# Plug-In SGD: Image Reconstruction in the Age of Machine Learning

Ulugbek S. Kamilov

Computational Imaging Group (CIG)
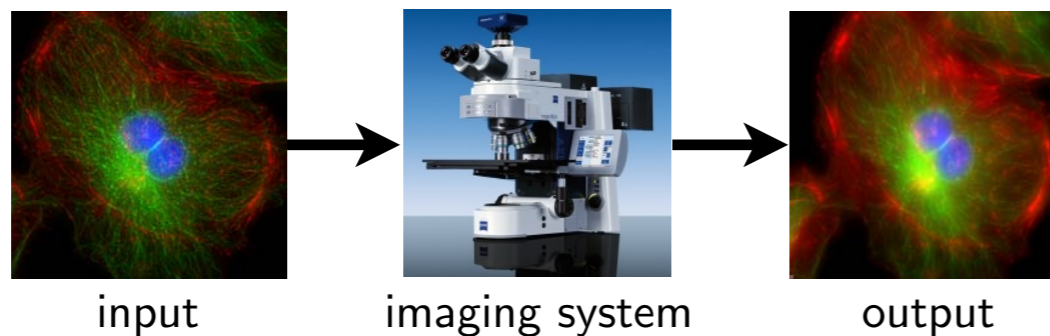Washington University, St. Louis, USA

cigroup.wustl.edu • @wustlcig • kamilov@wustl.edu

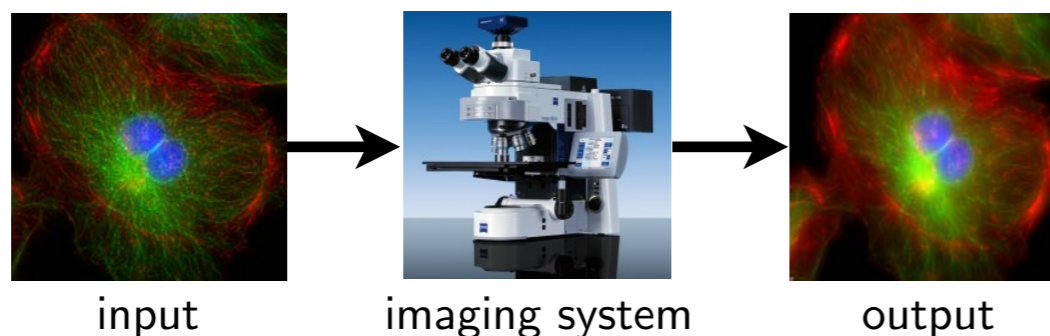# Imaging technology is going through a paradigm shift with computation at its core

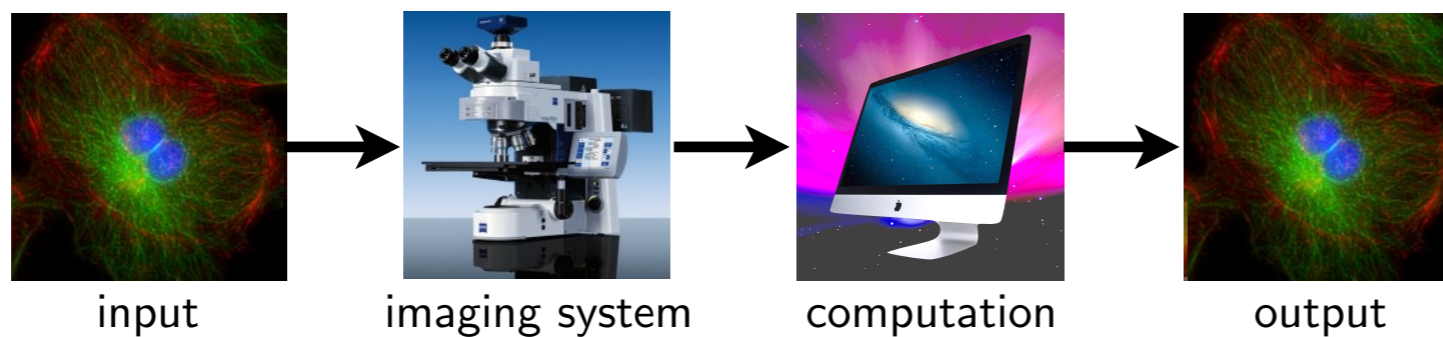# Imaging technology is going through a paradigm shift with computation at its core

**Inverse problems in bio-imaging**

**Past:** Can I see?

- Linear forward model

$$\mathbf{y} = \mathbf{Hs} + \mathbf{n}$$



s    input        imaging system        output: measurements y

Problem: Recover s from measurements y

- The easy scenario

Inverse problem is well-posed it.. $\exists$ for all $\mathbf{s}$.. $\alpha \|\mathbf{s}\| \leq \|\mathbf{Hs}\|$

$$\Rightarrow \quad \mathbf{s} \approx \mathbf{H}^{-1}\mathbf{y}$$

- Backprojection (pseudo-solution) $\mathbf{s} \approx \mathbf{H}^T\mathbf{y}$

Part 1:

Setting up
the problem

# Imaging technology is going through a paradigm shift with computation at its core

**Inverse problems in bio-imaging**

## Past: Can I see?

- Linear forward model

$$y = Hs + n$$



s      input      imaging system      output

**Inverse problems in bio-imaging**

## Present: Can I see better?

- Linear forward model

$$y = Hs + n$$

Inverse problem is **well posed** if $\exists c_0 > 0$ s.t., for all $s$, $c_0\|s\| \leq \|Hs\|$

$\Rightarrow \quad s \approx H^{-1}y$

- Backprojection (poor man's solution): $s \approx H^T y$



input      imaging system      computation      output

Problem: recover **s** from noisy measurements **y**

- The easy scenario

Inverse problem is **well posed** if $\exists c_0 > 0$ s.t., for all $s \in \mathcal{X}$, $c_0\|s\| \leq \|Hs\|$

$\Rightarrow \quad s \approx H^{-1}y$

- Backprojection (poor man's solution): $s \approx H^T y$

**Part 1:**

# Imaging technology is going through a paradigm shift with computation at its core

**Inverse problems in bio-imaging**

## Past: Can I see?

- Linear forward model

$$y = Hs + n$$

s input imaging system output

## Present: Can I see better?

- Linear forward model $$y = Hs + n$$

Inverse problem is **well posed** if $\exists c_0 > 0$ s.t., for all $s \in \mathcal{X}$, $c_0\|s\| \leq \|Hs\|$

- Backprojection

input imaging system computation output

Problem: recover **s** from noisy measurements **y**

## Future: Can I see more?

- Linear forward model $$y = Hs + n$$

Inverse problem is **well posed** if $\exists c_0 > 0$ s.t., for all $s \in \mathcal{X}$, $c_0\|s\| \leq \|Hs\|$

(backprojection): $s \approx H^T y$

input imaging system cloud computation output

Problem: recover **s** from noisy measurements **y**

- The easy scenario

# Today we will talk about

- **Forward models in imaging**
  Relating the unknowns to the measured data

- **Notions of ill-posedness and regularization**
  When measurements are not enough

