

# Plug-and-Play Models for Large-Scale Computational Imaging

**Ulugbek S. Kamilov**

Computational Imaging Group (CIG)

Twitter: [@ukmlv](https://twitter.com/ukmlv)

# Big Thank You to the members of the WashU **Computational Imaging Group (CIG)**

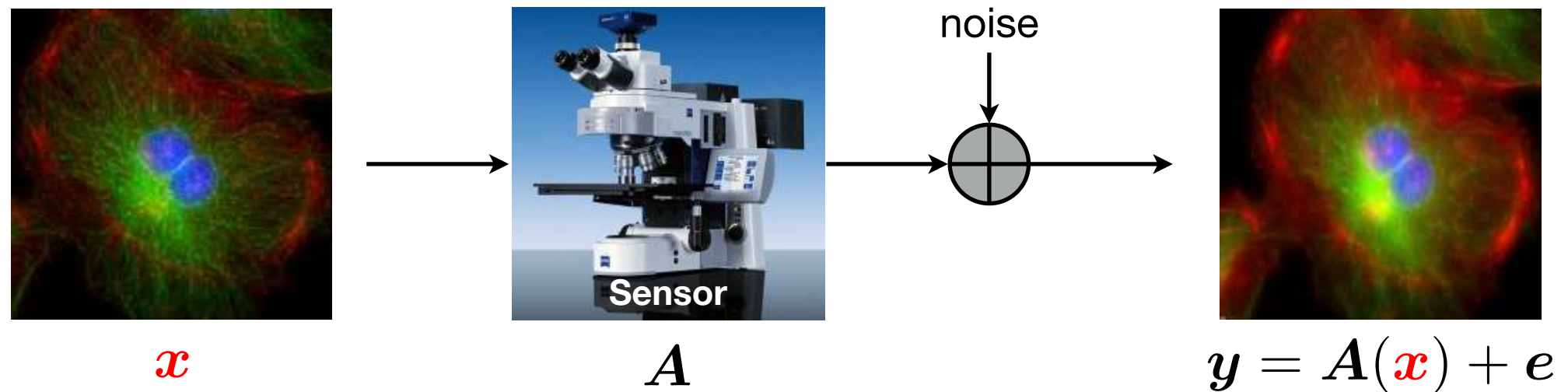


**From left to right in the photo:**

Weijie Gan, Tomas Kerepecky (alum), Xiaojian Xu (alum), Jiaming Liu,  
Flora Sun, Yu Sun (alum), Yuyang Hu, and Shirin Shoushtari.

Many **computational imaging** problems  
can be formulated as **inverse problems**

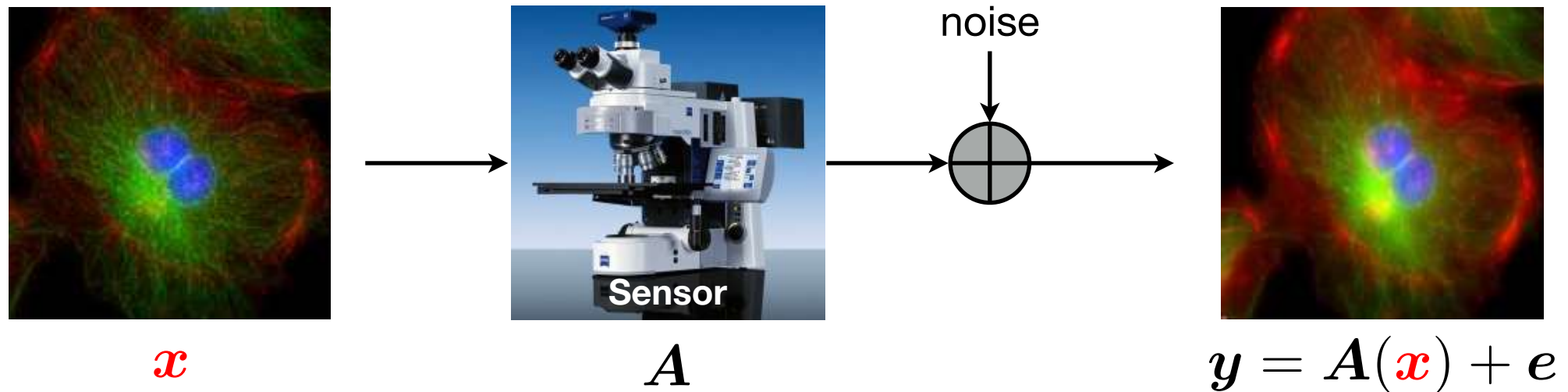
# Many computational imaging problems can be formulated as inverse problems



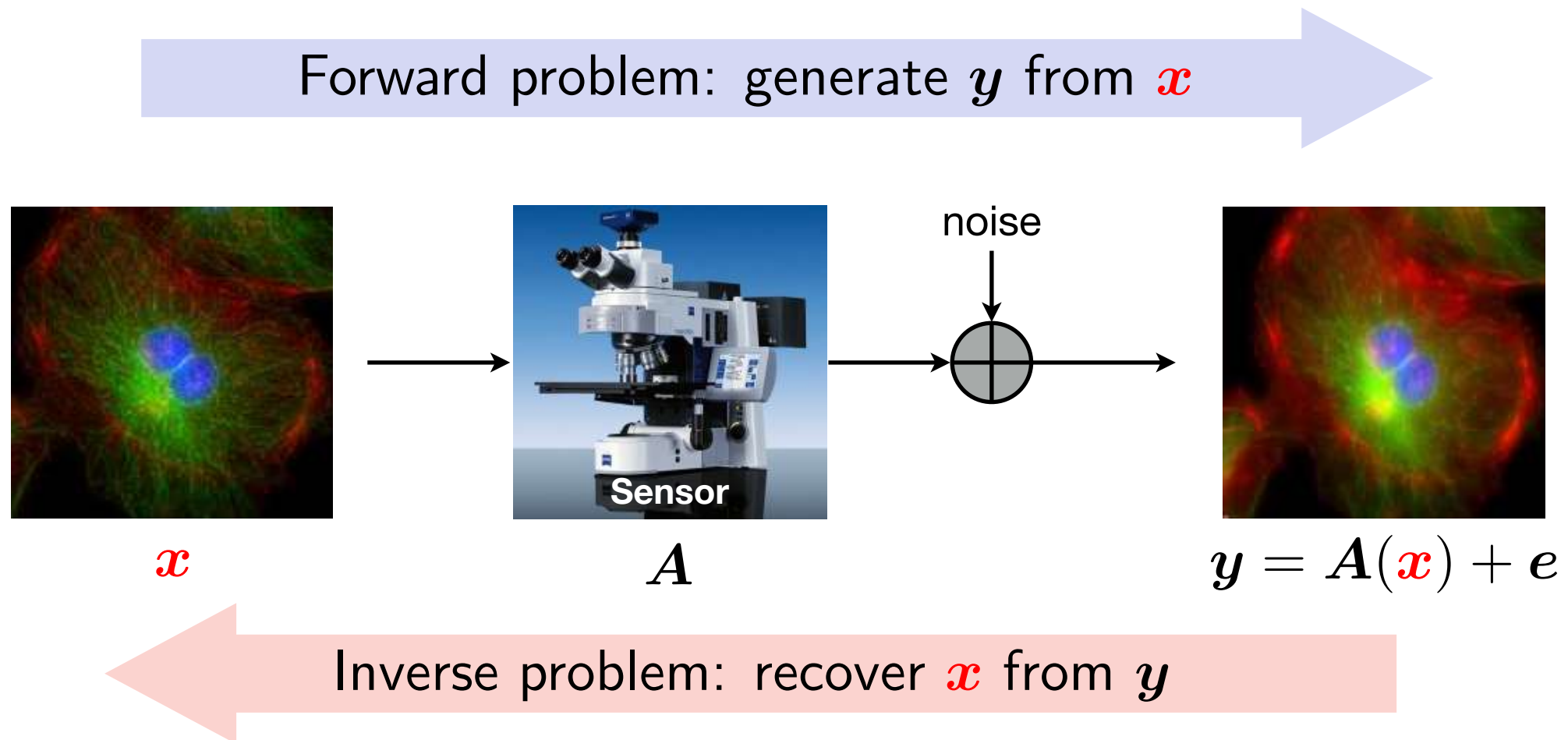


# Many computational imaging problems can be formulated as inverse problems

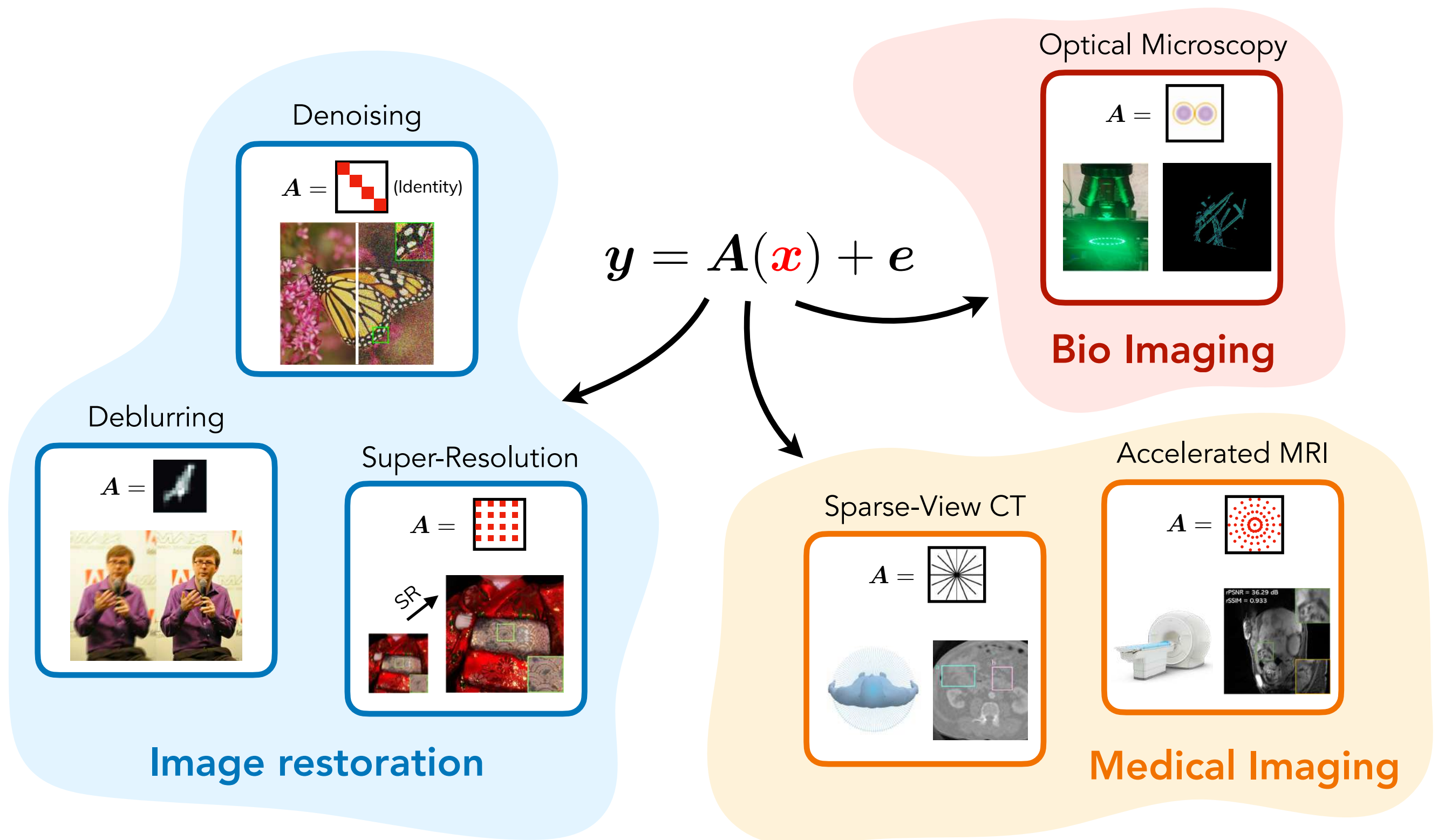
Forward problem: generate  $y$  from  $x$



# Many computational imaging problems can be formulated as inverse problems



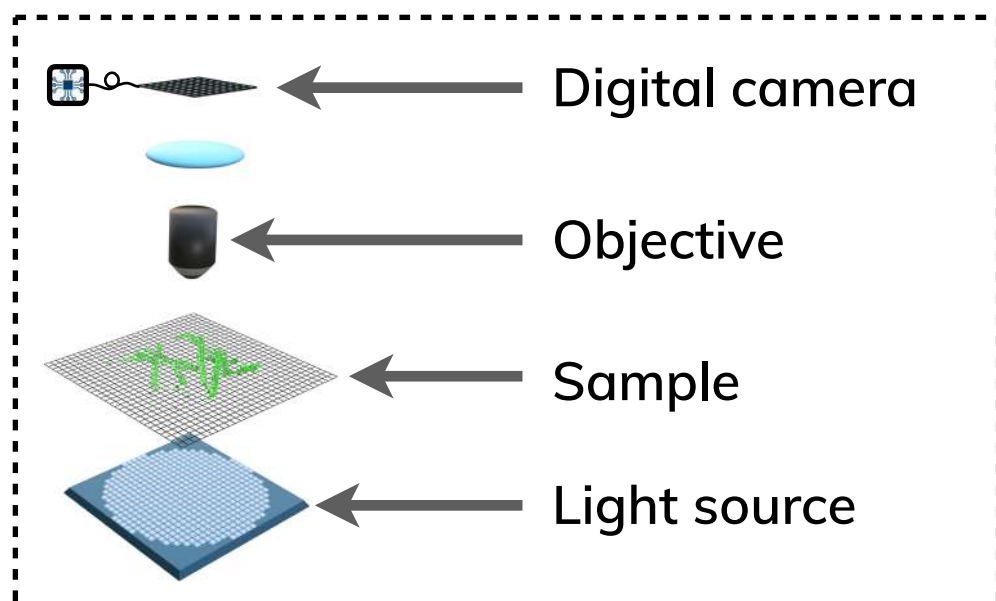
# Many computational imaging problems can be formulated as inverse problems



**Example:** Intensity diffraction tomography (IDT)  
collects intensity measurements of scattered light



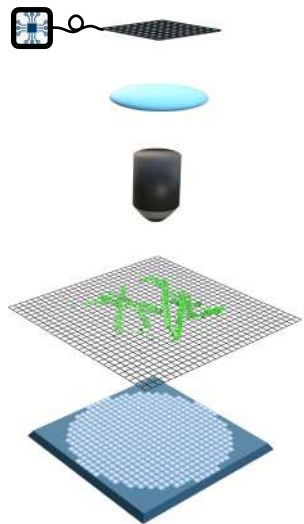
# Example: Intensity diffraction tomography (IDT) collects intensity measurements of scattered light



IDT is a relatively cheap and simple optical microscopy system

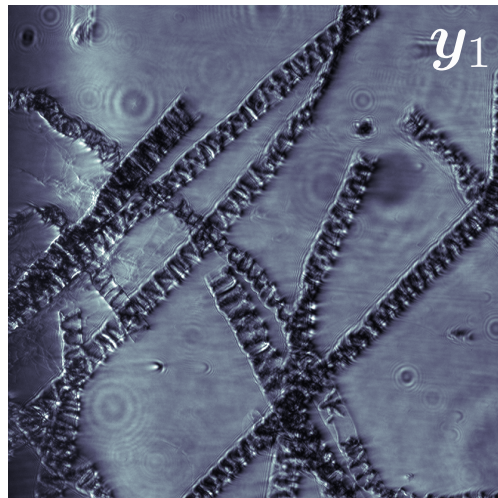
Liu et al, "Recovery of continuous 3D refractive index maps from discrete intensity-only measurements using neural fields," *Nature Machine Intelligence*, 2022

# Example: Intensity diffraction tomography (IDT) collects intensity measurements of scattered light

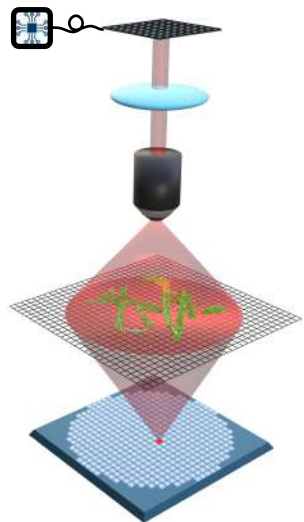


Liu et al, "Recovery of continuous 3D refractive index maps from discrete intensity-only measurements using neural fields," *Nature Machine Intelligence*, 2022

