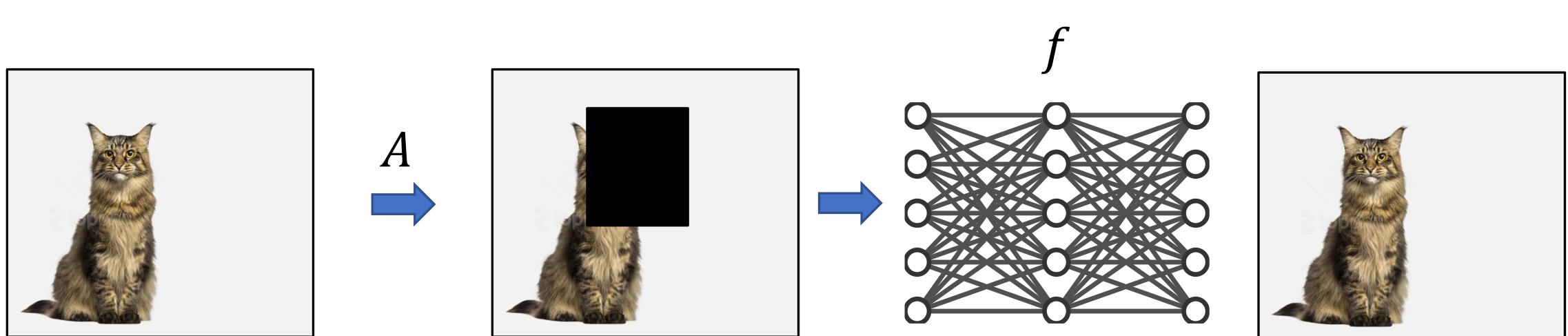
**Theorem** [T. et al., 2023]: The full rank condition requires that  $A$  **is not equivariant**:  $AT_g \neq \tilde{T}_g A$

$$\text{rank}(\mathbb{E}_g T_g^\top A^\top A T_g) = \text{rank}(A^\top (\mathbb{E}_g \tilde{T}_g^\top \tilde{T}_g) A) = \text{rank}(A^\top A) = m < n$$

# 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

How can we enforce equivariance in practice?

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

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

**Proposition** [Tachella & Pereyra, 2024]: *For linear and measurement consistent  $Af(Ax) = Ax$  reconstruction, we have*

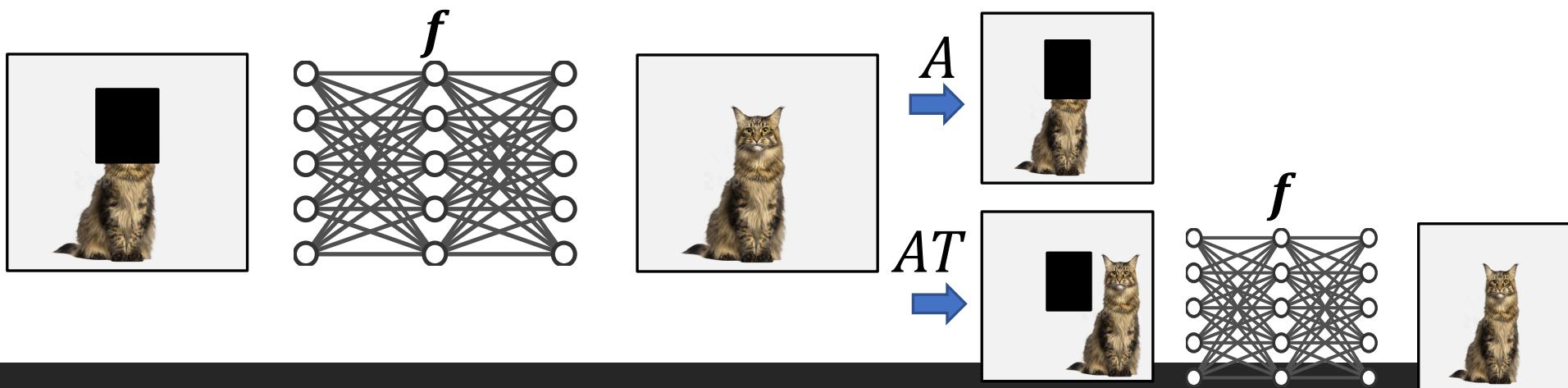
$$\mathcal{L}_{\text{EI}}(\mathbf{y}, f) = \| \mathbf{x} - f(\mathbf{y}) \|^2 + \text{bias}$$

where the **bias** term is small if  $f \circ A$  is **not** equivariant.