- **Optimization at large scales**
  When analytical solutions are not enough

- **Plug-and-Play Priors (PnP) at large scales**
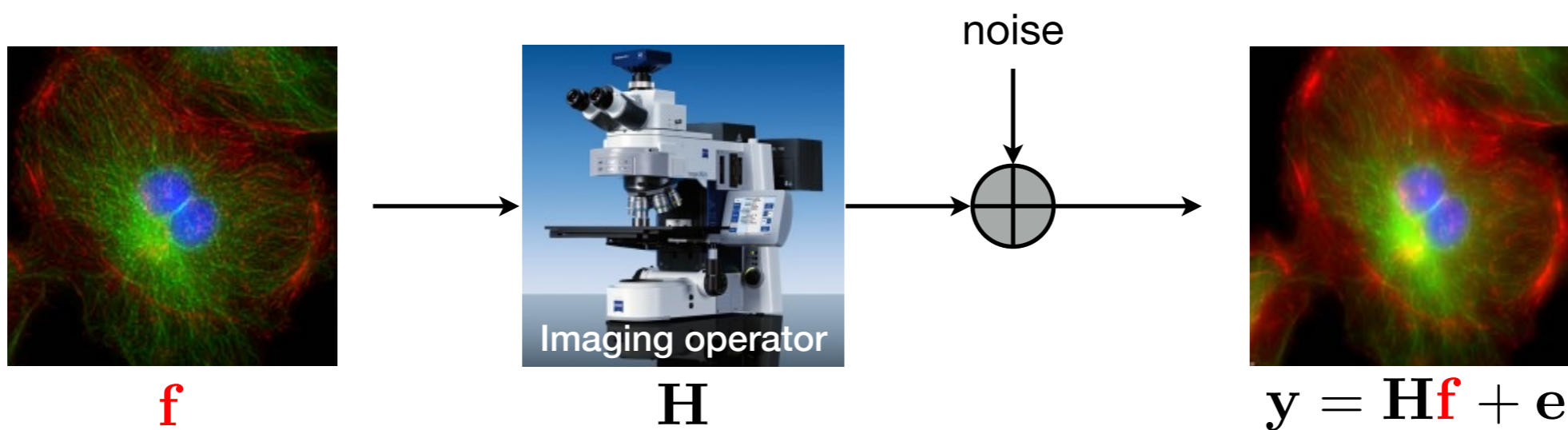  When traditional optimization is not enough

# Today we will talk about

◉ **Forward models in imaging**
  **Relating the unknowns to the measured data**

◉ Notions of ill-posedness and regularization
  When measurements are not enough

◉ Optimization at large scales
  When analytical solutions are not enough

◉ Plug-and-Play Priors (PnP) at large scales
  When traditional optimization is not enough

# Forward model relates
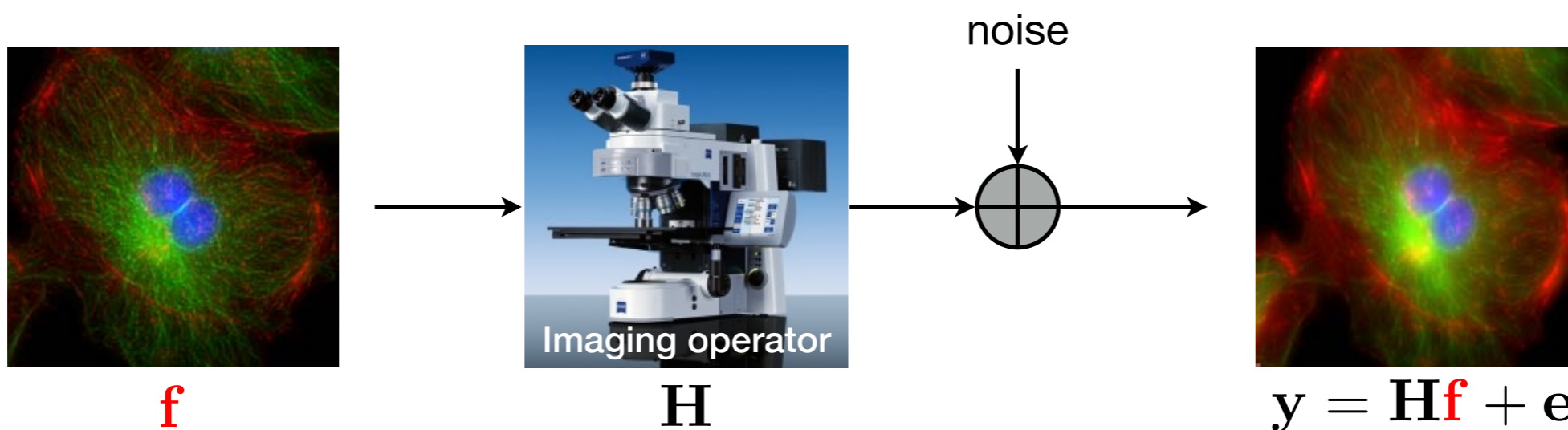# the unknown object to the observed data

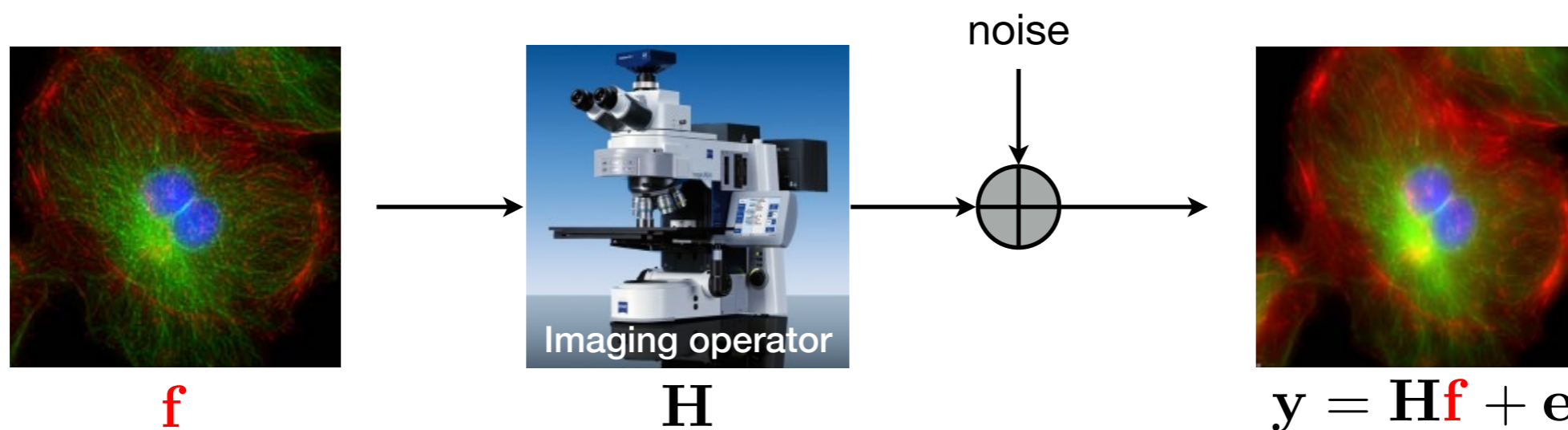# Forward model relates
# the unknown object to the observed data

**Inverse problems in bio-imaging**

- Linear forward model

$$y = Hs + n$$

noise

$$H$$

Imaging operator

$$n$$

$$f \quad s$$

Problem: recover $s$ from noisy measurements $y$

- The easy scenario

Inverse problem is **well posed** if $\exists \; c_0 > 0$ s.t. for all $s$, $c_0 \|s\| \leq \|Hs\|$

$$\Rightarrow \quad s \approx H^{-1} y$$

- Back projection (poor man's solution) $s \approx H^T y$

# Forward model relates
# the unknown object to the observed data

Forward problem: generate y from $\mathbf{f}$

**Inverse problems in imaging**

- Linear forward model

$$y = \mathbf{H}\mathbf{s} + \mathbf{n}$$



H

noise

Imaging operator

$\mathbf{n}$

$\mathbf{f}$  $\mathbf{s}$

Problem: recover $\mathbf{s}$ from noisy measurements $y$

- The easy scenario

Inverse problem is **well posed** if $\exists c_0 > 0$ s.t., for all $\mathbf{s}$, $c_0 \|\mathbf{s}\| \leq \|\mathbf{H}\mathbf{s}\|$