# Example: Intensity diffraction tomography (IDT) collects intensity measurements of scattered light

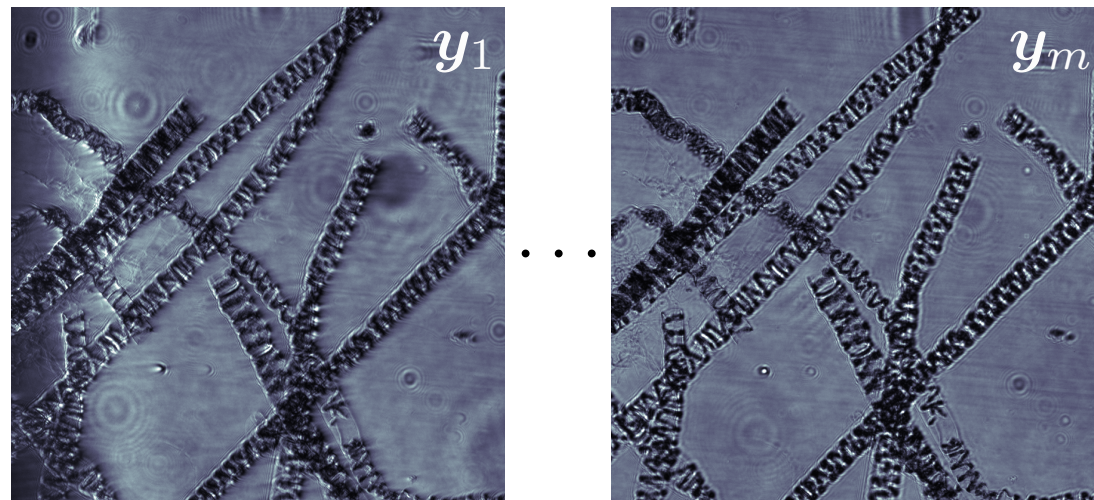


intensity image 1



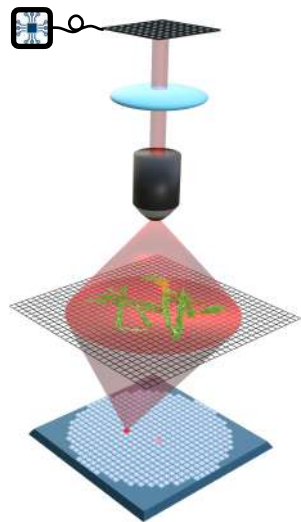
Liu et al, "Recovery of continuous 3D refractive index maps from discrete intensity-only measurements using neural fields," *Nature Machine Intelligence*, 2022

# Example: Intensity diffraction tomography (IDT) collects intensity measurements of scattered light



intensity image 1

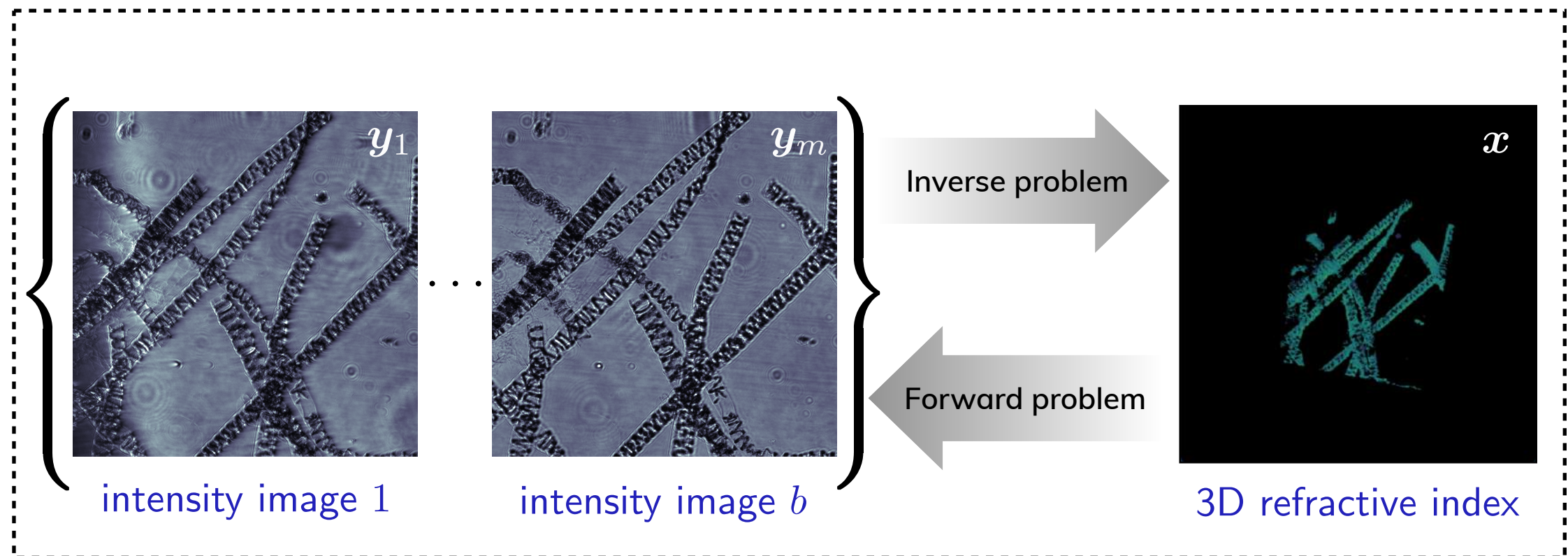
intensity image  $b$



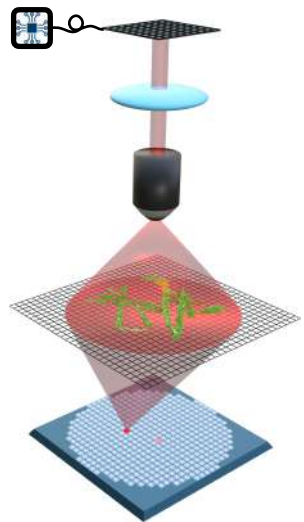
Liu et al, "Recovery of continuous 3D refractive index maps from discrete intensity-only measurements using neural fields," *Nature Machine Intelligence*, 2022



# Example: Intensity diffraction tomography (IDT) collects intensity measurements of scattered light

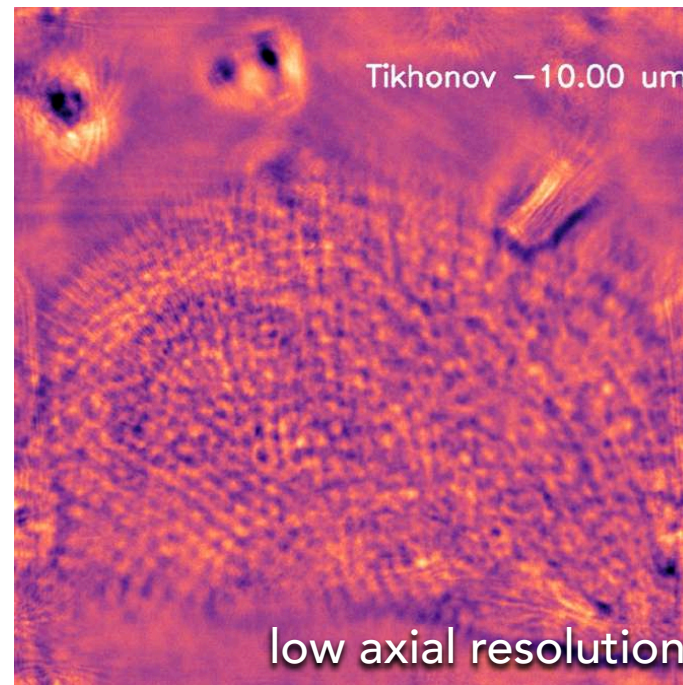


IDT is a data-intensive limited-angle tomography under **light-scattering** and **phase-loss**



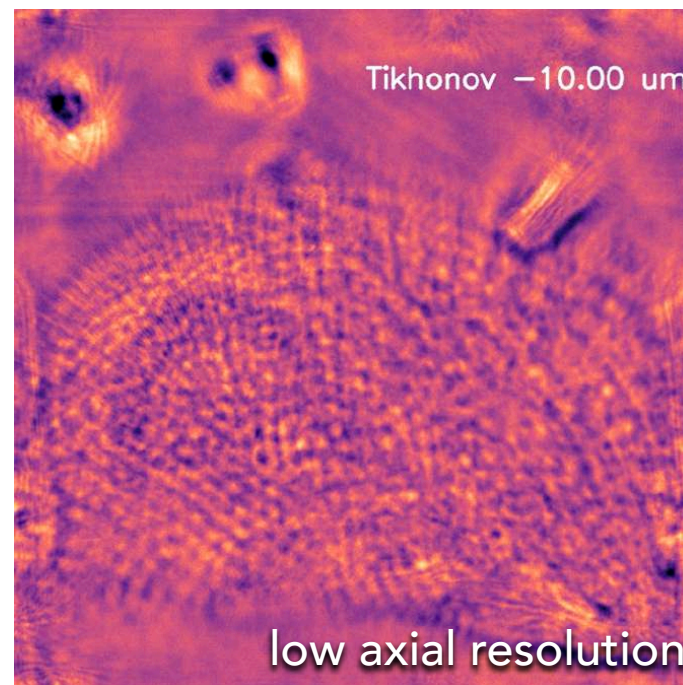
# Classical issues in the context of **inverse problems**: **costly acquisition, imaging artifacts, and big data**

# Classical issues in the context of inverse problems: costly acquisition, imaging artifacts, and big data





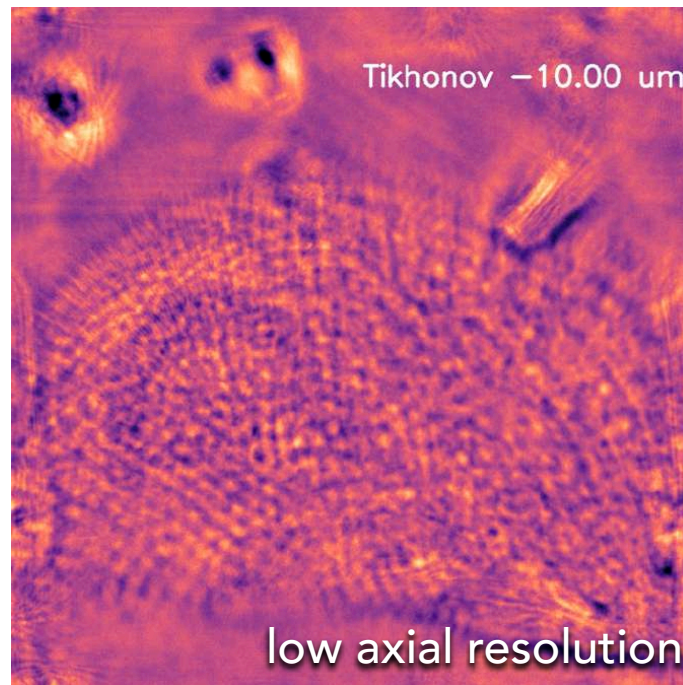
# Classical issues in the context of inverse problems: costly acquisition, imaging artifacts, and big data



**Challenge #1: Acquisition is too slow or costly:**  
Due to sequential and indirect acquisition of data



# Classical issues in the context of inverse problems: costly acquisition, imaging artifacts, and big data



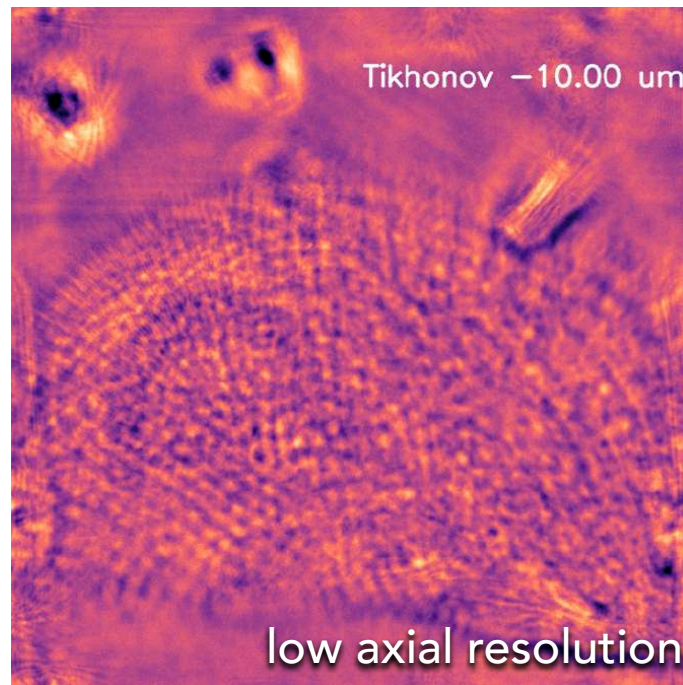
## **Challenge #1: Acquisition is too slow or costly:**

Due to sequential and indirect acquisition of data

## **Challenge #2: Reconstructed images contain artifacts:**

Due to undersampling, motion, model mismatch, and noise

# Classical issues in the context of inverse problems: costly acquisition, imaging artifacts, and big data



## **Challenge #1: Acquisition is too slow or costly:**

Due to sequential and indirect acquisition of data

## **Challenge #2: Reconstructed images contain artifacts:**

Due to undersampling, motion, model mismatch, and noise

## **Challenge #3: High computational/memory requirements:**

Due to large volumes of data to process in 3D, 4D, or 5D

# Outline for the rest of the talk

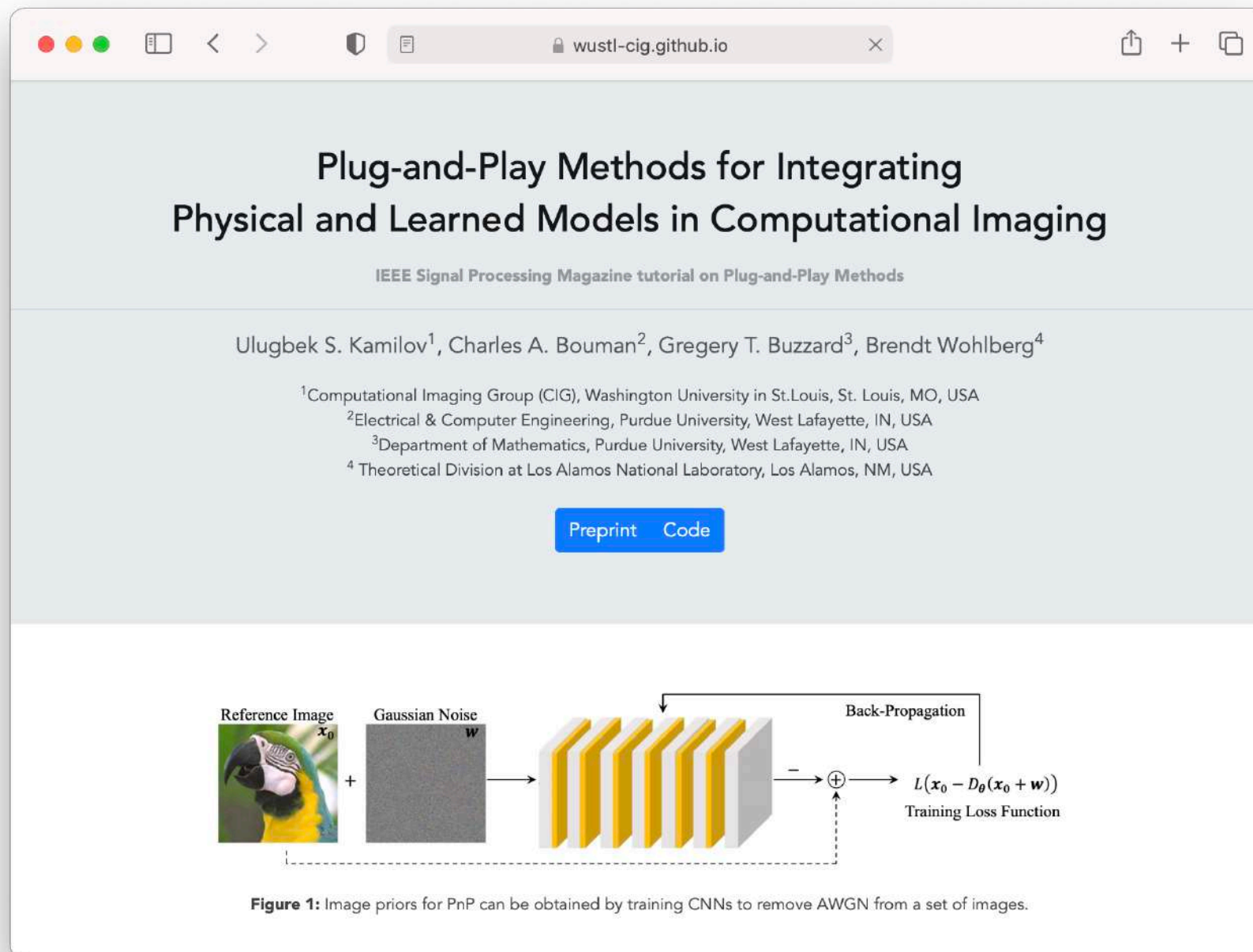
- Plug-and-Play Methods for Inverse Problems (**IEEE SPM 2022**)  
Integrating physical models and learned deep priors
- Online Deep Equilibrium Learning (**NeurIPS 2022**)  
A new PnP framework for efficient prior learning
- Deep Continuous Artifact-Free Fields (**Nature MI 2022**)  
A new PnP framework for continuous image recovery