# Equivariant Imaging

**Robust Equivariant Imaging** [Chen, Tachella and Davies, CVPR 2022]

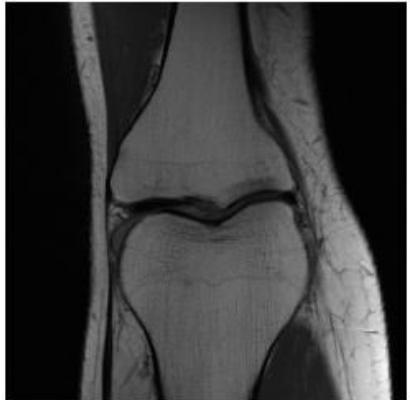
$$\mathcal{L}_{\text{REI}}(\mathbf{y}, f) = \underbrace{\|f(\mathbf{y}) - \mathbf{x}\|^2}_{\text{Replaced by SURE, R2R, etc}} + \underbrace{\mathcal{L}_{\text{EI}}(\mathbf{y}, f)}_{\text{enforces equivariance of } f \circ A}$$



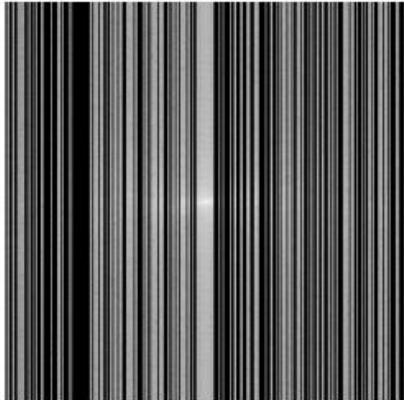
# MRI

- Operator  $A$  is a subset of Fourier measurements (x2 downsampling)
- Dataset is approximately **rotation invariant**