$$\Rightarrow \quad \mathbf{s} \approx \mathbf{H}^{-1}y$$

- Back-projection (poor man's solution) $\mathbf{s} \approx \mathbf{H}^{T}y$

# Forward model relates
# the unknown object to the observed data

**Inverse problems in imaging**

Forward problem: generate y from $\mathbf{f}$

- Linear forward model

$$y = \mathbf{H}\mathbf{f} + \mathbf{n}$$

noise

$\mathbf{H}$

Imaging operator

$\mathbf{n}$

$\mathbf{f}$

Problem: recover $\mathbf{f}$ from noisy measurements $y$

- The easy scenario

Inverse problem: recover $\mathbf{f}$ from $y$

Inverse problem is well-posed if $\|H\mathbf{f}\| \leq \|\mathbf{H}\mathbf{f}\|$

$$\Rightarrow \quad \mathbf{f} \approx \mathbf{H}^{-1}y$$

- Back-projection (poor man's solution) $\mathbf{f} \approx \mathbf{H}^{T}y$

# Forward model relates the unknown object to the observed data

**Inverse problems in imaging**

Forward problem: generate y from $\mathbf{f}$

- Linear forward model

$$y = \mathbf{H}\mathbf{f} + \mathbf{n}$$

noise

$\mathbf{H}$

Imaging operator

$\mathbf{n}$

$\mathbf{f}$    $\mathbf{s}$

Problem: recover $\mathbf{f}$ from noisy measurements $y$

$$\mathbf{y} = \mathbf{H}\mathbf{f} + \mathbf{e}$$

- The easy scenario

Inverse problem is **well-posed** if $\det \mathbf{H} \neq 0$, for all, s.t. $\|\mathbf{a}\|_0 \leq \|\mathbf{s}\|$  $\|\mathbf{H}\mathbf{s}\|$

$$\Rightarrow \quad \mathbf{s} \approx \mathbf{H}^{-1}y$$

Inverse problem: recover $\mathbf{f}$ from $\mathbf{y}$

- Back-projection (poor man's solution) $\approx \mathbf{H}^T y$

Question: Which problem is harder to solve?

# Forward models can be represented as integrals

# Forward models can be represented as integrals

# Forward models can be represented as integrals

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

**Inverse problems in imaging**

- Linear forward model $\quad y = Hs + n$



$s \qquad \qquad y = Hf$

Problem: recover $s$ from noisy measurements $y$

- The easy scenario

# Forward models can be represented as integrals

**Unknown molecular/anatomical map:** $f(\boldsymbol{r}), \quad \boldsymbol{r} = (x, y, z, t) \in \mathbb{R}^d$

defined over a continuum in space-time



Molecular Imaging
Basic Principles and Applications
in Biomedical Research
··· Markus Rudin
Imperial College Press

**Inverse problems in imaging**

- Linear forward model

$$y = Hs + n$$



Problem: recover $s$ from noisy measurements $y$

- The easy scenario

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Forward models can be represented as integrals

Unknown molecular/anatomical map: $\quad f(\boldsymbol{r}), \quad \boldsymbol{r} = (x, y, z, t) \in \mathbb{R}^d$

Space of finite-energy functions: $\quad \boxed{f \in L_2(\mathbb{R}^d)}$



Molecular Imaging
Basic Principles and Applications in Biomedical Research

··· Markus Rudin

Imperial College Press

Inverse problems in imaging

- Linear forward model $\quad \mathbf{y} = \mathbf{Hs} + \mathbf{n}$



Problem: recover $s$ from noisy measurements $y$

- The easy scenario

# Forward models can be represented as integrals

Unknown molecular/anatomical map:   $f(\boldsymbol{r}), \quad \boldsymbol{r} = (x, y, z, t) \in \mathbb{R}^d$

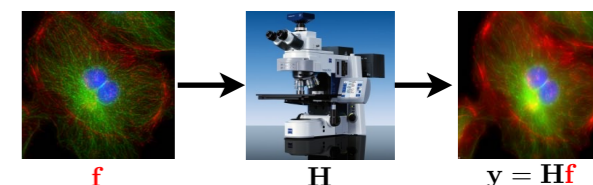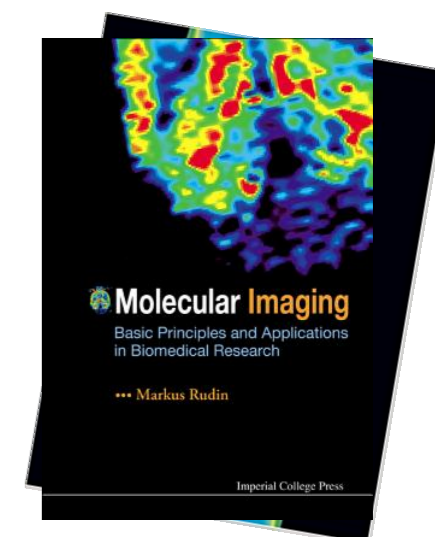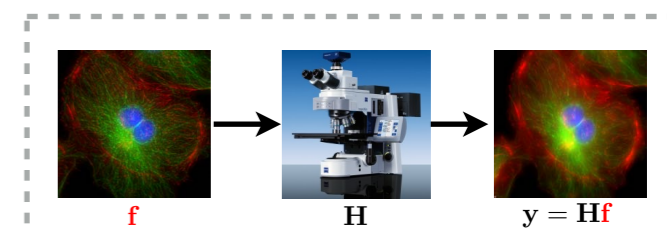Space of finite-energy functions:   $f \in L_2(\mathbb{R}^d)$

Imaging operator:   $\mathsf{H} : s \mapsto \mathbf{y} = (y_1, \ldots, y_m) = \mathsf{H}\{f\}$

from continuum to finite dimensional:   $\mathsf{H} : L_2(\mathbb{R}^d) \to \mathbb{R}^m$



Molecular Imaging
Basic Principles and Applications in Biomedical Research
••• Markus Rudin

Imperial College Press

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

Inverse problems in imaging

■ Linear forward model   $\mathbf{y} = \mathsf{H}\mathbf{s} + \mathbf{n}$



s   Problem: recover s from noisy measurements y

■ The easy scenario

# Forward models can be represented as integrals

Unknown molecular/anatomical map: $\quad f(\boldsymbol{r}), \quad \boldsymbol{r} = (x, y, z, t) \in \mathbb{R}^d$

Space of finite-energy functions: $\quad f \in L_2(\mathbb{R}^d)$
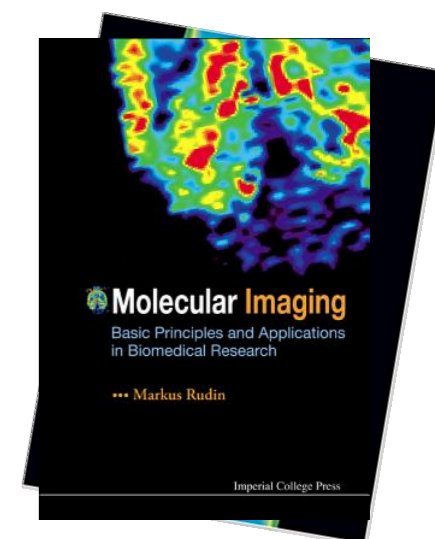
Imaging operator: $\quad \mathsf{H} : s \mapsto \mathbf{y} = (y_1, \ldots, y_m) = \mathsf{H}\{f\}$

Linearity assumption: $\quad \forall \alpha_1, \alpha_2 \in \mathbb{R}, \quad \forall f_1, f_2 \in L_2(\mathbb{R}^d)$

$$\mathsf{H}\{\alpha_1 f_1 + \alpha_2 f_2\} = \alpha_1 \mathsf{H}\{f_1\} + \alpha_2 \mathsf{H}\{f_2\}$$

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Forward models can be represented as integrals