# Outline for the rest of the talk

- **Plug-and-Play Methods for Inverse Problems (IEEE SPM 2022)**  
Integrating physical models and learned deep priors
- Online Deep Equilibrium Learning (NeurIPS 2022)  
A new PnP framework for efficient prior learning
- Deep Continuous Artifact-Free Fields (Nature MI 2022)  
A new PnP framework for continuous image recovery



# Plug-and-Play Methods (PnP) tutorial for IEEE Signal Processing Magazine is available online



Scan the QR code  
to see the tutorial

**Bayesian perspective** is common, but the view on how to best represent image priors is still evolving

# Bayesian perspective is common, but the view on how to best represent image priors is still evolving

## Probabilistic interpretation of an inverse problem

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \mathbf{e}$$

forward model

$$\mathbf{x} \sim p_{\mathbf{x}}$$

image prior

$$\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

noise model

# Bayesian perspective is common, but the view on how to best represent image priors is still evolving

## Probabilistic interpretation of an inverse problem

$$y = A(x) + e$$

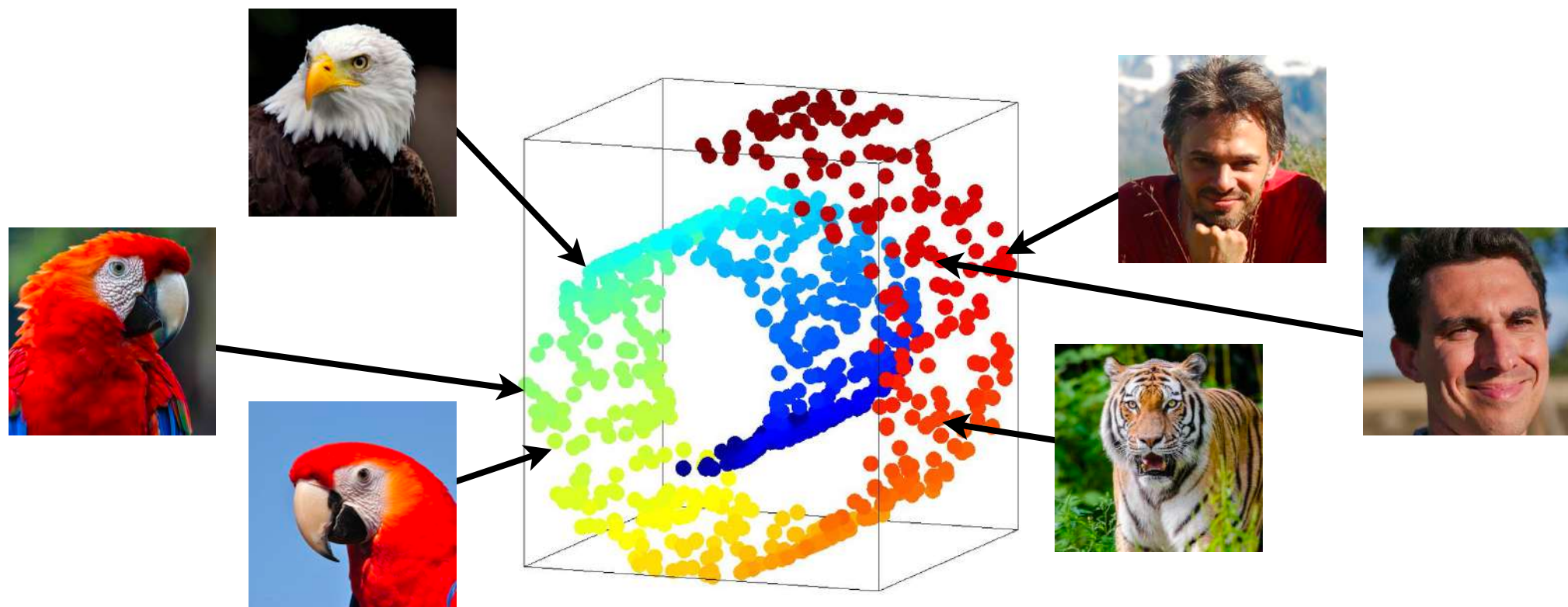
forward model

$$x \sim p_x$$

image prior

$$e \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

noise model



**Remark:** Generating images from  $p_x$  is equivalent to sampling from a lower-dimensional and non-linear subset  $\mathcal{X} \subset \mathbb{R}^n$  of a high-dimensional space

# Bayesian perspective is common, but the view on how to best represent image priors is still evolving

Probabilistic interpretation of an inverse problem

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \mathbf{e} \quad \mathbf{x} \sim p_{\mathbf{x}} \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

**MAP** and **MMSE** statistical estimators can be expressed as **model-based optimization**

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{A}(\mathbf{x})\|_2^2 + \sigma^2 h(\mathbf{x}) \right\}$$

$$h_{\text{MAP}}(\mathbf{x}) = -\log(p_{\mathbf{x}}(\mathbf{x}))$$

**Remark:** Maximum a posteriori probability (MAP) estimator returns an image that maximizes  $p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$



# Bayesian perspective is common, but the view on how to best represent image priors is still evolving

Probabilistic interpretation of an inverse problem

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \mathbf{e} \quad \mathbf{x} \sim p_{\mathbf{x}} \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

**MAP** and **MMSE** statistical estimators can be expressed as **model-based optimization**

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{A}(\mathbf{x})\|_2^2 + \sigma^2 h(\mathbf{x}) \right\}$$

$$h_{\text{MAP}}(\mathbf{x}) = -\log(p_{\mathbf{x}}(\mathbf{x}))$$

$$h_{\text{MMSE}}(\mathbf{x}) = \text{some expression}$$

**Remark:** Minimum mean squared error (MMSE) estimator returns the conditional mean  $\mathbb{E}[\mathbf{x}|\mathbf{y}]$ , which is the solution that maximizes the signal-to-noise ratio (SNR)

# Bayesian perspective is common, but the view on how to best represent image priors is still evolving

Probabilistic interpretation of an inverse problem

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \mathbf{e} \quad \mathbf{x} \sim p_{\mathbf{x}} \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

MAP and MMSE statistical estimators can be expressed as model-based optimization

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{A}(\mathbf{x})\|_2^2 + \sigma^2 h(\mathbf{x}) \right\} \quad \begin{aligned} h_{\text{MAP}}(\mathbf{x}) &= -\log(p_{\mathbf{x}}(\mathbf{x})) \\ h_{\text{MMSE}}(\mathbf{x}) &= \text{some expression} \end{aligned}$$

Image denoising is a special inverse problem where the forward model is an identity

$$\mathbf{D}_{\sigma}(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 + \sigma^2 h(\mathbf{x}) \right\}$$

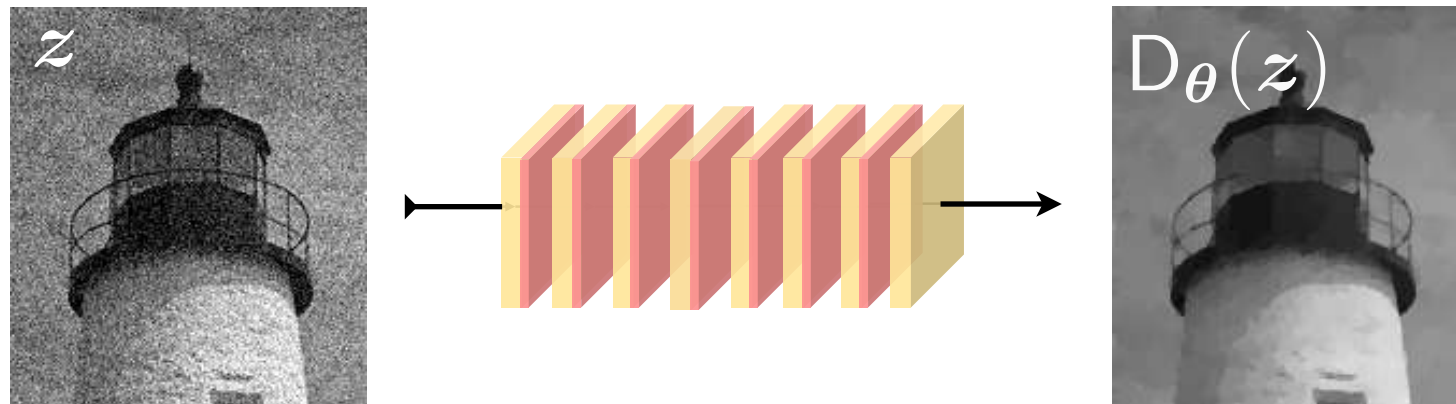
**Proximal operator** is a convenient proxy for the prior from the algorithmic perspective

**Plug-and-Play Methods (PnP)** are flexible deep models that use image denoisers as image priors

# Plug-and-Play Methods (PnP) are flexible deep models that use image denoisers as image priors

**Learned model:** Pre-trained **image denoising** neural network

$D_\theta$  : more noisy image  $\mapsto$  less noisy image



**Remark:** PnP is a self-supervised learning framework since the image prior is learned on a “pretext task”

# Plug-and-Play Methods (PnP) are flexible deep models that use image denoisers as image priors

Learned model: Pre-trained image denoising neural network

$D_\theta$  : more noisy image  $\mapsto$  less noisy image

**Physical model:** Infuses information from the **forward model**

$I - \gamma \nabla g$  : less measurement consistent  $\mapsto$  more measurement consistent

**Remark:** Smaller values of  $g(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}(\mathbf{x})\|_2^2$   
are more consistent with the measured data



# Plug-and-Play Methods (PnP) are flexible deep models that use image denoisers as image priors

Learned model: Pre-trained image denoising neural network

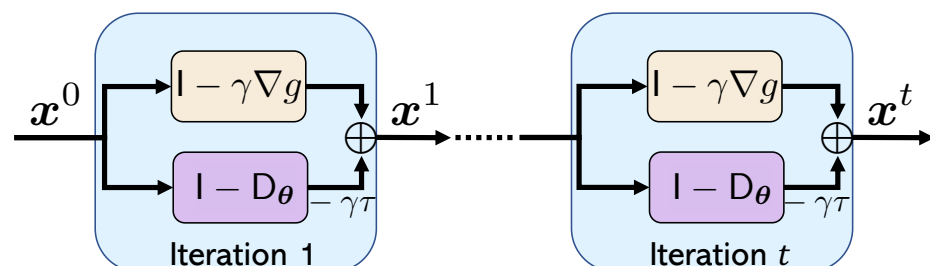
$D_\theta$  : more noisy image  $\mapsto$  less noisy image

Physical model: Infuses information from the forward model

$I - \gamma \nabla g$  : less measurement consistent  $\mapsto$  more measurement consistent

**PnP methods** integrate both models into a **deep model-based architecture (DMBA)** that can have infinitely many layers

$$x^t \leftarrow x^{t-1} - \gamma G(x^{t-1}) \quad G(x) := \nabla g(x) + \tau(x - D_\theta(x))$$



# Plug-and-Play Methods (PnP) are flexible deep models that use image denoisers as image priors

Learned model: Pre-trained image denoising neural network

$D_\theta$  : more noisy image  $\mapsto$  less noisy image

Physical model: Infuses information from the forward model

$I - \gamma \nabla g$  : less measurement consistent  $\mapsto$  more measurement consistent

PnP methods integrate both models into a deep model-based architecture (DMBA) that can have infinitely many layers

$$\mathbf{x}^t \leftarrow \mathbf{x}^{t-1} - \gamma \mathbf{G}(\mathbf{x}^{t-1}) \quad \mathbf{G}(\mathbf{x}) := \nabla g(\mathbf{x}) + \tau(\mathbf{x} - D_\theta(\mathbf{x}))$$

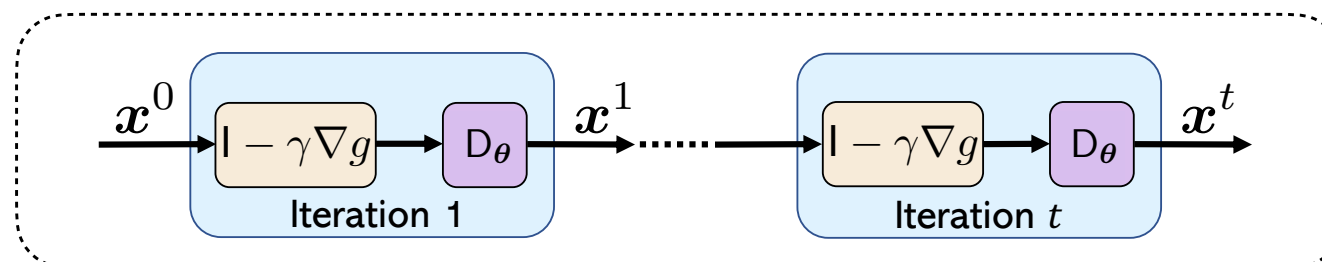
PnP methods have become influential in the context of inverse problems due to their mathematical elegance, flexibility, and performance

**PnP** does not have to be restricted to **PnP-ADMM**,  
**MAP**, **implicit priors**, **denoisers**, and **point estimates**

# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)

## Example: Architecture of PnP-ISTA



# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)
- While original PnP is interpreted as using implicit priors, there are formulations based on explicit regularizers (such as RED)

**Example:** RED regularizer and its gradient

$$h(\boldsymbol{x}) = \boldsymbol{x}^T (\boldsymbol{x} - \mathbf{D}_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \Rightarrow \quad \nabla h(\boldsymbol{x}) = \boldsymbol{x} - \mathbf{D}_{\boldsymbol{\theta}}(\boldsymbol{x})$$



# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)
- While original PnP is interpreted as using implicit priors, there are formulations based on explicit regularizers (such as RED)
- While the MAP view of the proximal operator is common, CNNs are often trained to act as MMSE estimators over the training data

$$D_{\sigma}^*(z) = \arg \min_D \mathbb{E} [\|x - D(z)\|_2^2]$$

MMSE estimator  
for  $z = x + e$

# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)
- While original PnP is interpreted as using implicit priors, there are formulations based on explicit regularizers (such as RED)
- While the MAP view of the proximal operator is common, CNNs are often trained to act as MMSE estimators over the training data

**Theorem:** PnP-ISTA using  $D_{\sigma}^*$  monotonically converges to a stationary point of the following composite function

$$f(\mathbf{x}) = g(\mathbf{x}) + h_{\text{mmse}}(\mathbf{x})$$

- generally not convex
- $\neq h_{\text{map}}(\mathbf{x}) = -\log(p_{\mathbf{x}}(\mathbf{x}))$
- MAP estimator for the prior:  $p(\mathbf{x}) \propto \exp(-h_{\text{mmse}}(\mathbf{x}))$

# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

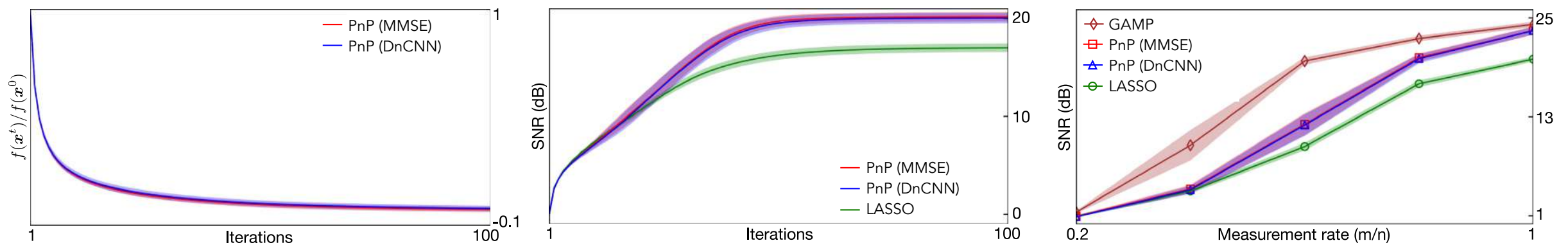
- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)
- While original PnP is interpreted as using implicit priors, there are formulations based on explicit regularizers (such as RED)
- While the MAP view of the proximal operator is common, CNNs are often trained to act as MMSE estimators over the training data

$$\|D_{\sigma}(z) - D_{\sigma}^*(z)\|_2 \leq \sigma\varepsilon \quad \Rightarrow \quad \min_{1 \leq i \leq t} \|G^*(x^i)\|_2^2 \leq \frac{C_1}{t} + \tau\sigma\varepsilon C_2$$

**Remark:** Even when the denoiser is not exact,  
PnP converges up to an error bound!

# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)
- While original PnP is interpreted as using implicit priors, there are formulations based on explicit regularizers (such as RED)
- While the MAP view of the proximal operator is common, CNNs are often trained to act as MMSE estimators over the training data

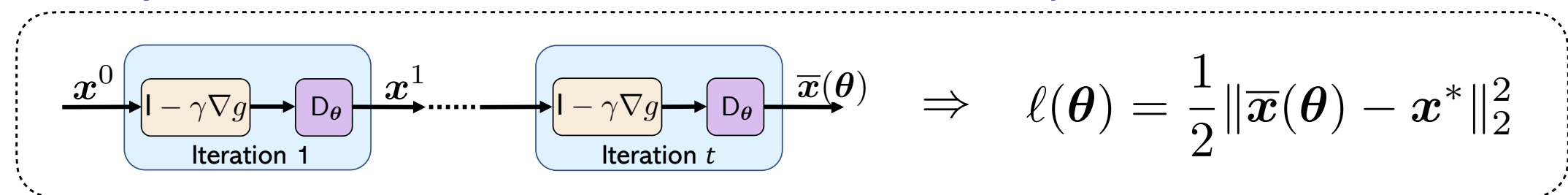


Note how DnCNN, trained to approximate the MMSE denoiser, perfectly agrees with the results using the true MMSE denoiser!

# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)
- While original PnP is interpreted as using implicit priors, there are formulations based on explicit regularizers (such as RED)
- While the MAP view of the proximal operator is common, CNNs are often trained to act as MMSE estimators over the training data
- While original PnP was restricted to image denoisers, there are tools for efficiently learning priors end-to-end (such as DEQ)

**Example:** DEQ formulation for the end-to-end optimization of PnP-ISTA





# PnP does not have to be restricted to PnP-ADMM, MAP, implicit priors, denoisers, and point estimates

- While PnP-ADMM was the first PnP method, there are variants based on other algorithms (such as PnP-FISTA)
- While original PnP is interpreted as using implicit priors, there are formulations based on explicit regularizers (such as RED)
- While the MAP view of the proximal operator is common, CNNs are often trained to act as MMSE estimators over the training data
- While original PnP was restricted to image denoisers, there are tools for efficiently learning priors end-to-end (such as DEQ)
- While original PnP was restricted to point estimates, there are tools for efficiently sampling from priors (such as PnP-ULA)

# Training PnP denoisers as **diffusion models** can enable sampling realistic images with little data

Sinogram ( $\theta \leq 60^\circ$ )

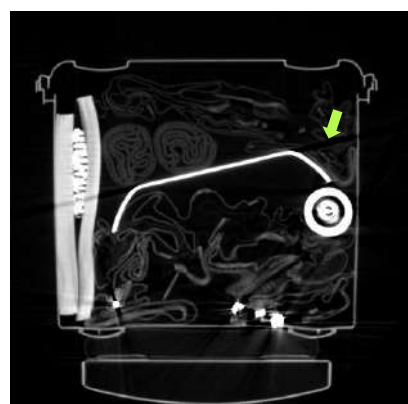
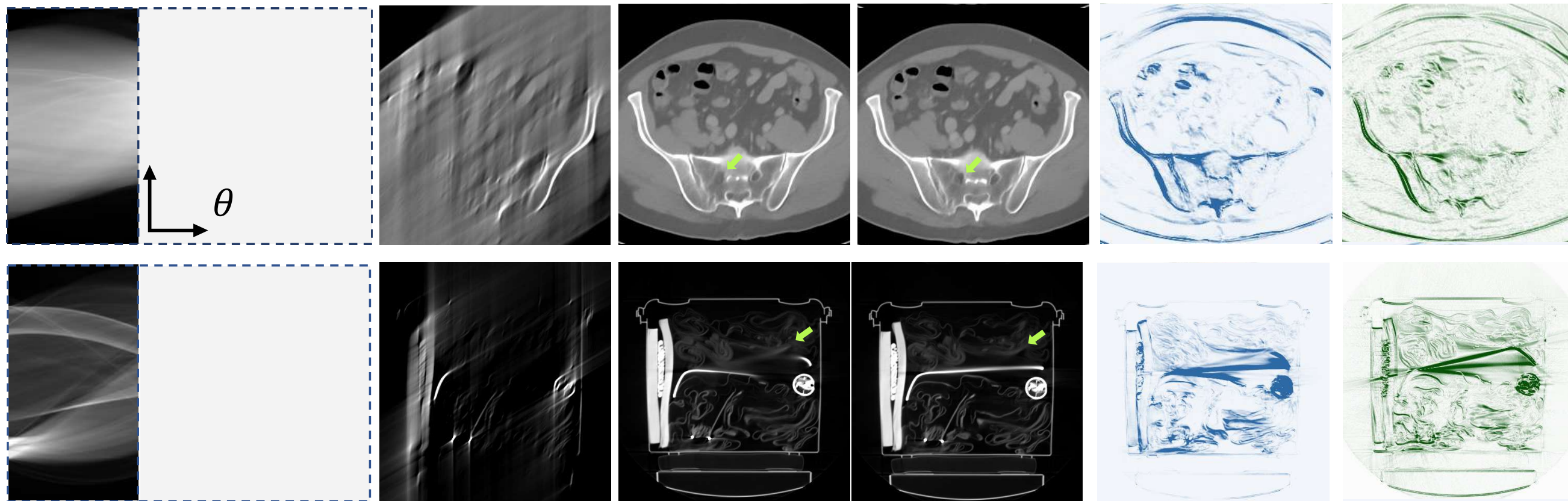
FBP

Sample #1

Sample #2

Variance

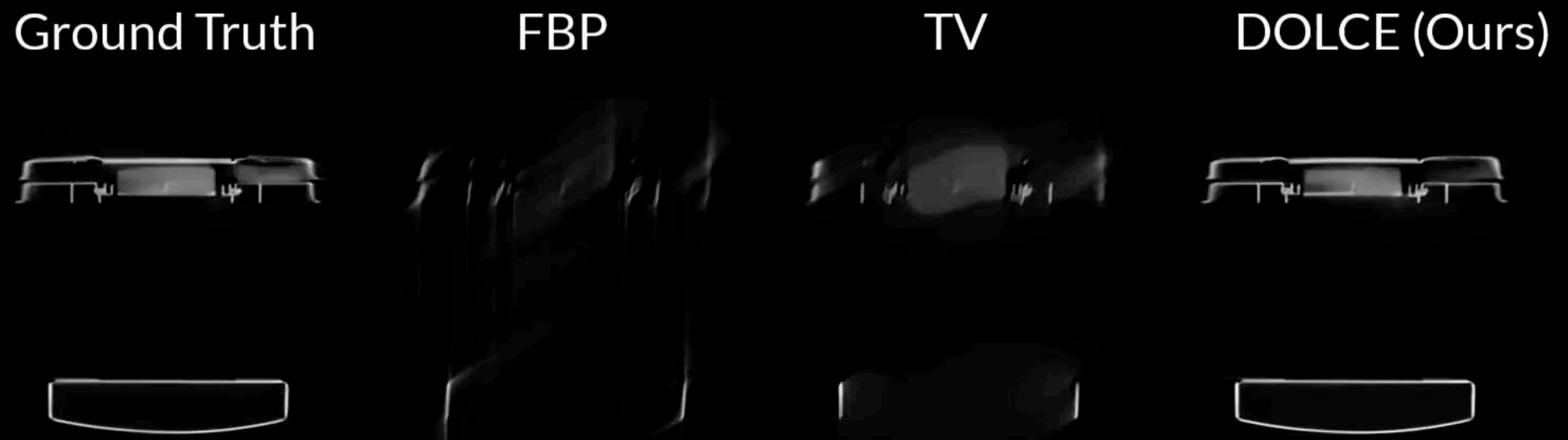
Abs. Error



Ground-truth suitcase

Liu et al, "DOLCE: A Model-Based Probabilistic Diffusion Framework for Limited-Angle CT Reconstruction," [arXiv:2211.12340](#), 2022

# Training PnP denoisers as diffusion models can enable sampling realistic images with little data

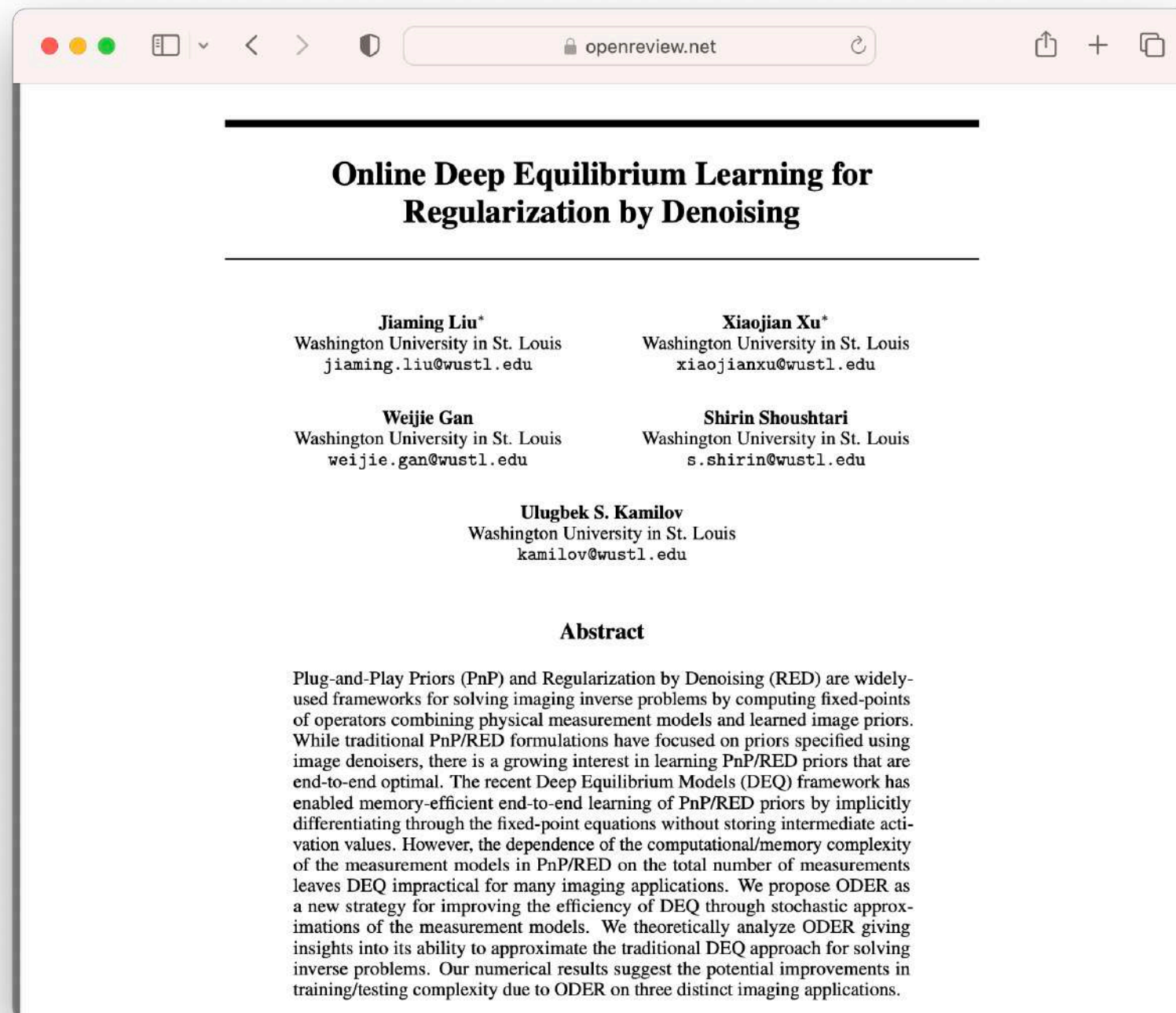


Limited-angle CT with  $\theta_{\max} \leq 60^\circ$ . Note how DOLCE can synthesize realistic-looking 3D reconstruction from very little data!

# Outline for the rest of the talk

- Plug-and-Play Methods for Inverse Problems (IEEE SPM 2022)  
Integrating physical models and learned deep priors
- **Online Deep Equilibrium Learning (NeurIPS 2022)**  
A new PnP framework for efficient prior learning
- Deep Continuous Artifact-Free Fields (Nature MI 2022)  
A new PnP framework for continuous image recovery

# ODER is a PnP method for learning end-to-end optimal image priors for large-scale problems



Scan the QR code  
to see the paper

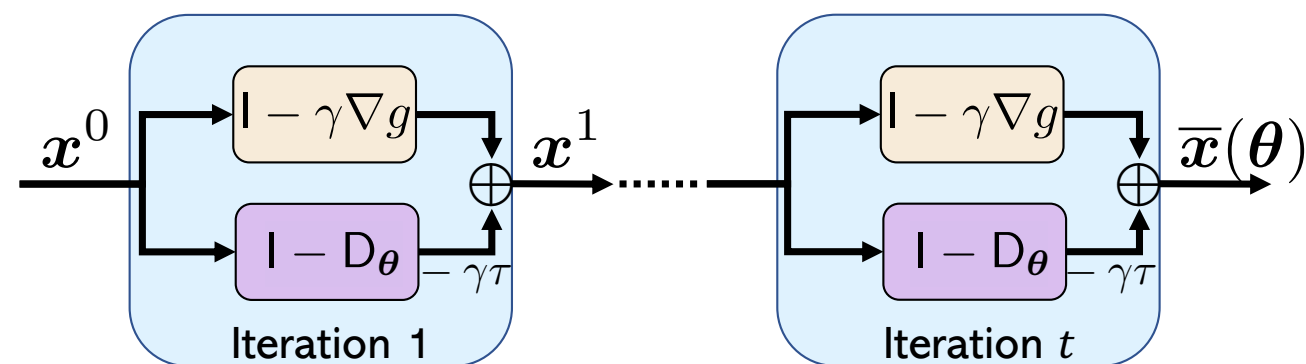


**PnP models** achieve **SOTA performance**  
when the image prior is trained at the fixed points

# PnP models achieve SOTA performance when the image prior is trained at the fixed points

PnP can be interpreted as a **model-based implicit network**

$$\mathbf{x}^t \leftarrow \mathcal{T}_{\theta}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla g(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathcal{D}_{\theta}(\mathbf{x}^{t-1})) \right)$$



Two types of operations:

physical forward model

learned CNN prior

# PnP models achieve SOTA performance when the image prior is trained at the fixed points

PnP can be interpreted as a model-based implicit network

$$\mathbf{x}^t \leftarrow \mathsf{T}_{\theta}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla g(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathsf{D}_{\theta}(\mathbf{x}^{t-1})) \right)$$

Performance of PnP models is significantly improved by  
going beyond AWGN denoisers to **artifact removing (AR) priors**

$\mathsf{D}_{\theta}$  : more image artifacts  $\mapsto$  less image artifacts

**Remark:** AWGN denoisers are suboptimal at intermediate PnP layers!

# PnP models achieve SOTA performance when the image prior is trained at the fixed points

PnP can be interpreted as a model-based implicit network

$$\mathbf{x}^t \leftarrow \mathsf{T}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla g(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathsf{D}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1})) \right)$$

Performance of PnP models is significantly improved by going beyond AWGN denoisers to artifact removing (AR) priors

$\mathsf{D}_{\boldsymbol{\theta}}$  : more image artifacts  $\mapsto$  less image artifacts