Signal  $x$

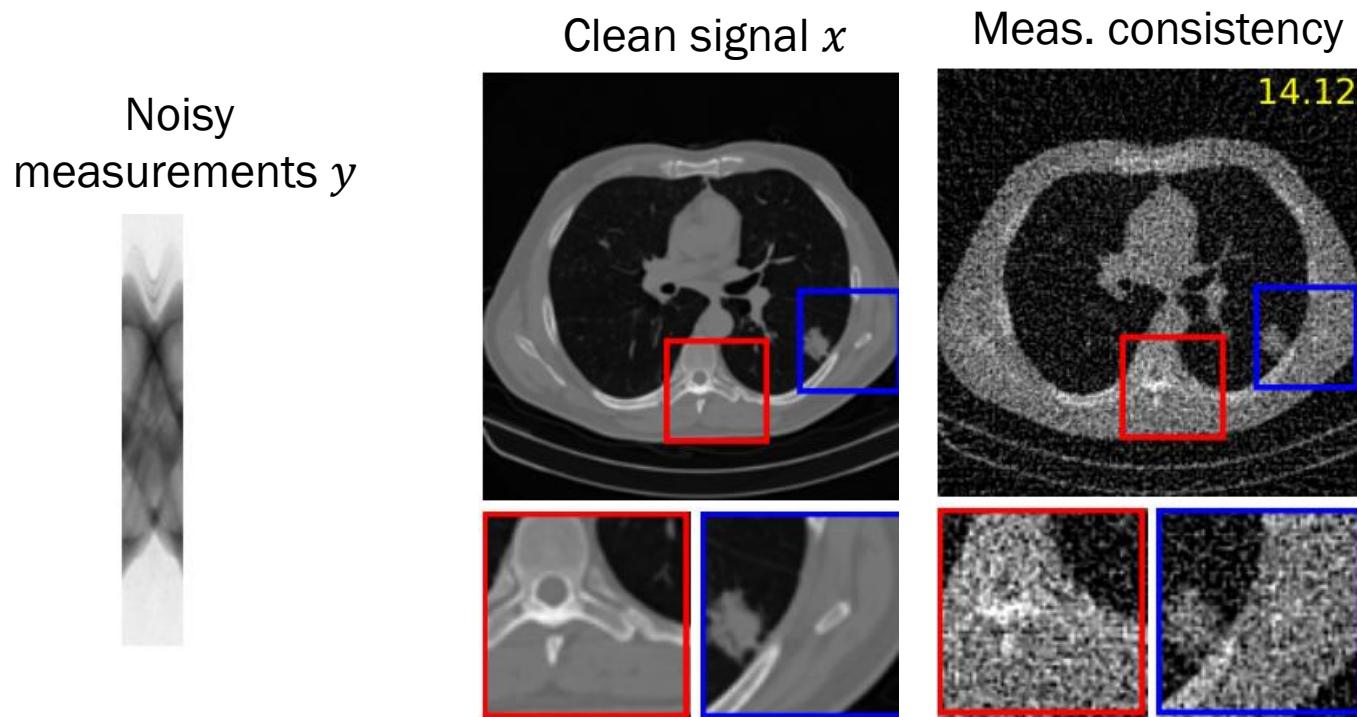


Measurements  $y$



# Computed Tomography

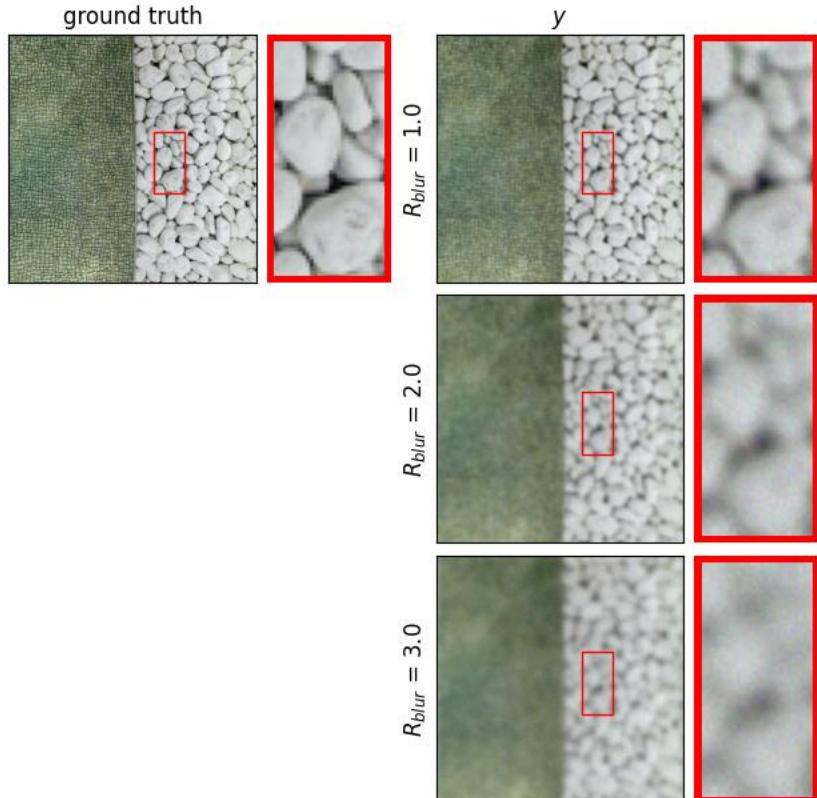
- Operator  $A$  is (non-linear variant) sparse radon transform
- Mixed Poisson-Gaussian noise
- Dataset is approximately **rotation invariant**