**Unknown molecular/anatomical map:** $\quad f(\boldsymbol{r}), \quad \boldsymbol{r} = (x, y, z, t) \in \mathbb{R}^d$

**Space of finite-energy functions:** $\quad f \in L_2(\mathbb{R}^d)$

**Imaging operator:** $\quad \mathsf{H} : s \mapsto \mathbf{y} = (y_1, \ldots, y_m) = \mathsf{H}\{f\}$

**Linearity assumption:** $\quad \forall \alpha_1, \alpha_2 \in \mathbb{R}, \quad \forall f_1, f_2 \in L_2(\mathbb{R}^d)$

$$\mathsf{H}\{\alpha_1 f_1 + \alpha_2 f_2\} = \alpha_1 \mathsf{H}\{f_1\} + \alpha_2 \mathsf{H}\{f_2\}$$

$$\Rightarrow \quad [\mathbf{y}]_m = y_m = \langle h_m, f \rangle = \int_{\mathbb{R}^d} h_m(\boldsymbol{r}) f(\boldsymbol{r}) \, \mathrm{d}\boldsymbol{r}$$

by the Riesz representation theorem

**Molecular Imaging**
Basic Principles and Applications in Biomedical Research
••• Markus Rudin
Imperial College Press

Inverse problems in imaging

- Linear forward model $\qquad \mathbf{y} = \mathbf{H}\mathbf{s} + \mathbf{n}$

H

n

Problem: recover $s$ from noisy measurements $y$

- The easy scenario

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Forward models can be represented as integrals

Unknown molecular/anatomical map:  $f(\boldsymbol{r}), \quad \boldsymbol{r} = (x, y, z, t) \in \mathbb{R}^d$

Space of finite-energy functions:  $f \in L_2(\mathbb{R}^d)$

Imaging operator:  $\mathsf{H} : s \mapsto \mathbf{y} = (y_1, \ldots, y_m) = \mathsf{H}\{f\}$

Linearity assumption:  $\forall \alpha_1, \alpha_2 \in \mathbb{R}, \quad \forall f_1, f_2 \in L_2(\mathbb{R}^d)$

$$\mathsf{H}\{\alpha_1 f_1 + \alpha_2 f_2\} = \alpha_1 \mathsf{H}\{f_1\} + \alpha_2 \mathsf{H}\{f_2\}$$

impulse response of $m$th detector

$$\Rightarrow \quad [\mathbf{y}]_m = y_m = \langle h_m, f \rangle = \int_{\mathbb{R}^d} h_m(\boldsymbol{r}) f(\boldsymbol{r}) \, \mathrm{d}\boldsymbol{r}$$

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Example imaging operator:
# Fourier transform is extensively used in MRI

# Example imaging operator:
# Fourier transform is extensively used in MRI



*"Images are obviously made of sine waves..."*

# Example imaging operator:
# Fourier transform is extensively used in MRI

**Fourier transform:** $\mathcal{F} : L_2(\mathbb{R}^d) \to L_2(\mathbb{R}^d)$

$$\hat{f}(\boldsymbol{\omega}) = \mathcal{F}\{f\} = \int_{\mathbb{R}^d} f(\boldsymbol{r})\, e^{-j\langle \boldsymbol{\omega}, \boldsymbol{r} \rangle}\, d\boldsymbol{r}$$

# Example imaging operator:
# Fourier transform is extensively used in MRI

**Fourier transform:** $\mathcal{F} : L_2(\mathbb{R}^d) \to L_2(\mathbb{R}^d)$

$$\hat{f}(\boldsymbol{\omega}) = \mathcal{F}\{f\} = \int_{\mathbb{R}^d} f(\boldsymbol{r}) \, \mathrm{e}^{-\mathrm{j}\langle \boldsymbol{\omega}, \boldsymbol{r} \rangle} \, \mathrm{d}\boldsymbol{r}$$

**Inverse Fourier transform (reconstruction formula)**

$$f(\boldsymbol{r}) = \mathcal{F}^{-1}\{f\} = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \hat{f}(\boldsymbol{\omega}) \, \mathrm{e}^{\mathrm{j}\langle \boldsymbol{\omega}, \boldsymbol{r} \rangle} \, \mathrm{d}\boldsymbol{\omega} \qquad \text{(a.e.)}$$

# Example imaging operator:
# Fourier transform is extensively used in MRI

**Fourier transform:** $\mathcal{F} : L_2(\mathbb{R}^d) \to L_2(\mathbb{R}^d)$

$$\hat{f}(\boldsymbol{\omega}) = \mathcal{F}\{f\} = \int_{\mathbb{R}^d} f(\boldsymbol{r})\, \mathrm{e}^{-\mathrm{j}\langle \boldsymbol{\omega}, \boldsymbol{r}\rangle}\, \mathrm{d}\boldsymbol{r}$$

**Inverse Fourier transform (reconstruction formula)**

$$f(\boldsymbol{r}) = \mathcal{F}^{-1}\{f\} = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \hat{f}(\boldsymbol{\omega})\, \mathrm{e}^{\mathrm{j}\langle \boldsymbol{\omega}, \boldsymbol{r}\rangle}\, \mathrm{d}\boldsymbol{\omega}$$

**As a measurement function:** $h_m(\boldsymbol{r}) = \mathrm{e}^{-\mathrm{j}\langle \boldsymbol{\omega}_m, \boldsymbol{r}\rangle}$   (*complex sinusoid*)

$$y_m = \langle h_m, f\rangle = \int_{\mathbb{R}^d} h_m(\boldsymbol{r}) f(\boldsymbol{r})\, \mathrm{d}\boldsymbol{r}$$

# Example imaging operator:
# Fourier transform is extensively used in MRI

# Example imaging operator:
# Fourier transform is extensively used in MRI

# Example imaging operator:
# Fourier transform is extensively used in MRI
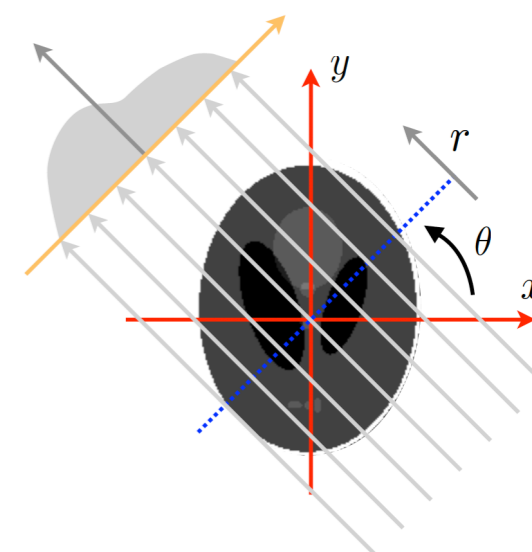
**Linear forward model for MRI**

$$\hat{s}(\boldsymbol{\omega}_m) = \int_{\mathbb{R}^3} s(\boldsymbol{r}) \mathrm{e}^{-\mathrm{j}\langle \boldsymbol{\omega}_m, \boldsymbol{r} \rangle} \, \mathrm{d}\boldsymbol{r}$$

sampling of Fourier transform

$$\boldsymbol{r} = (x, y, z) \quad \boldsymbol{\omega} = (\omega_x, \omega_y, \omega_z)$$

# Example imaging operator:
# Fourier transform is extensively used in MRI

**Linear forward model for MRI**

$$\hat{s}(\boldsymbol{\omega}_m) = \int_{\mathbb{R}^3} s(\boldsymbol{r}) \mathrm{e}^{-\mathrm{j}\langle \boldsymbol{\omega}_m, \boldsymbol{r}\rangle} \, \mathrm{d}\boldsymbol{r}$$

sampling of Fourier transform

$$\boldsymbol{r} = (x, y, z) \quad \boldsymbol{\omega} = (\omega_x, \omega_y, \omega_z)$$



(A) Completely sampled

(B) Uniformly under-sampled

(C) Incoherently under-sampled

(D) Variable density incoherently under-sampled

[Source]

# Example imaging operator:
# Fourier transform is extensively used in MRI

**Linear forward model for MRI**

$$\hat{s}(\boldsymbol{\omega}_m) = \int_{\mathbb{R}^3} s(\boldsymbol{r}) \mathrm{e}^{-\mathrm{j}\langle \boldsymbol{\omega}_m, \boldsymbol{r} \rangle} \, \mathrm{d}\boldsymbol{r}$$

**Extended forward model with coil sensitivity**

$$\hat{s}_w(\boldsymbol{\omega}_m) = \int_{\mathbb{R}^3} w(\boldsymbol{r}) s(\boldsymbol{r}) \mathrm{e}^{-\mathrm{j}\langle \boldsymbol{\omega}_m, \boldsymbol{r} \rangle} \, \mathrm{d}\boldsymbol{r}$$