**Deep equilibrium models (DEQ)** enable efficient training of AR operators at PnP fixed points using implicit differentiation

$$\ell(\boldsymbol{\theta}) = \frac{1}{2} \|\bar{\mathbf{x}}(\boldsymbol{\theta}) - \mathbf{x}^*\|_2^2 \quad \bar{\mathbf{x}}(\boldsymbol{\theta}) \in \text{Fix}(\mathsf{T}_{\boldsymbol{\theta}})$$

$$\Rightarrow \nabla \ell(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\mathbf{x}}))^{\top} (\mathbf{I} - \nabla_{\mathbf{x}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\mathbf{x}}))^{-\top} (\bar{\mathbf{x}} - \mathbf{x}^*)$$

# PnP models achieve SOTA performance when the image prior is trained at the fixed points

PnP can be interpreted as a model-based implicit network

$$\mathbf{x}^t \leftarrow \mathsf{T}_{\theta}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla g(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathsf{D}_{\theta}(\mathbf{x}^{t-1})) \right) \quad \text{forward pass}$$

Performance of PnP models is significantly improved by going beyond AWGN denoisers to artifact removing (AR) priors

$\mathsf{D}_{\theta}$  : more image artifacts  $\mapsto$  less image artifacts

**Deep equilibrium models (DEQ)** enable efficient training of AR operators at PnP fixed points using implicit differentiation

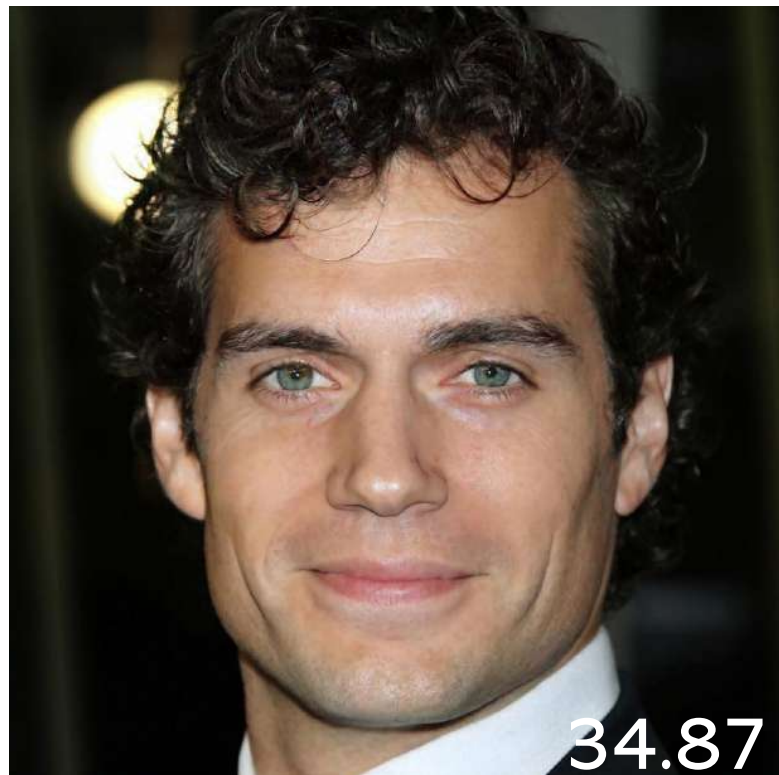
$$\ell(\theta) = \frac{1}{2} \|\bar{\mathbf{x}}(\theta) - \mathbf{x}^*\|_2^2 \quad \bar{\mathbf{x}}(\theta) \in \text{Fix}(\mathsf{T}_{\theta})$$

$$\Rightarrow \nabla \ell(\theta) = (\nabla_{\theta} \mathsf{T}_{\theta}(\bar{\mathbf{x}}))^{\top} (\mathbf{I} - \nabla_{\mathbf{x}} \mathsf{T}_{\theta}(\bar{\mathbf{x}}))^{-\top} (\bar{\mathbf{x}} - \mathbf{x}^*)$$

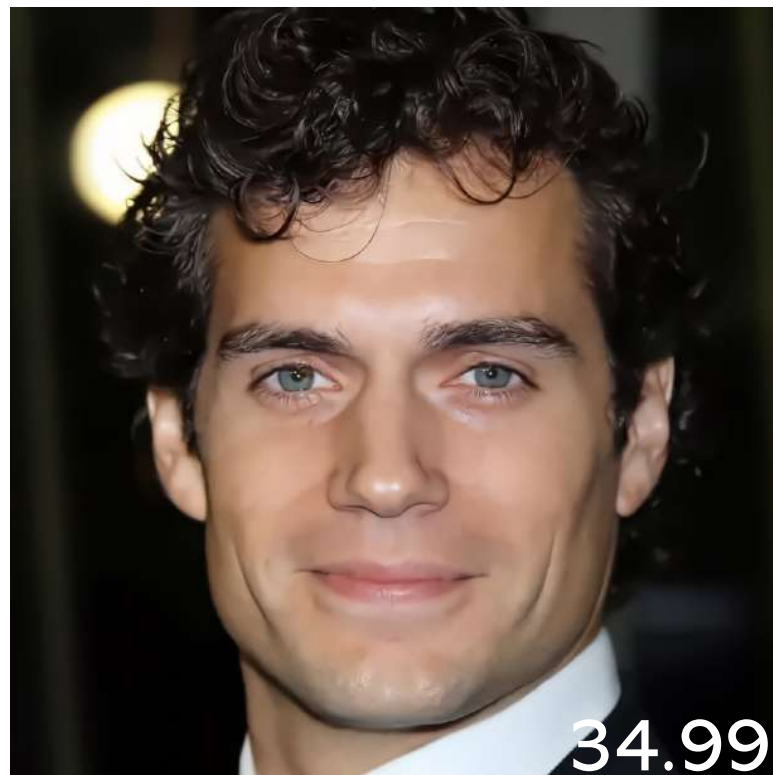
$$\Rightarrow \mathbf{b}^t \leftarrow \mathsf{F}(\mathbf{b}^{t-1}) = (\nabla_{\mathbf{x}} \mathsf{T}_{\theta}(\bar{\mathbf{x}}))^{\top} \mathbf{b}^{t-1} + (\bar{\mathbf{x}} - \mathbf{x}^*) \quad \text{backward pass}$$

# PnP models achieve SOTA performance when the image prior is trained at the fixed points

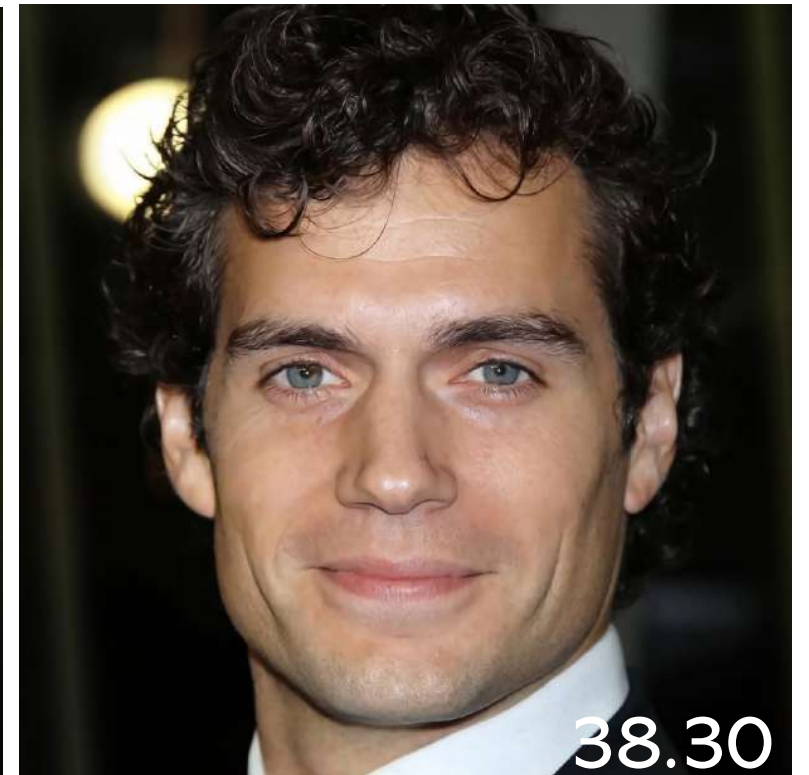
ILO (ICML 2021)



PnP (AWGN)



PnP (AR)



**Example:** In compressive sensing from random projections with 10% subsampling, PnP (AR) significantly outperforms PnP (AWGN)!



# PnP models achieve SOTA performance when the image prior is trained at the fixed points

ILO (ICML 2021)



PnP (AWGN)



PnP (AR)



**Example:** In compressive sensing from random projections with 10% subsampling, PnP (AR) significantly outperforms PnP (AWGN)!

# PnP models achieve SOTA performance when the image prior is trained at the fixed points

Table 3: Average PSNR (dB) values for several algorithms on test images from CelebA HQ.

Method \ CS Ratio						
	10%	20%	30%	40%	50%	
TV	32.13	35.24	37.41	39.35	41.29	
PULSE [36]	27.45	29.98	33.06	34.25	34.77	
ILO [37]	36.15	40.98	43.46	47.89	48.21	
PnP (denoising)	35.61	41.51	45.71	48.05	52.24	
PnP (AR)	<b>39.19</b>	<b>44.20</b>	<b>48.66</b>	<b>51.32</b>	<b>53.89</b>	

**Example:** In compressive sensing from random projections with 10% subsampling, PnP (AR) significantly outperforms PnP (AWGN)!

Training **deep equilibrium models** becomes expensive as the **number of measurements grows**

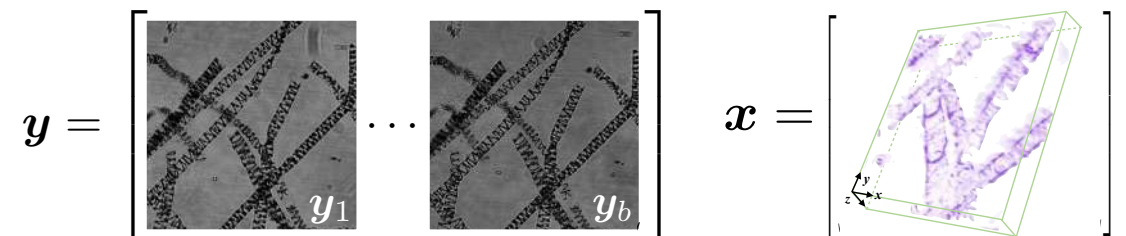
# Training deep equilibrium models becomes expensive as the number of measurements grows

Consider **data-fidelity terms** that can be expressed as

$$g(\mathbf{x}) = \frac{1}{b} \sum_{i=1}^b g_i(\mathbf{x})$$

each term  $g_i$  depends  
on a subset of measurements  $\mathbf{y}_i$

Example in IDT:  $g_i(\mathbf{x}) = \frac{1}{2} \|\mathbf{y}_i - \mathbf{A}_i(\mathbf{x})\|_2^2$



# Training deep equilibrium models becomes expensive as the number of measurements grows

Consider data-fidelity terms that can be expressed as

$$g(\mathbf{x}) = \frac{1}{b} \sum_{i=1}^b g_i(\mathbf{x})$$

**DEQ gradient** becomes impractical to compute when the number of measurements grows

$$\nabla \ell(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\mathbf{x}}))^{\top} \bar{\mathbf{b}}$$

$$\mathbf{x}^t \leftarrow \mathsf{T}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla g(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathsf{D}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1})) \right)$$

$$\mathbf{b}^t \leftarrow \mathsf{F}(\mathbf{b}^{t-1}) = (\nabla_{\mathbf{x}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\mathbf{x}}))^{\top} \mathbf{b}^{t-1} + (\bar{\mathbf{x}} - \mathbf{x}^*)$$

impractical per-iteration computational and memory complexity when  $b \rightarrow \infty$

# Training deep equilibrium models becomes expensive as the number of measurements grows

Consider data-fidelity terms that can be expressed as

$$g(\mathbf{x}) = \frac{1}{b} \sum_{i=1}^b g_i(\mathbf{x})$$

DEQ gradient becomes impractical to compute when the number of measurements grows

$$\nabla \ell(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\mathbf{x}}))^{\top} \bar{\mathbf{b}}$$

$$\mathbf{x}^t \leftarrow \mathsf{T}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla g(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathsf{R}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1})) \right)$$

$$\mathbf{b}^t \leftarrow \mathsf{F}(\mathbf{b}^{t-1}) = (\nabla_{\mathbf{x}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\mathbf{x}}))^{\top} \mathbf{b}^{t-1} + (\bar{\mathbf{x}} - \mathbf{x}^*)$$

**Question:** Can we train AR priors for PnP with complexity independent of  $b$ ?



**Stochastic PnP methods** are useful for  
reducing the per-iteration cost **during inference**

# Stochastic PnP methods are useful for reducing the per-iteration cost during inference

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x})$$

- $w \ll b$  is the mini-batch size
- complexity independent of  $b$

# Stochastic PnP methods are useful for reducing the per-iteration cost during inference

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x})$$

$$\mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x})$$

if  $\{i_1, \dots, i_w\}$  are iid  
uniform random variables

# Stochastic PnP methods are useful for reducing the per-iteration cost during inference

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x}) \quad \mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x}) \quad \text{if } \{i_1, \dots, i_w\} \text{ are iid uniform random variables}$$

**SIMBA** is a type of PnP-SGD algorithm

$$\mathbf{x}^t \leftarrow \hat{\mathbf{T}}_{\theta}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla \hat{g}(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathbf{D}_{\theta}(\mathbf{x}^{t-1})) \right)$$

**Remark:** Complexity of the SIMBA iteration is independent of  $b$ !

# Stochastic PnP methods are useful for reducing the per-iteration cost during inference

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x}) \quad \mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x}) \quad \text{if } \{i_1, \dots, i_w\} \text{ are iid uniform random variables}$$

**SIMBA** is a type of PnP-SGD algorithm

$$\mathbf{x}^t \leftarrow \hat{T}_{\theta}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla \hat{g}(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - D_{\theta}(\mathbf{x}^{t-1})) \right)$$

A theoretical error bound on the convergence of SIMBA

$$\mathbb{E} [\|\mathbf{x}^k - \bar{\mathbf{x}}\|_2^2] \leq \eta^k R + \frac{\gamma \nu}{(1 - \eta) \sqrt{w}}$$

- $0 < \eta < 1$
- $\bar{\mathbf{x}} \in \text{Fix}(\mathbf{T})$

# Stochastic PnP methods are useful for reducing the per-iteration cost during inference

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x}) \quad \mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x}) \quad \text{if } \{i_1, \dots, i_w\} \text{ are iid uniform random variables}$$

**SIMBA** is a type of PnP-SGD algorithm

$$\mathbf{x}^t \leftarrow \hat{T}_{\theta}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma \left( \nabla \hat{g}(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - D_{\theta}(\mathbf{x}^{t-1})) \right)$$

A theoretical error bound on the convergence of SIMBA

$$\mathbb{E} [\|\mathbf{x}^k - \bar{\mathbf{x}}\|_2^2] \leq \eta^k R + \frac{\gamma \nu}{(1 - \eta) \sqrt{w}}$$

**Remark:** SIMBA can in expectation approximate the fixed points of the full PnP algorithm up to an error term that depends on  $\gamma$  and  $w$ !