Chen, T., Davies, CVPR 2022

# Image Deblurring

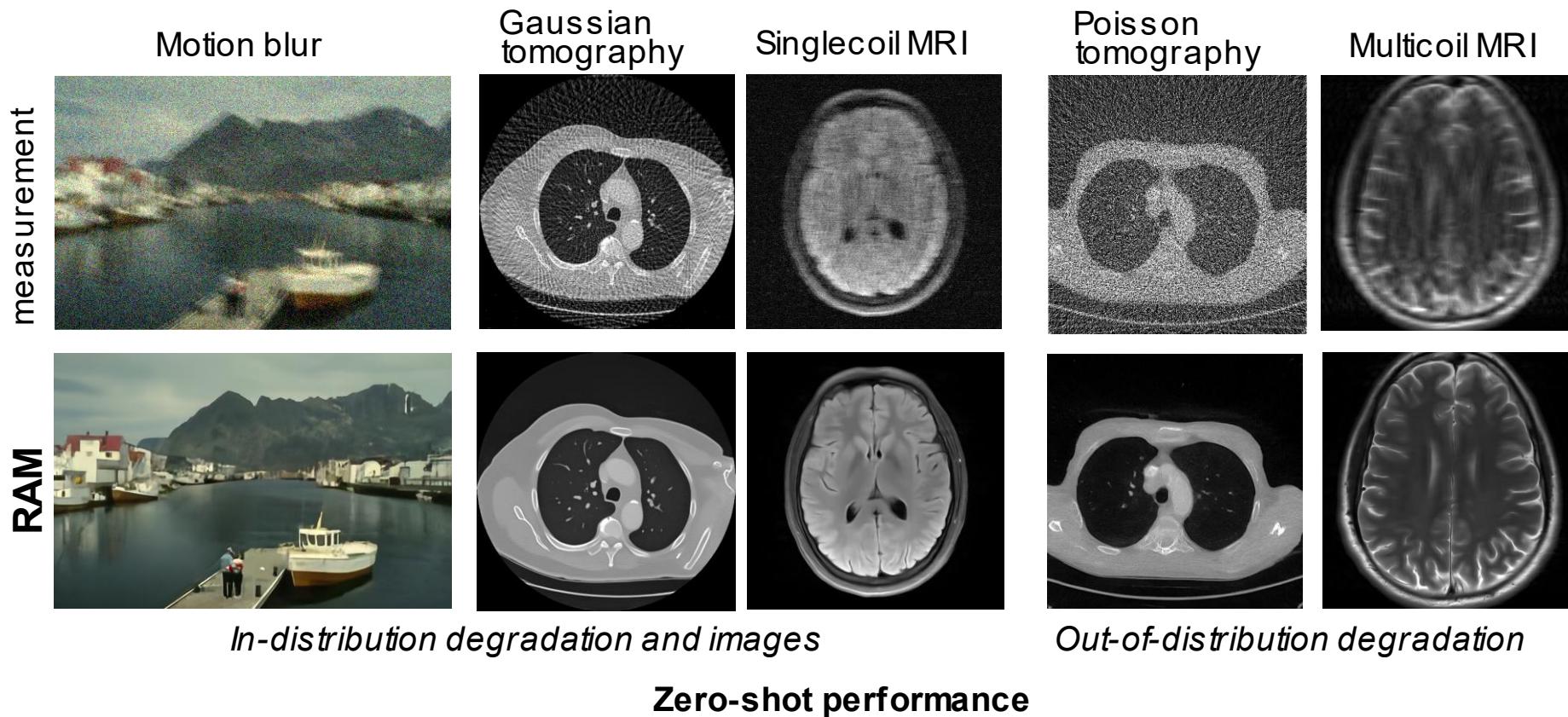
- Operator  $A$  is isotropic blur with Gaussian noise
- Dataset is approximately **scale invariant**



# Bonus: Finetuning foundation models

# Reconstruct Anything Model

- We trained a model that can solve many inverse problems at once [Terris et al., 2025]



# Finetuning

- The model can be finetuned with self-supervised losses on up to a single  $y$  ( $N = 1$ )
- Finetuning can be done in a **few seconds**

	Compressed Sensing	DRUNet	RAM	Reference
$A^\top y$				
Demosaicing				
Compressed Sensing		Demosaicing		
$N = 1$	43	60		
$N = 10$	80	107		
$N = 100$	228	267		

Table 5. Self-supervised finetuning time in seconds.

# Conclusions

Self-supervised learning for imaging problems

- **Theory:** Necessary & sufficient conditions for learning
  - Unbiased risk estimators
  - Number of measurements
  - Interplay between forward operator & data invariance
- **Practice:** self-supervised losses
  - Can be applied to any model (including foundation ones!)
  - Losses can be combined together

## LINEAR IMAGING



- Poor reconstructions
- Cannot handle noise/missing data



Quickstart [Examples](#) User Guide API Finding Help More ▾

Home > Examples

## Section Navigation

- Basics
- Optimization
- Plug-and-Play
- Sampling
- Unfolded
- Patch Priors
- Self-Supervised Learning
- Adversarial Learning
- Advanced

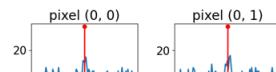
# Examples

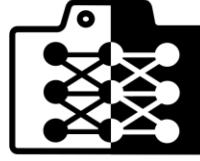
All the examples have a download link at the end. You can load the example's notebook on [Google Colab](#) and run them by adding the line

```
pip install git+https://github.com/deepinv/deepinv.git#egg=deepinv
```

to the top of the notebook (e.g., [as in here](#)).

# Basics





*Deep  
Inverse*



*inria*



PSL



**EPFL**



# References



**Paper references:**

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

**Code examples:**

<https://andrewwango.github.io/deepinv-selfsup-fastmri/demo>

[https://deepinv.github.io/deepinv/auto\\_examples/self-supervised-learning/index.html](https://deepinv.github.io/deepinv/auto_examples/self-supervised-learning/index.html)

**YouTube version (3 hours):**

<https://youtu.be/gf-WCHXAdfk?si=bRC6Pq0WpZHNrRLU>