(A) Completely sampled

(B) Uniformly under-sampled

(C) Incoherently under-sampled

(D) Variable density incoherently under-sampled

[Source]

# Example imaging operator:
# Radon transform is extensively used in tomography

# Example imaging operator:
# Radon transform is extensively used in tomography



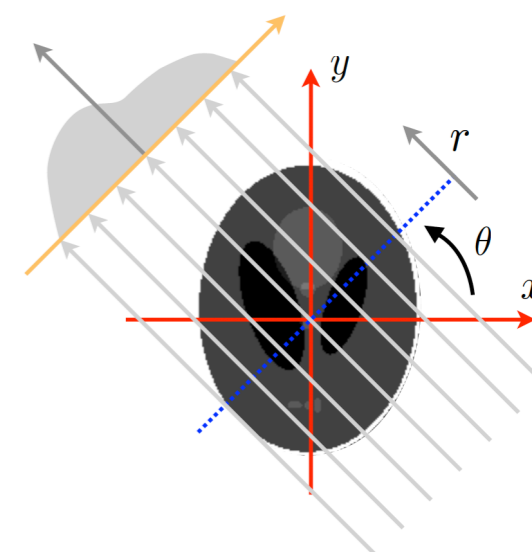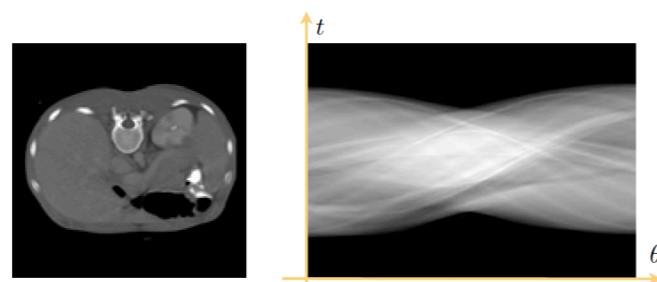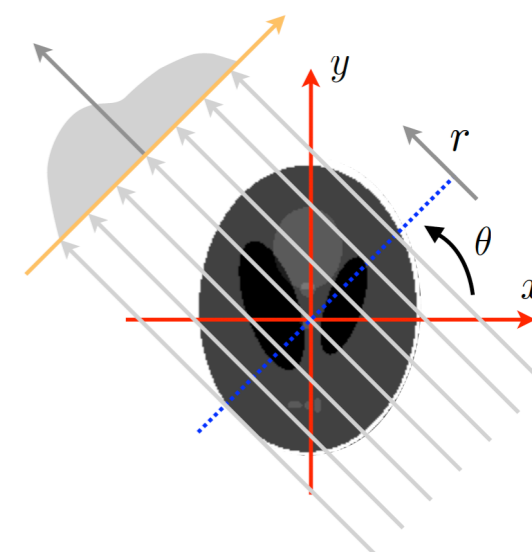Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Example imaging operator:
# Radon transform is extensively used in tomography

**Projection geometry:** $\boldsymbol{r} = t\boldsymbol{\theta} + r\boldsymbol{\theta}^{\perp}, \quad \boldsymbol{\theta} = (\cos\theta, \sin\theta)$



Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Example imaging operator:
# Radon transform is extensively used in tomography

**Projection geometry:** $\quad \boldsymbol{r} = t\boldsymbol{\theta} + r\boldsymbol{\theta}^{\perp}, \quad \boldsymbol{\theta} = (\cos\theta, \sin\theta)$

**Radon transform computes
line integrals of the object:**

$$R_{\theta}\{f(\boldsymbol{r})\}(t) = \int_{\mathbb{R}} f(t\boldsymbol{\theta} + r\boldsymbol{\theta}^{\perp}) \, \mathrm{d}r$$

$$= \int_{\mathbb{R}^2} f(\boldsymbol{r})\delta(t - \langle \boldsymbol{r}, \boldsymbol{\theta} \rangle) \, \mathrm{d}\boldsymbol{r}$$

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Example imaging operator:
# Radon transform is extensively used in tomography

**Projection geometry:** $\quad \boldsymbol{r} = t\boldsymbol{\theta} + r\boldsymbol{\theta}^{\perp}, \quad \boldsymbol{\theta} = (\cos\theta, \sin\theta)$

**Radon transform computes line integrals of the object:**

$$\mathsf{R}_\theta\{f(\boldsymbol{r})\}(t) = \int_{\mathbb{R}} f(t\boldsymbol{\theta} + r\boldsymbol{\theta}^{\perp})\,\mathrm{d}r$$

$$= \int_{\mathbb{R}^2} f(\boldsymbol{r})\delta(t - \langle \boldsymbol{r}, \boldsymbol{\theta}\rangle)\,\mathrm{d}\boldsymbol{r}$$

image and its sinogram

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Example imaging operator:
# Radon transform is extensively used in tomography

**Projection geometry:** $\quad \boldsymbol{r} = t\boldsymbol{\theta} + r\boldsymbol{\theta}^{\perp}, \quad \boldsymbol{\theta} = (\cos\theta, \sin\theta)$

**Radon transform computes line integrals of the object:**

$$\mathsf{R}_{\theta}\{f(\boldsymbol{r})\}(t) = \int_{\mathbb{R}} f(t\boldsymbol{\theta} + r\boldsymbol{\theta}^{\perp}) \, \mathrm{d}r$$

$$= \int_{\mathbb{R}^2} f(\boldsymbol{r})\delta(t - \langle \boldsymbol{r}, \boldsymbol{\theta} \rangle) \, \mathrm{d}\boldsymbol{r}$$



**As a measurement function:** $\quad h_m(\boldsymbol{r}) = \delta(t_m - \langle \boldsymbol{r}, \theta_m \rangle)$

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Central slice theorem relates projections to the Fourier transform of the object

# Central slice theorem relates projections to the Fourier transform of the object

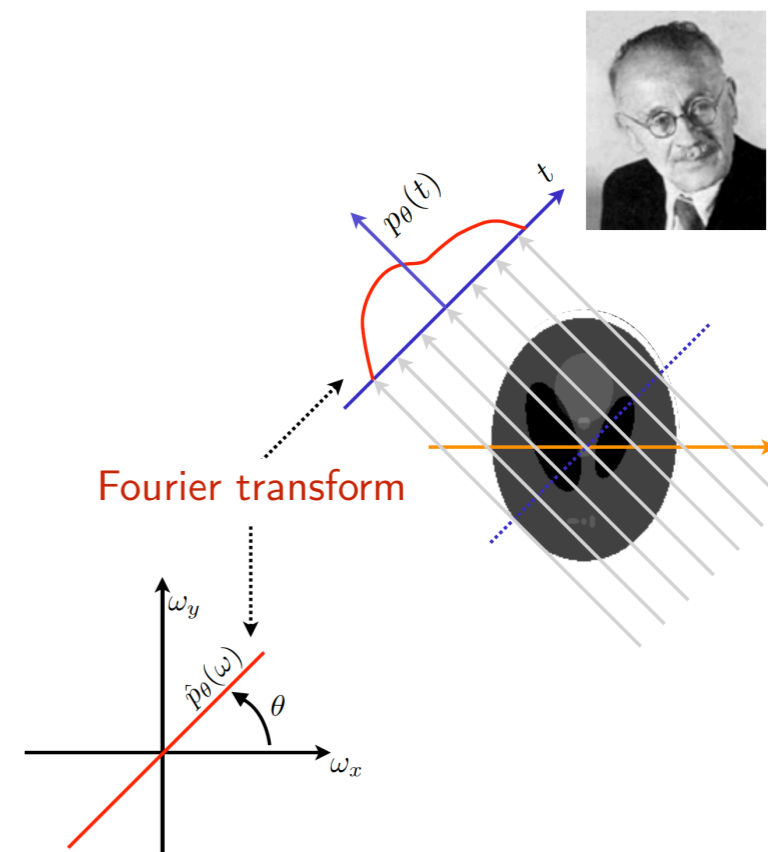**Radon transform:** $p_\theta(t) = R_\theta\{f\}(t, \theta)$