**ODER** extends stochastic PnP for **efficient training** of AR priors with theoretical guarantees

# ODER extends stochastic PnP for efficient training of AR priors with theoretical guarantees

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\boldsymbol{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\boldsymbol{x})$$

$$\mathbb{E} [\hat{g}(\boldsymbol{x})] = g(\boldsymbol{x})$$

# ODER extends stochastic PnP for efficient training of AR priors with theoretical guarantees

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x}) \quad \mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x})$$

**ODER** seeks to approximate the implicit gradient as

$$\nabla \hat{\ell}(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^{\top} \mathbf{b}^T \quad T \geq 1 \text{ forward and backward iterations}$$

$$\mathbf{x}^t \leftarrow \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma (\nabla \hat{g}(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathbf{R}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1})))$$

$$\mathbf{b}^t \leftarrow \hat{\mathbf{F}}(\mathbf{b}^{t-1}) = (\nabla_{\mathbf{x}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^{\top} \mathbf{b}^{t-1} + (\mathbf{x}^T - \mathbf{x}^*)$$

# ODER extends stochastic PnP for efficient training of AR priors with theoretical guarantees

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x}) \quad \mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x})$$

ODER seeks to approximate the implicit gradient as

$$\nabla \hat{\ell}(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^T \mathbf{b}^T$$

$$\mathbf{x}^t \leftarrow \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma (\nabla \hat{g}(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathbf{R}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1})))$$

$$\mathbf{b}^t \leftarrow \hat{\mathbf{F}}(\mathbf{b}^{t-1}) = (\nabla_{\mathbf{x}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^T \mathbf{b}^{t-1} + (\mathbf{x}^T - \mathbf{x}^*)$$

**Remark:** Complexity of ODER is independent of  $b$ !

# ODER extends stochastic PnP for efficient training of AR priors with theoretical guarantees

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x}) \quad \mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x})$$

ODER seeks to approximate the implicit gradient as

$$\nabla \hat{\ell}(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^T \mathbf{b}^T$$

$$\mathbf{x}^t \leftarrow \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma (\nabla \hat{g}(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathbf{R}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1})))$$

$$\mathbf{b}^t \leftarrow \hat{\mathbf{F}}(\mathbf{b}^{t-1}) = (\nabla_{\mathbf{x}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^T \mathbf{b}^{t-1} + (\mathbf{x}^T - \mathbf{x}^*)$$

**A theoretical error bound on SGD learning using ODER gradients**

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} [\|\nabla \ell(\boldsymbol{\theta}^k)\|_2^2] \leq \frac{2(\ell(\boldsymbol{\theta}^0) - \ell(\boldsymbol{\theta}^*))}{\beta K} + \frac{C_2}{\sqrt{w}} + \beta C_2$$

- $\beta > 0$  learning rate
- $K \geq 1$  SGD iterations

# ODER is a framework for efficiently training AR priors for PnP with theoretical guarantees

Consider a mini-batch approximation of the data-fidelity term

$$\hat{g}(\mathbf{x}) = \frac{1}{w} \sum_{s=1}^w g_{i_s}(\mathbf{x}) \quad \mathbb{E} [\hat{g}(\mathbf{x})] = g(\mathbf{x})$$

ODER seeks to approximate the implicit gradient as

$$\nabla \hat{\ell}(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^T \mathbf{b}^T$$

$$\mathbf{x}^t \leftarrow \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1}) = \mathbf{x}^{t-1} - \gamma (\nabla \hat{g}(\mathbf{x}^{t-1}) + \tau(\mathbf{x}^{t-1} - \mathbf{R}_{\boldsymbol{\theta}}(\mathbf{x}^{t-1})))$$

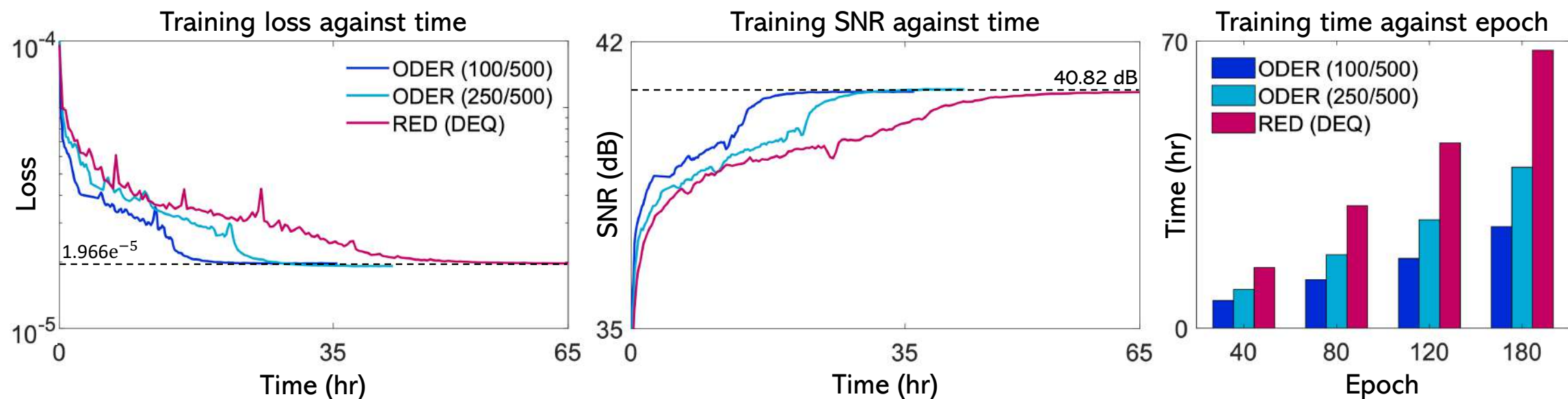
$$\mathbf{b}^t \leftarrow \hat{\mathbf{F}}(\mathbf{b}^{t-1}) = (\nabla_{\mathbf{x}} \hat{\mathbf{T}}_{\boldsymbol{\theta}}(\mathbf{x}^T))^T \mathbf{b}^{t-1} + (\mathbf{x}^T - \mathbf{x}^*)$$

**A theoretical error bound** on SGD learning using ODER gradients

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} [\|\nabla \ell(\boldsymbol{\theta}^k)\|_2^2] \leq \frac{2(\ell(\boldsymbol{\theta}^0) - \ell(\boldsymbol{\theta}^*))}{\beta K} + \frac{C_2}{\sqrt{w}} + \beta C_2$$

ODER can approximate the DEQ learning with controllable accuracy!

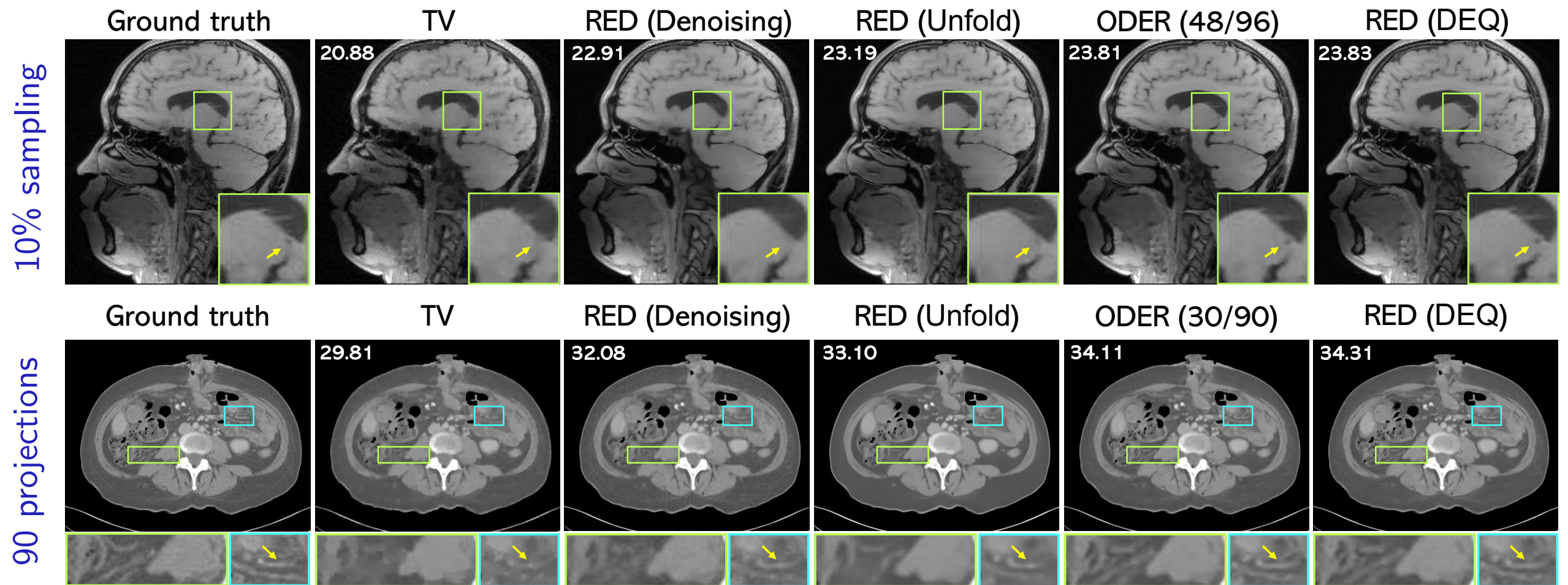
# ODER achieves nearly 2.5x improvement in training time on IDT for the same image quality



**Remark:**  $2.5\times$  Faster training in IDT for a similar loss and SNR!



# ODER achieves SOTA reconstruction performance in parallel MRI and sparse-view CT



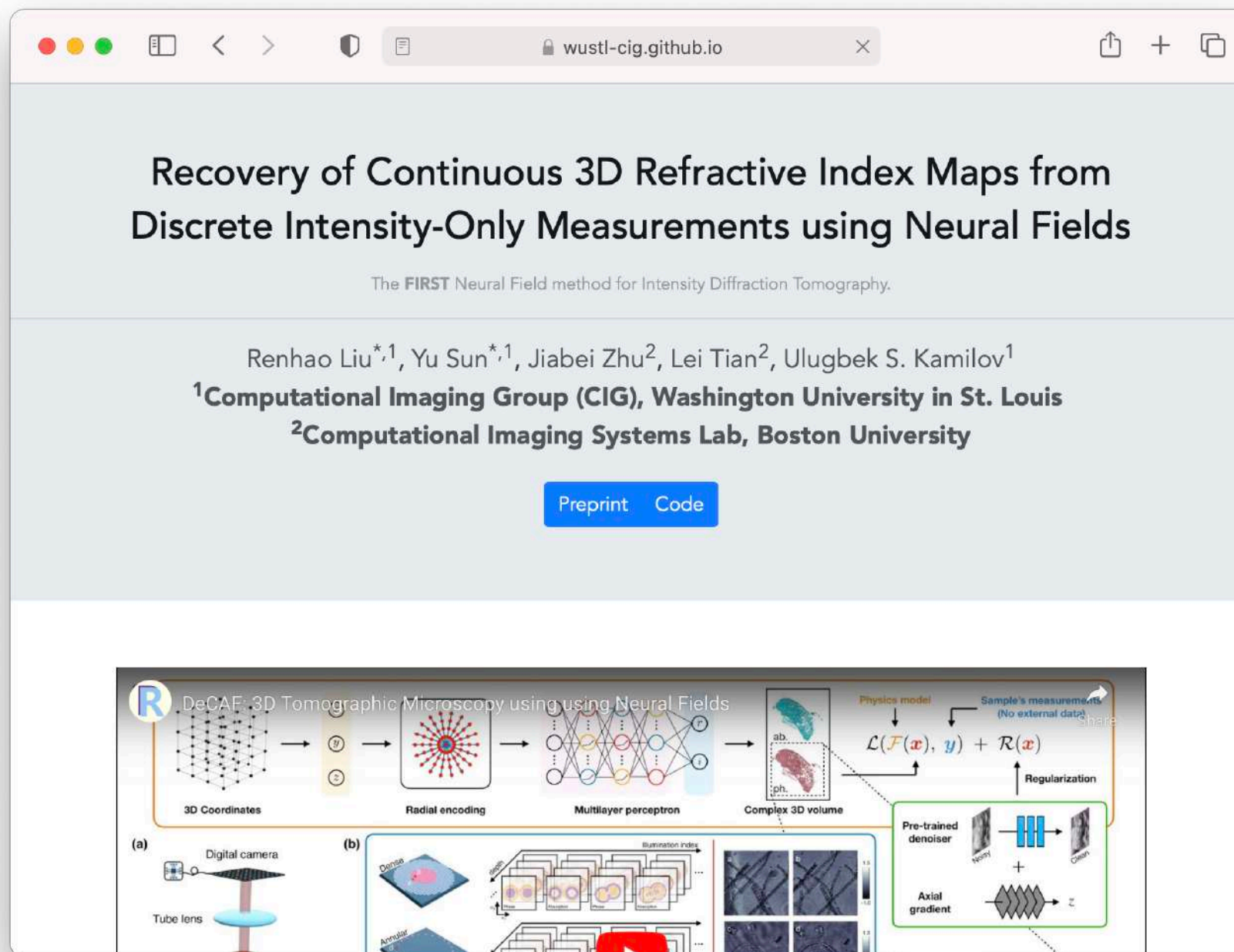
Note the similar performance of ODER and RED (DEQ), and the improvement over RED (Denoising) and RED (Unfold) due to DEQ learning!

# Outline for the rest of the talk

- ◉ Plug-and-Play Methods for Inverse Problems (IEEE SPM 2022)  
Integrating physical models and learned deep priors
- ◉ Online Deep Equilibrium Learning (NeurIPS 2022)  
A new PnP framework for efficient prior learning
- ◉ Deep Continuous Artifact-Free Fields (**Nature MI 2022**)  
A new PnP framework for continuous image recovery

**DeCAF** is a PnP method that enables the recovery of **continuous images** represented by neural fields

# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields



Scan the QR code  
to see the project

**DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields**

# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

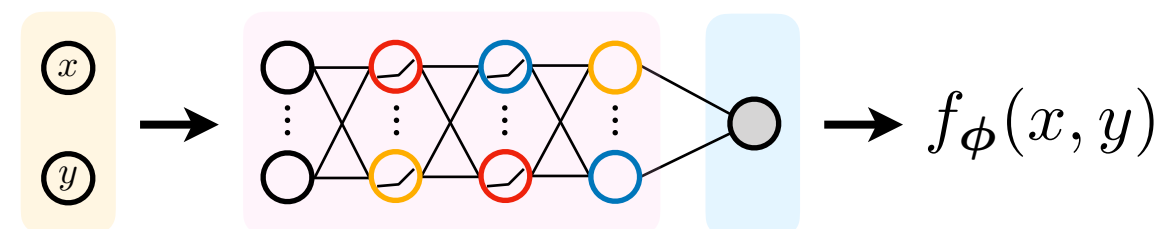
- Traditional PnP formulations seek to reconstruct images represented using a **pre-defined pixel grid**

$$\mathbf{x} = (x_1, \dots, x_n)$$



# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

- Traditional PnP formulations seek to reconstruct images represented using a pre-defined pixel grid
- DeCAF is a variant of PnP that continuously represents the desired image by using a **coordinate-based neural network**





# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

- Traditional PnP formulations seek to reconstruct images represented using a pre-defined pixel grid
- DeCAF is a variant of PnP that continuously represents the desired image by using a coordinate-based neural network
- Continuous representation in DeCAF **decouples** the representation of the solution from any pre-defined pixel grid

# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

- Traditional PnP formulations seek to reconstruct images represented using a pre-defined pixel grid
- DeCAF is a variant of PnP that continuously represents the desired image by using a coordinate-based neural network
- Continuous representation in DeCAF decouples the representation of the solution from any pre-defined pixel grid
- We tested DeCAF on a large-scale 3D IDT problem, where it is difficult to have ground-truth training data for the refractive index

$$y = \left[ \begin{array}{c} \text{Image } y_1 \\ \dots \\ \text{Image } y_b \end{array} \right] \quad x = \left[ \begin{array}{c} \text{3D Volume } x \end{array} \right]$$

# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

**DeCAF** maps **coordinates** to **image values** using a deep network

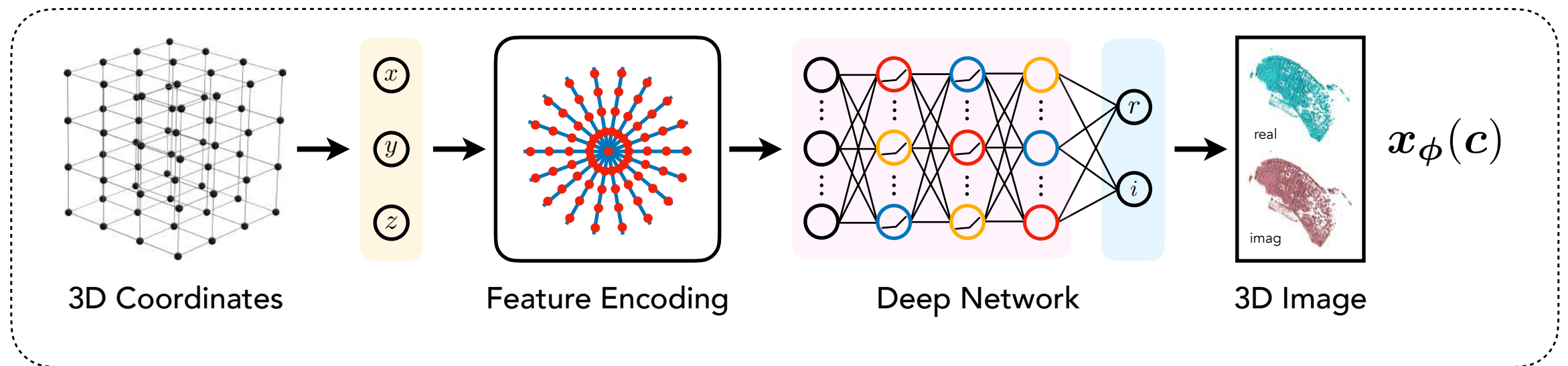
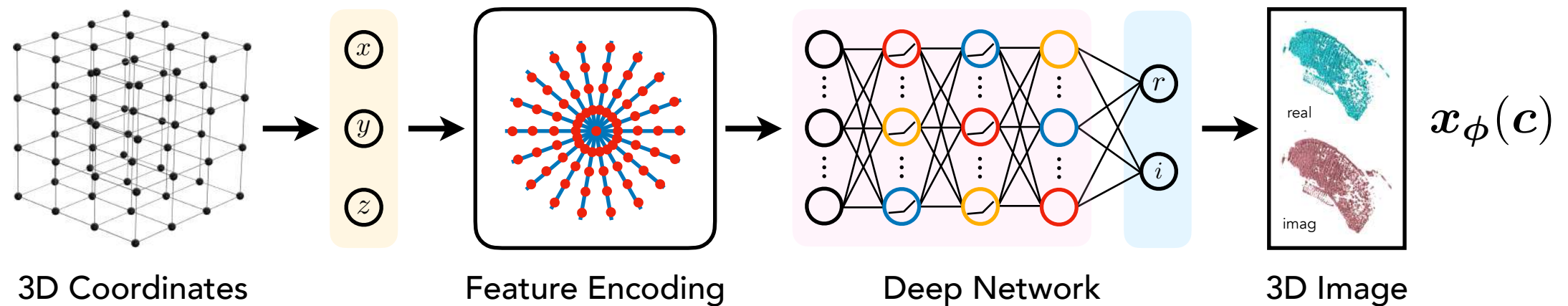


Image is represented as  $x_\phi(c)$  where  $\phi$  are weights and  $c = \{(x_i, y_i, z_i)\}$  are coordinates

# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

DeCAF maps coordinates to image values using a deep network



**Image formation** is formulated as a stochastic PnP algorithm

$$\phi^* = \arg \min_{\phi} \left\{ \frac{1}{b} \sum_{i=1}^b g_i(\mathbf{x}_\phi) + h(\mathbf{x}_\phi) \right\}$$

Minimization over network weights; not image pixel values

Quadratic data-fidelity

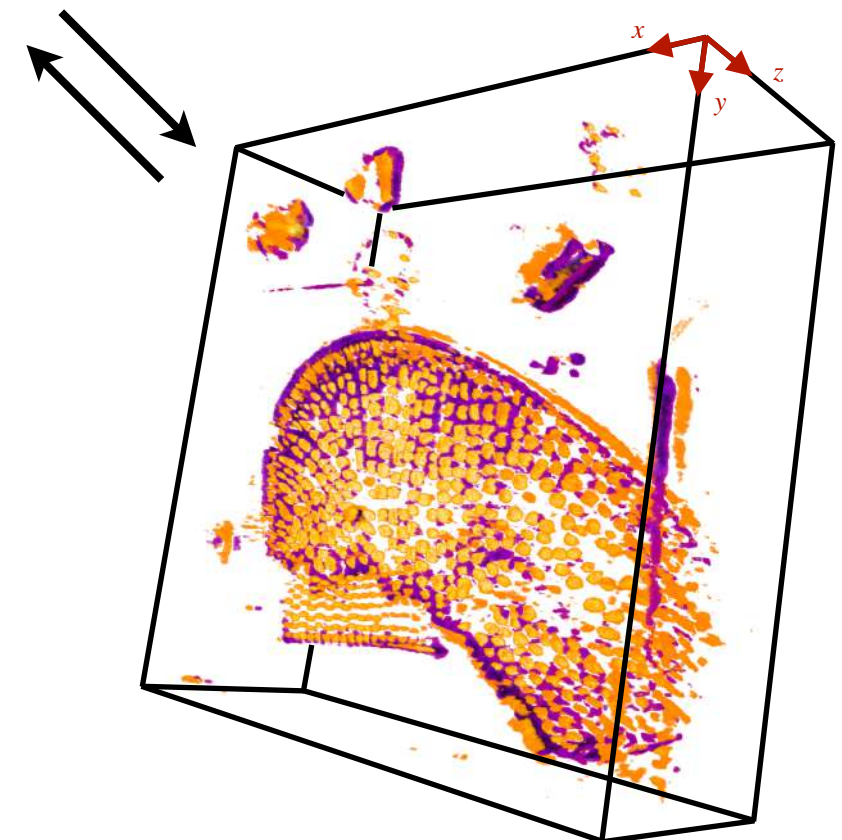
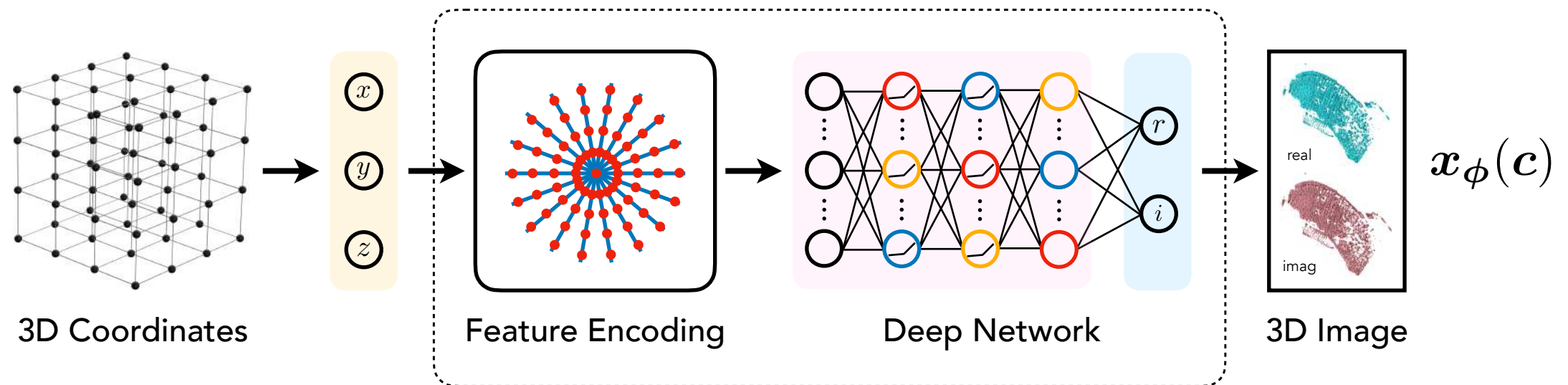
$$g(\mathbf{x}_\phi) = \frac{1}{2} \|\mathbf{y}_i - \mathbf{A}_i(\mathbf{x}_\phi)\|_2^2$$

Pre-trained deep prior

$$h(\mathbf{x}_\phi) = \frac{\tau}{2} \|\mathbf{x}_\phi - \mathbf{D}_\theta(\mathbf{x}_\phi)\|_2^2$$

# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

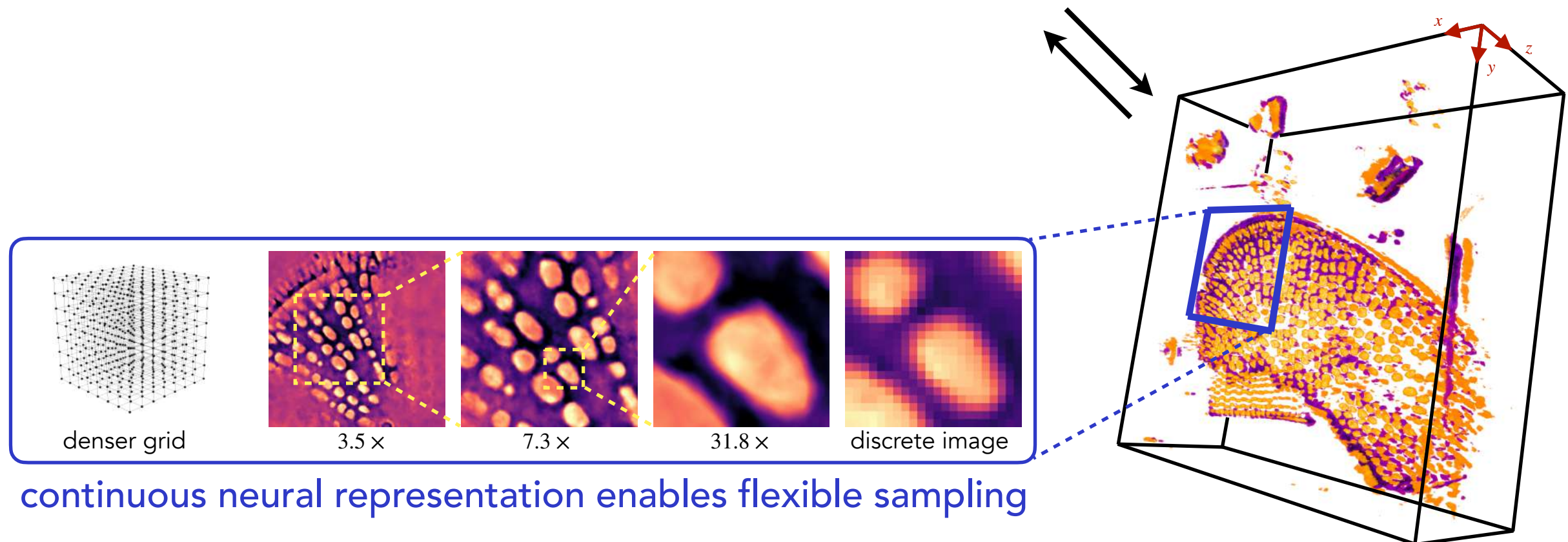
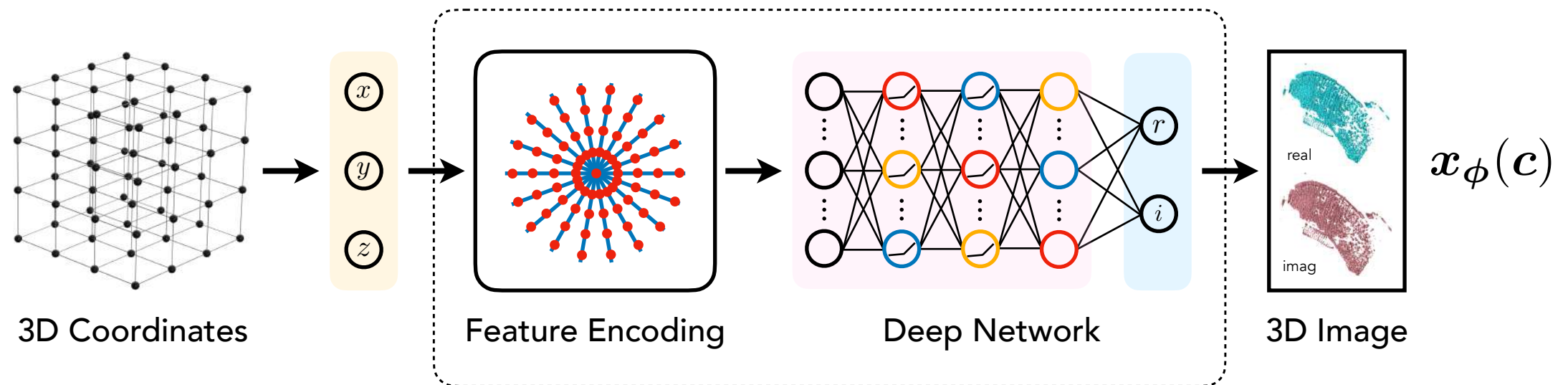
DeCAF maps coordinates to image values using a deep network





# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

DeCAF maps coordinates to image values using a deep network



# DeCAF is a PnP method that enables the recovery of continuous images represented by neural fields

DeCAF maps coordinates to image values using a deep network

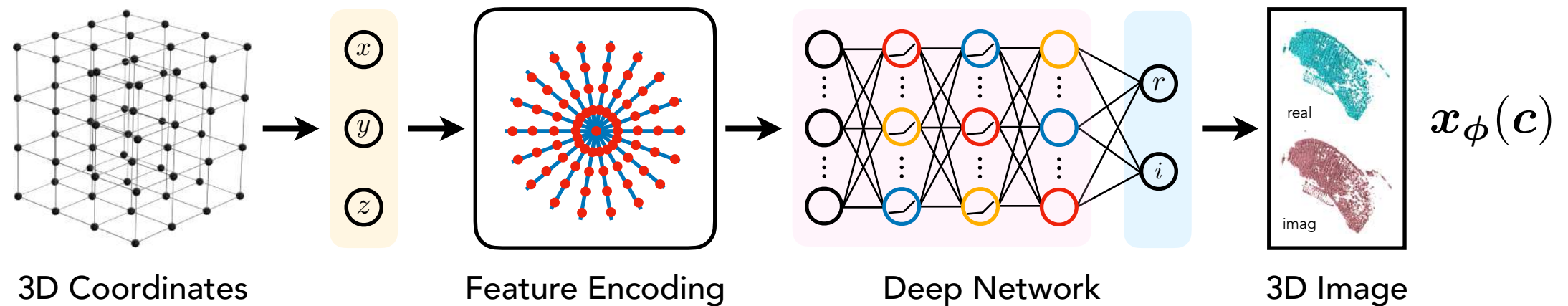


Image formation is performed by integrating DeCAF into PnP

$$\phi^* = \arg \min_{\phi} \{g(\mathbf{x}_\phi) + h(\mathbf{x}_\phi)\}$$

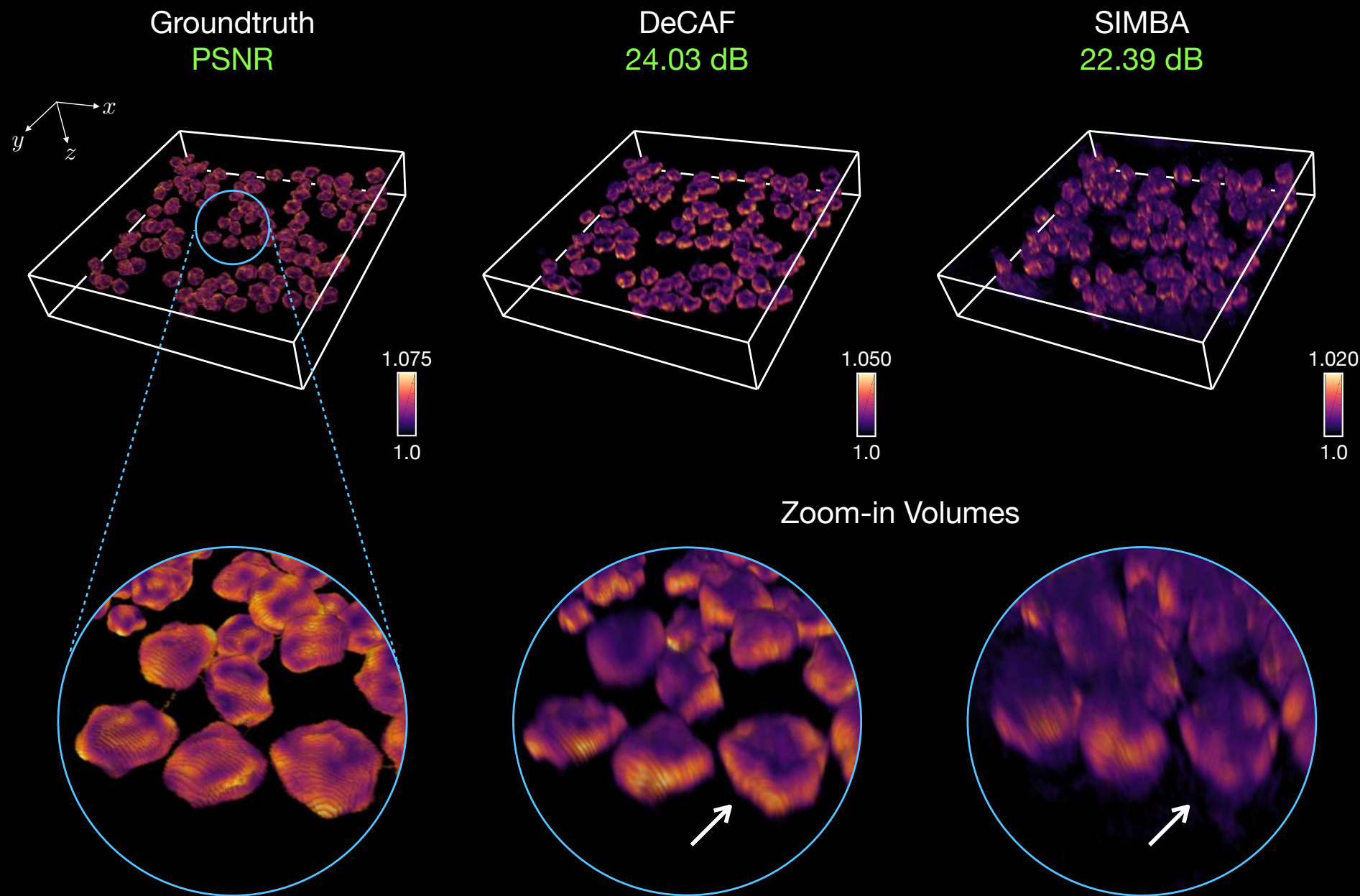
The key benefits of DeCAF include:

- (a) representation is decoupled from any pre-defined voxel grid
- (b) subimages can be synthesized at will during optimization
- (c) image prior is trained on natural images



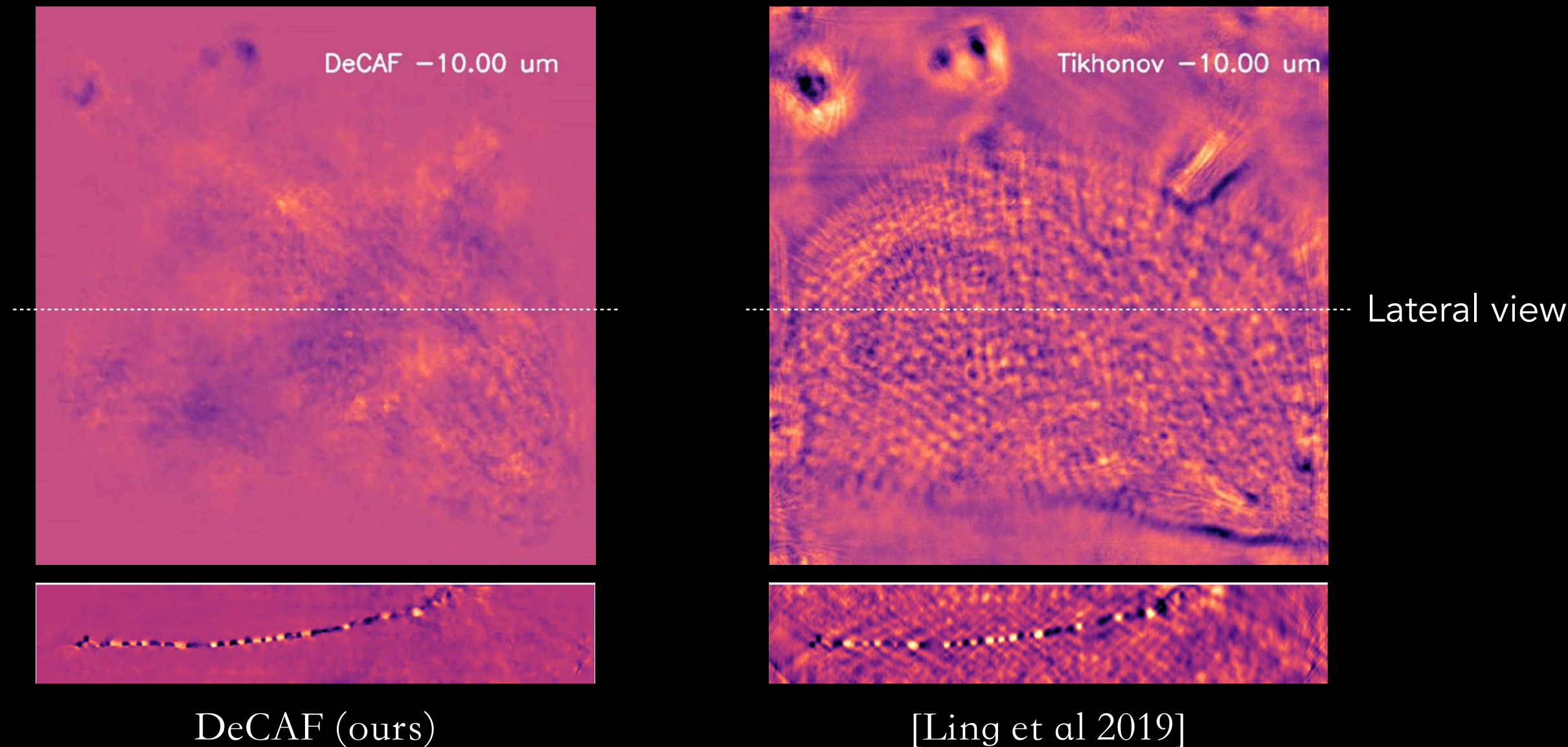
**DeCAF can recover high-quality 3D RI volumes  
from experimentally collected 2D IDT images**

# DeCAF can recover high-quality 3D RI volumes from experimentally collected 2D IDT images



Note how DeCAF outperforms SIMBA by about  $1.5\times$  dB on simulated data;  
SIMBA is a traditional (discrete) online PnP method!

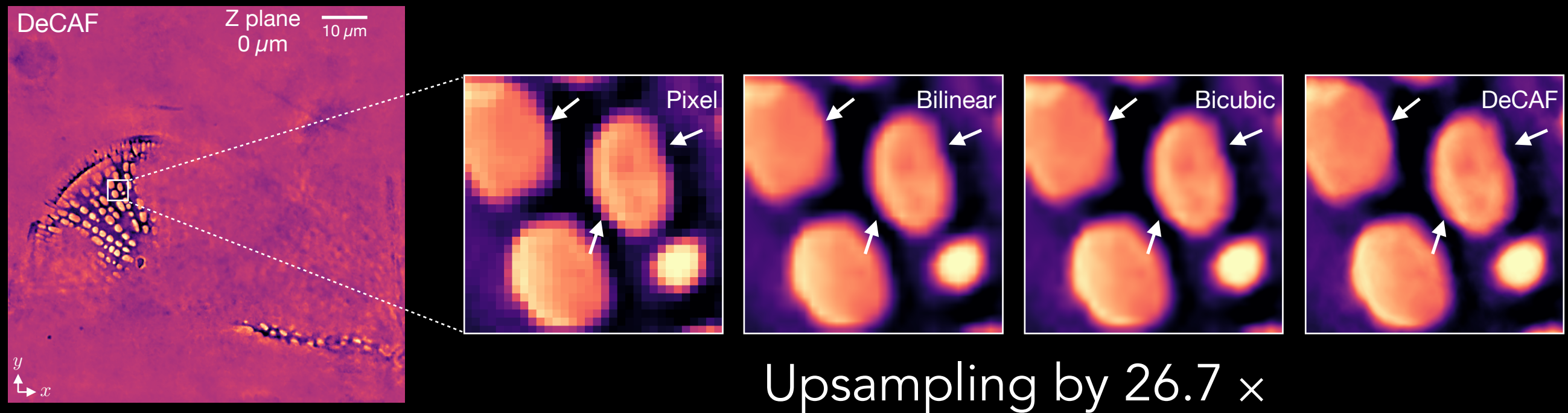
# DeCAF can recover high-quality 3D RI volumes from experimentally collected 2D IDT images



Note the optical sectioning ability of DeCAF on this Diatom Algae sample!

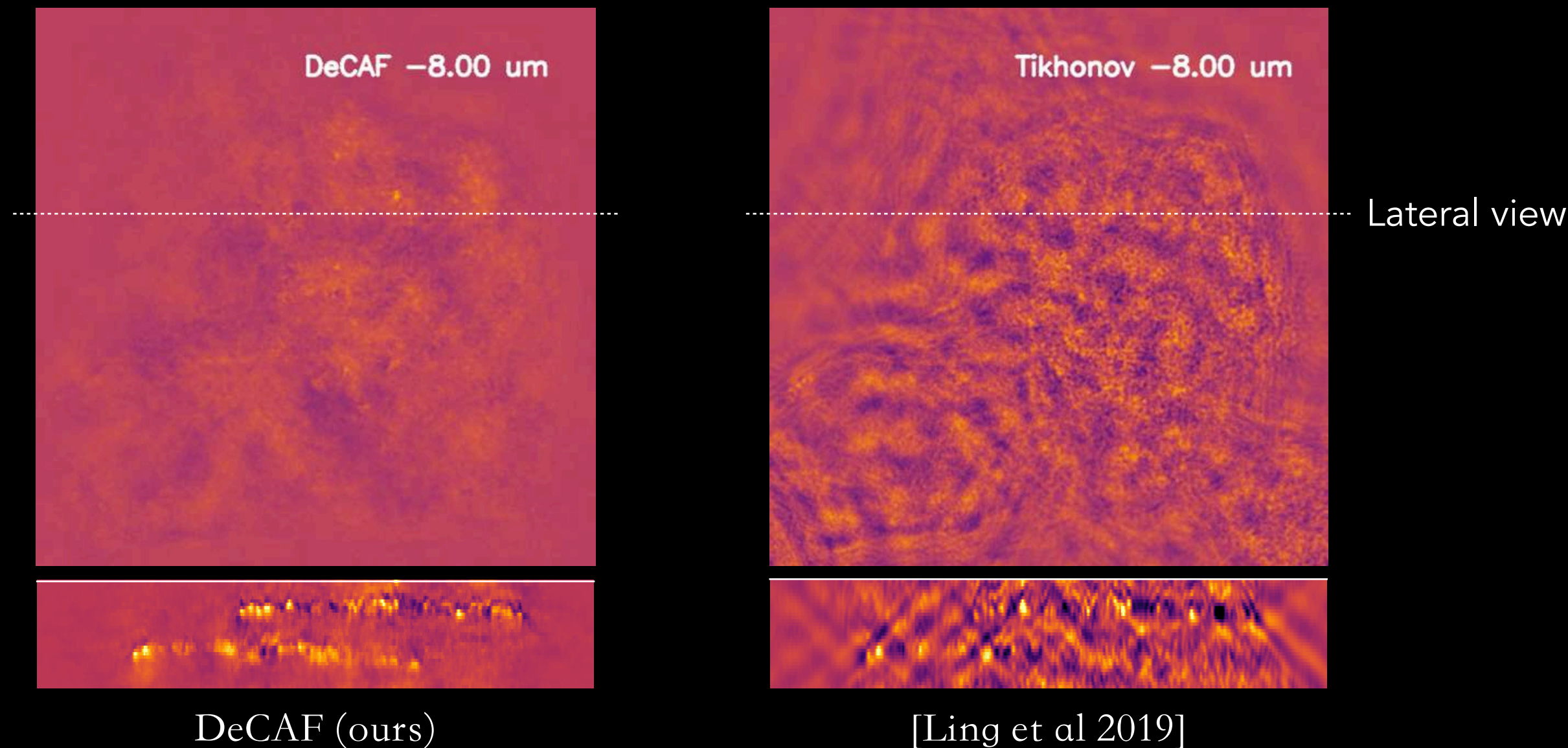


# DeCAF can recover high-quality 3D RI volumes from experimentally collected 2D IDT images



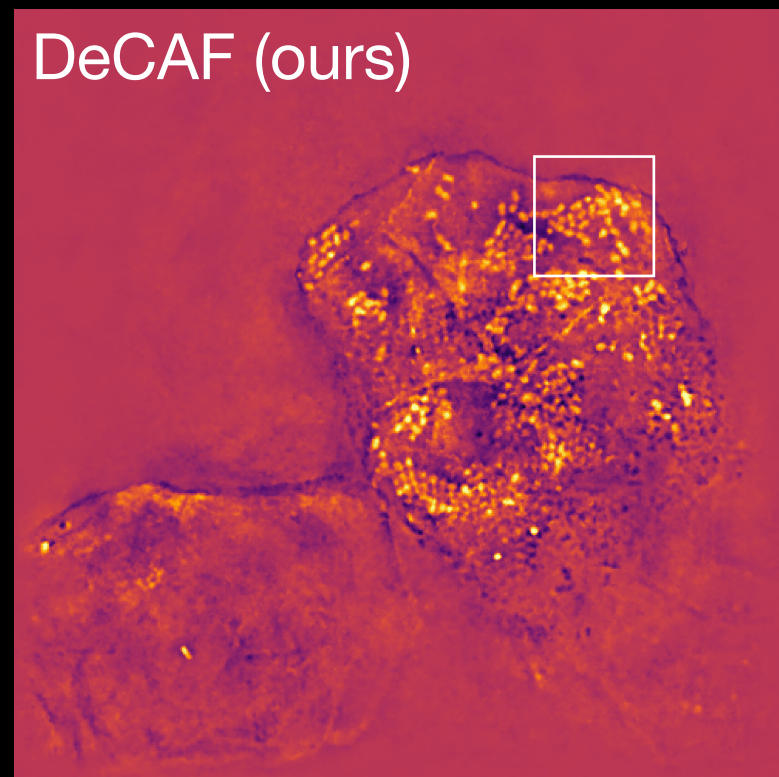
Note how DeCAF can synthesize different parts of this Diatom Algae sample on any grid!

# DeCAF can recover high-quality 3D RI volumes from experimentally collected 2D IDT images

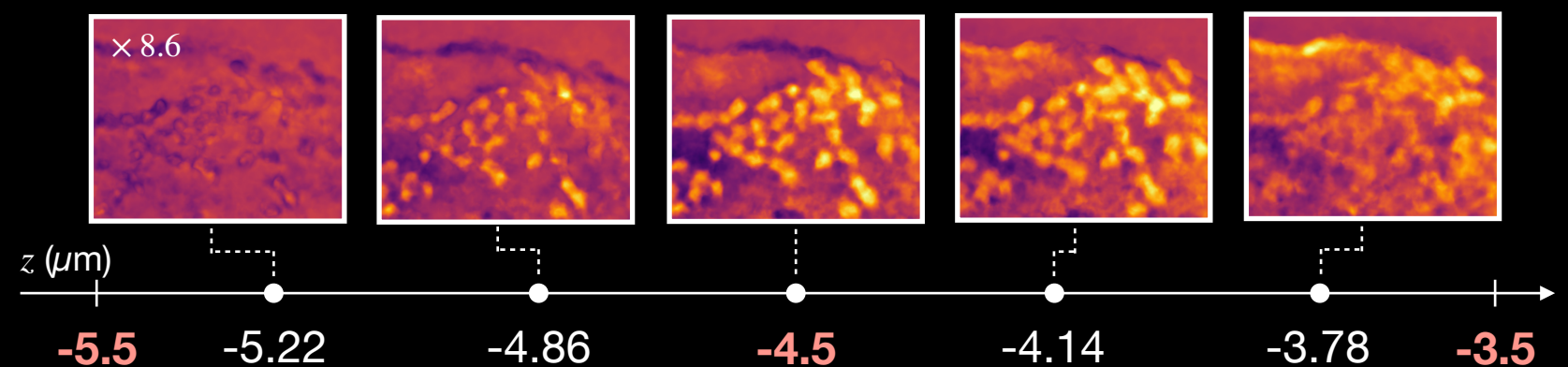


Note the optical sectioning ability of DeCAF on this Human Buccal Epithelial Cells sample!

# DeCAF can recover high-quality 3D RI volumes from experimentally collected 2D IDT images



Continuous interpolation along  $z$  dimension



Note how DeCAF can be used to synthesize any slice of this Human Buccal Epithelial Cells sample along the  $z$  dimension!

# To conclude

- **PnP** is one of the most influential class of methods across computational imaging and inverse problems
- We presented three recent extensions of the PnP models for solving large-scale inverse problems in IDT, CT, and MRI
- **ODER** is a framework for training implicit online neural nets with theoretical error bounds on the accuracy
- **DeCAF** is a PnP framework for directly reconstructing an image continuously represented by a neural field
- **DOLCE** is a conditional diffusion model for sampling realistic-looking images from few limited-angle CT data



# Big Thank You to my team at the WashU Computational Imaging Group (CIG)



My Twitter: [@ukmlv](#)

Web: [cigroup.wustl.edu](http://cigroup.wustl.edu)

Group Twitter: [@wustlcig](#)

GitHub: [github.com/wustl-cig](https://github.com/wustl-cig)