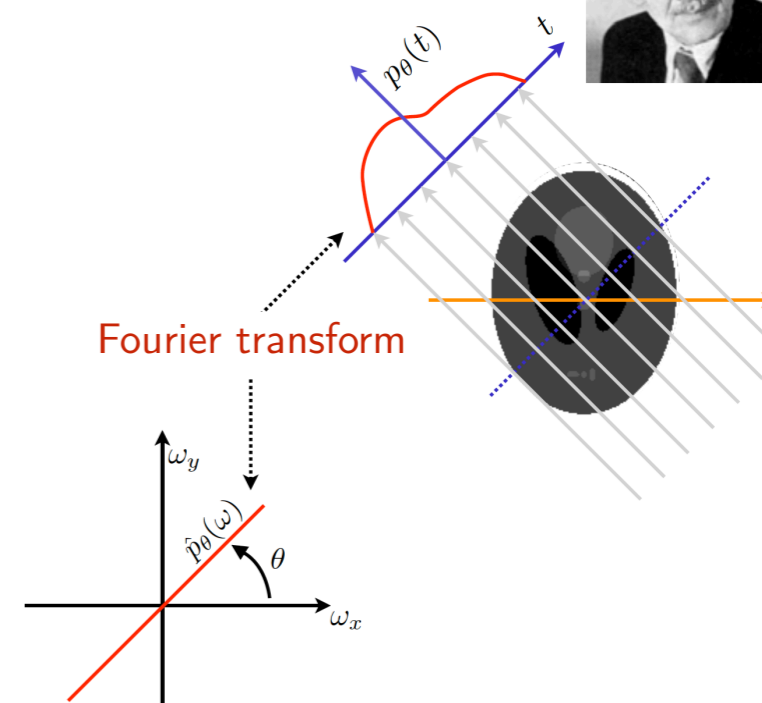Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Central slice theorem relates projections to the Fourier transform of the object

**Radon transform:** $\quad p_\theta(t) = \mathsf{R}_\theta\{f\}(t, \theta)$

**1D and 2D Fourier relationships:**

$$\hat{p}_\theta(\omega) = \mathcal{F}_{1D}\{p_\theta\}(\omega)$$ 1D Fourier of data

$$\hat{f}(\boldsymbol{\omega}) = \mathcal{F}_{2D}\{f\}(\boldsymbol{\omega}) = \hat{f}_{pol}(\omega, \theta)$$ 2D Fourier of image

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

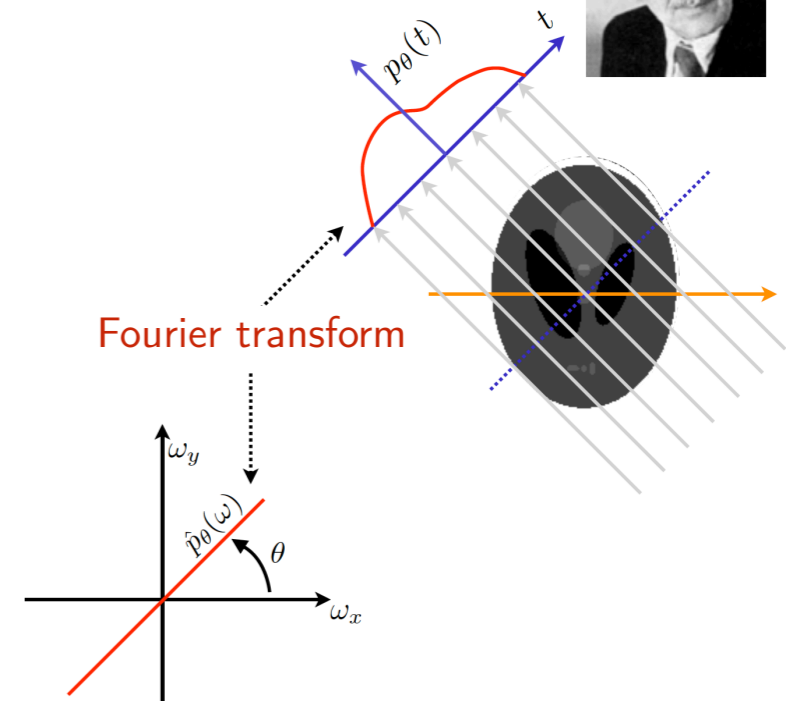# Central slice theorem relates projections to the Fourier transform of the object

**Radon transform:** $p_\theta(t) = \mathsf{R}_\theta\{f\}(t, \theta)$

**1D and 2D Fourier relationships:**

$$\hat{p}_\theta(\omega) = \mathcal{F}_{\mathrm{1D}}\{p_\theta\}(\omega)$$

$$\hat{f}(\boldsymbol{\omega}) = \mathcal{F}_{\mathrm{2D}}\{f\}(\boldsymbol{\omega}) = \hat{f}_{\mathrm{pol}}(\omega, \theta)$$

**Central-slice theorem relates projections to Fourier sampling:**

$$\hat{p}_\theta(\omega) = \hat{f}(\omega\cos\theta, \omega\sin\theta) = \hat{f}_{\mathrm{pol}}(\omega, \theta)$$

Establishes Fourier relationship between data and image



Fourier transform

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Central slice theorem relates projections to the Fourier transform of the object



**Radon transform:** $\quad p_\theta(t) = \mathsf{R}_\theta\{f\}(t, \theta)$

**1D and 2D Fourier relationships:**

$$\hat{p}_\theta(\omega) = \mathcal{F}_{\text{1D}}\{p_\theta\}(\omega)$$

$$\hat{f}(\boldsymbol{\omega}) = \mathcal{F}_{\text{2D}}\{f\}(\boldsymbol{\omega}) = \hat{f}_{\text{pol}}(\omega, \theta)$$

**Central-slice theorem relates projections to Fourier sampling:**

$$\hat{p}_\theta(\omega) = \hat{f}(\omega\cos\theta, \omega\sin\theta) = \hat{f}_{\text{pol}}(\omega, \theta)$$

Establishes Fourier relationship between data and image

**Proof for angle zero:**

$$\hat{f}(\omega, 0) = \int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty} f(x, y)\,\mathrm{e}^{-\mathrm{j}\omega x}\,\mathrm{d}x\mathrm{d}y = \int_{-\infty}^{+\infty}\underbrace{\left(\int_{-\infty}^{+\infty} f(x, y)\,\mathrm{d}y\right)}_{p_0(x)}\mathrm{e}^{-\mathrm{j}\omega x}\,\mathrm{d}x = \hat{p}_0(x)$$

# Central slice theorem relates projections to the Fourier transform of the object



**Radon transform:** $\quad p_\theta(t) = \mathrm{R}_\theta\{f\}(t, \theta)$

**1D and 2D Fourier relationships:**

$$\hat{p}_\theta(\omega) = \mathcal{F}_{\text{1D}}\{p_\theta\}(\omega)$$

$$\hat{f}(\boldsymbol{\omega}) = \mathcal{F}_{\text{2D}}\{f\}(\boldsymbol{\omega}) = \hat{f}_{\text{pol}}(\omega, \theta)$$

**Central-slice theorem relates projections to Fourier sampling:**

$$\hat{p}_\theta(\omega) = \hat{f}(\omega\cos\theta, \omega\sin\theta) = \hat{f}_{\text{pol}}(\omega, \theta)$$

Establishes Fourier relationship between data and image

**Proof for angle zero:**   Question: How to generalize to other angles?

$$\hat{f}(\omega, 0) = \int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty} f(x, y)\, \mathrm{e}^{-\mathrm{j}\omega x}\, \mathrm{d}x\mathrm{d}y = \int_{-\infty}^{+\infty}\underbrace{\left(\int_{-\infty}^{+\infty} f(x, y)\, \mathrm{d}y\right)}_{p_0(x)} \mathrm{e}^{-\mathrm{j}\omega x}\, \mathrm{d}x = \hat{p}_0(x)$$

# Most imaging systems can be characterized with a forward model

# Most imaging systems can be characterized with a forward model

| Modality | Radiation | Forward model | Variations |
|---|---|---|---|
| 2D or 3D tomography | coherent x-ray | $y_i = \mathrm{R}_{\boldsymbol{\theta}_i} x$ | parallel, cone beam, spiral sampling |
| 3D deconvolution microscopy | fluorescence | $y = \mathrm{H} x$ | brightfield, confocal, light sheet |
| structured illumination microscopy (SIM) | fluorescence | $y_i = \mathrm{H} \mathrm{W}_i x$<br><br>$\mathrm{H}$: PSF of microscope<br>$\mathrm{W}_i$: illumination pattern | full 3D reconstruction, non-sinusoidal patterns |
| Positron Emission Tomography (PET) | gamma rays | $y_i = \mathrm{H}_{\boldsymbol{\theta}_i} x$ | list mode with time-of-flight |
| Magnetic resonance imaging (MRI) | radio frequency | $y = \mathrm{F} x$ | uniform or non-uniform sampling in k space |
| Cardiac MRI (parallel, non-uniform) | radio frequency | $y_{t,i} = \mathrm{F}_t \mathrm{W}_i x$<br><br>$\mathrm{W}_i$: coil sensitivity | gated or not, retrospective registration |
| Optical diffraction tomography | coherent light | $y_i = \mathrm{W}_i \mathrm{F}_i x$ | with holography or grating interferometry |

# Most imaging systems can be characterized with a forward model

| Modality | Radiation | Forward model | Variations |
|---|---|---|---|
| 2D or 3D tomography | coherent x-ray | $y_i = \mathrm{R}_{\boldsymbol{\theta}_i} x$ | parallel, cone beam, spiral sampling |
| 3D deconvolution microscopy | fluorescence | $y = \mathrm{H} x$ | brightfield, confocal, light sheet |
| structured illumination microscopy (SIM) | fluorescence | $y_i = \mathrm{H}\mathrm{W}_i x$<br>H: PSF of microscope<br>$\mathrm{W}_i$: illumination pattern | full 3D reconstruction, non-sinusoidal patterns |
| Positron Emission Tomography (PET) | gamma rays | $y_i = \mathrm{H}_{\boldsymbol{\theta}_i} x$ | list mode with time-of-flight |
| Magnetic resonance imaging (MRI) | radio frequency | $y = \mathrm{F} x$ | uniform or non-uniform sampling in k space |
| Cardiac MRI (parallel, non-uniform) | radio frequency | $y_{t,i} = \mathrm{F}_t \mathrm{W}_i x$<br>$\mathrm{W}_i$: coil sensitivity | gated or not, retrospective registration |
| Optical diffraction tomography | coherent light | $y_i = \mathrm{W}_i \mathrm{F}_i x$ | with holography or grating interferometry |

Currently active collaborations at CIG

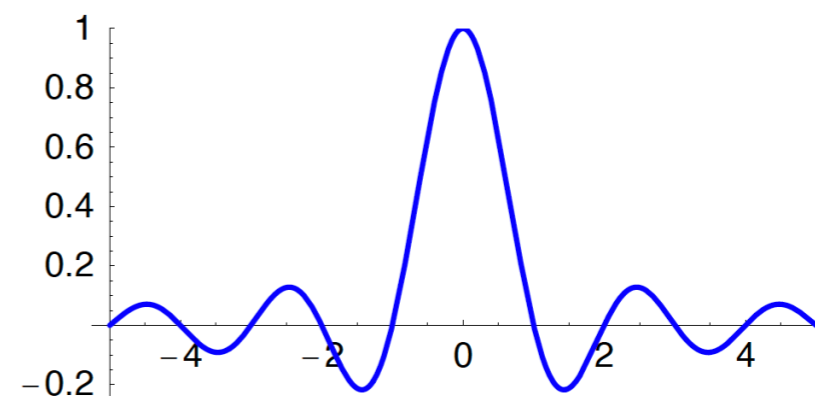# Discretization: Continuous domain formalism easily reduces to a noisy linear system
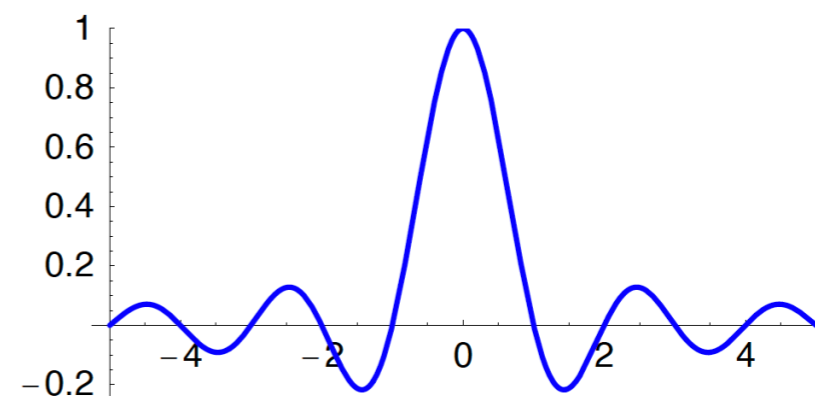
# Discretization: Continuous domain formalism easily reduces to a noisy linear system

**Representation with basis functions:**

$$f(\boldsymbol{r}) = \sum_{\boldsymbol{k}\in\Omega} f[\boldsymbol{k}]\beta_{\boldsymbol{k}}(\boldsymbol{r})$$

Question: What type of representation is offered by **sinc**?

# Discretization: Continuous domain formalism easily reduces to a noisy linear system

**Representation with basis functions:**

$$f(\boldsymbol{r}) = \sum_{\boldsymbol{k} \in \Omega} f[\boldsymbol{k}] \beta_{\boldsymbol{k}}(\boldsymbol{r})$$

**Signal vector:** $\mathbf{f} = \{f[\boldsymbol{k}]\}_{\boldsymbol{k} \in \Omega} \in \mathbb{R}^n$

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Discretization: Continuous domain formalism easily reduces to a noisy linear system

**Representation with basis functions:**

$$f(\boldsymbol{r}) = \sum_{\boldsymbol{k} \in \Omega} f[\boldsymbol{k}] \beta_{\boldsymbol{k}}(\boldsymbol{r})$$



**Signal vector:** $\quad \mathbf{f} = \{f[\boldsymbol{k}]\}_{\boldsymbol{k} \in \Omega} \in \mathbb{R}^n$

**Discretized measurement model:**

$$y_i = \int_{\mathbb{R}^d} f(\boldsymbol{r}) h_i(\boldsymbol{r}) \, \mathrm{d}\boldsymbol{r} + e_i = \langle f, h_i \rangle + e_i, \quad (i = 1, \ldots, m)$$

Question: What are the sources of noise?

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# Discretization: Continuous domain formalism easily reduces to a noisy linear system

**Representation with basis functions:**

$$f(\boldsymbol{r}) = \sum_{\boldsymbol{k} \in \Omega} f[\boldsymbol{k}] \beta_{\boldsymbol{k}}(\boldsymbol{r})$$

**Signal vector:** $\quad \mathbf{f} = \{f[\boldsymbol{k}]\}_{\boldsymbol{k} \in \Omega} \in \mathbb{R}^n$

**Discretized measurement model:**

$$y_i = \int_{\mathbb{R}^d} f(\boldsymbol{r}) h_i(\boldsymbol{r}) \, \mathrm{d}\boldsymbol{r} + e_i = \langle f, h_i \rangle + e_i, \quad (i = 1, \dots, m)$$

$$\Rightarrow \quad \boxed{\mathbf{y} = \mathbf{H}\mathbf{f} + \mathbf{e}} \qquad [\mathbf{H}]_{i,\boldsymbol{k}} = \langle h_i, \beta_{\boldsymbol{k}} \rangle = \int_{\mathbb{R}^d} h_m(\boldsymbol{r}) \beta_{\boldsymbol{k}}(\boldsymbol{r}) \, \mathrm{d}\boldsymbol{r}$$

linear system of equations

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# To conclude "forward models"

Many imaging problems reduce to solving large and noisy linear systems

$$y = Hf + e$$

Setting up the right forward model is a big step towards being able to form high quality images

## Inverse problems in bio-imaging

- Linear forward model

$$y = Hs + n$$
$$y = Hf + e$$

f

H

noise



Integral operator

s

Problem: recovers from noisy measurements y

- The easy scenario

Inverse problem is well posed if $t > 0$, for all, for all, $s$, $\|x\| \leq 0 \|Hs\| \|Hs\|$

# Today we will talk about

◉ Forward models in imaging
  Relating the unknowns to the measured data

◉ **Notions of ill-posedness and regularization**
  When measurements are not enough

◉ Optimization at large scales
  When analytical solutions are not enough

◉ Plug-and-Play Priors (PnP) at large scales
  When traditional optimization is not enough

# What makes imaging inverse problems difficult?

# Inverse problems in bio-imaging

- Linear forward model

$$y = Hf + e$$

noise

**H**

Integral operator

**n**

f

s

Problem: recover f from noisy measurements y

- The easy scenario

Inverse problem is **well posed** if $0 < c_0 \leq t$, for all s $\frac{\|Hs\|}{\|s\|}$

$$\Rightarrow \quad s \approx H^{-1}y$$

- Back projection (poor man's solution) $s \approx H^T y$

# Inverse problems in bio-imaging
# What makes imaging inverse problems difficult?

- Linear forward model

$$y = Hs + n$$
$$y = Hf + e$$

noise

**f**

**H**

**H**

Integral operator

**n**

s s

Problem: recover **f** from noisy measurements **y**

- The easy scenario

Question: Why can't we simply compute the inverse $\mathbf{f} = \mathbf{H}^{-1}\mathbf{y}$?

$$\Rightarrow \quad \mathbf{s} \approx \mathbf{H}^{-1}\mathbf{y}$$

- Back projection (poor man's solution) $\mathbf{s} \approx \mathbf{H}^T\mathbf{y}$

# What makes imaging inverse problems difficult?

- Linear forward model

$$y = \mathbf{H}\mathbf{f} + \mathbf{n}$$
$$y = \mathbf{H}\mathbf{f} + \mathbf{e}$$

noise

$\mathbf{f}$

$\mathbf{H}$

Integral operator

$\mathbf{n}$

$\mathbf{s}$

Problem: recover $\mathbf{f}$ from noisy measurements $y$

- The easy scenario

Question: Why can't we simply compute the inverse $\mathbf{f} = \mathbf{H}^{-1}y$?

$$\Rightarrow \quad \mathbf{f} \approx \mathbf{H}^{-1}y$$

- Back-projection (poor resolution) $\approx \mathbf{H}^T y$

**1) Difficult to invert the matrix as it is non-square or too large**

# Inverse problems in bio-imaging
# What makes imaging inverse problems difficult?

- Linear forward model

$$y = Hs + n$$
$$y = Hf + e$$

noise

**H**   **H**   **H**

Integral operator

**n**

**f**

**s**   **s**

Problem: recover **f** from noisy measurements **y**

- The easy scenario

Question: Why can't we simply compute the inverse $f = H^{-1}y$?

$$\Rightarrow \quad s \approx H^{-1}y$$

- Backprojection (poor resolution) $\approx H^T y$

**1) Difficult to invert the matrix as it is non-square or too large**

**2) Measurements do not uniquely describe the object**

# Inverse problems in bio-imaging
# What makes imaging inverse problems difficult?

- Linear forward model

$$y = \mathbf{H}s + n$$
$$y = \mathbf{H}\mathbf{f} + \mathbf{e}$$

noise

$f$

$\mathbf{H}$

$\mathbf{H}$

Integral operator

$\mathbf{n}$

$s$     $s$

Problem: recover $\mathbf{f}$ from noisy measurements $y$

- The easy scenario

Question: Why can't we simply compute the inverse $\mathbf{f} = \mathbf{H}^{-1}y$?

$$\Rightarrow \quad s \approx \mathbf{H}^{-1}y$$

- Back projection (poor's man solution) $\approx \mathbf{H}^{T}y$

1) Difficult to invert the matrix as it is non-square or too large

2) Measurements do not uniquely describe the object

$\mathcal{F}$

# Regularization framework enables the selection of the most suitable solution among alternatives

# Regularization framework enables the selection of the most suitable solution among alternatives

Consider a noisy linear system with noise of bounded norm

# Regularization framework enables the selection of the most suitable solution among alternatives

**Consider a noisy linear system with noise of bounded norm**



$$\mathbf{y} = \mathbf{Hf} + \mathbf{e}$$ such that $$\|\mathbf{y} - \mathbf{Hf}\|_{\ell_2}^2 \leq \sigma^2$$

# Regularization framework enables the selection of the most suitable solution among alternatives

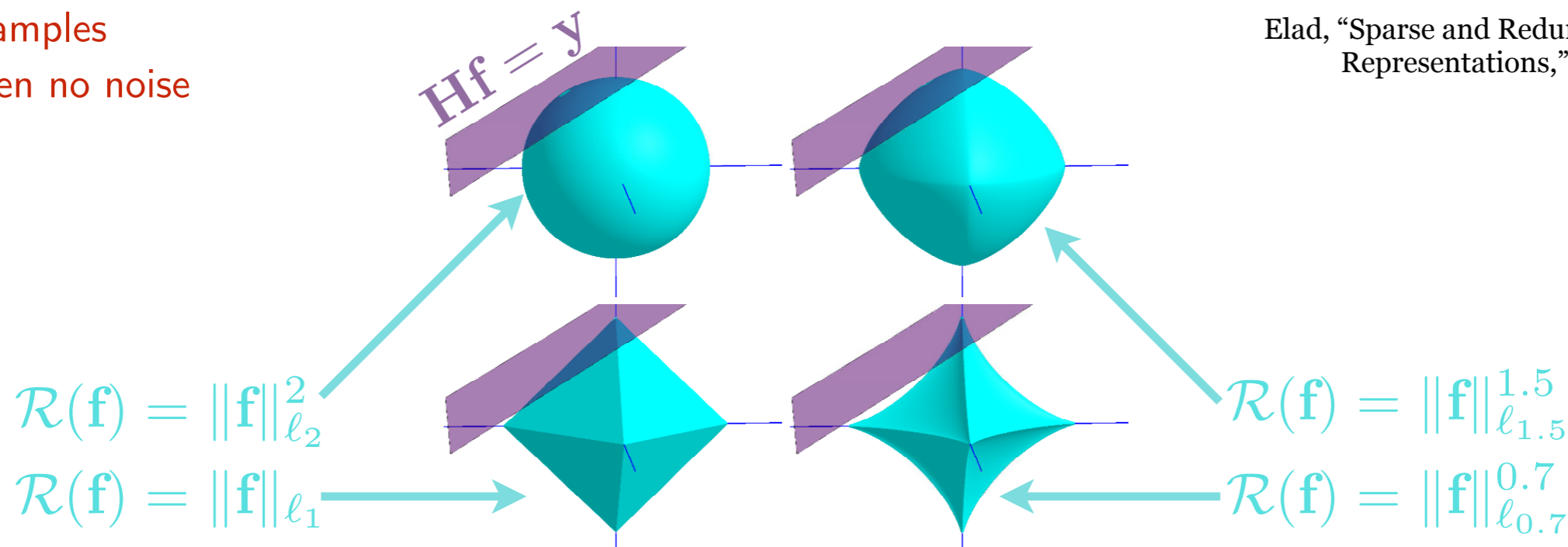**Consider a noisy linear system with noise of bounded norm**



$$\mathbf{y} = \mathbf{Hf} + \mathbf{e} \quad \text{such that} \quad \|\mathbf{y} - \mathbf{Hf}\|_{\ell_2}^2 \leq \sigma^2$$
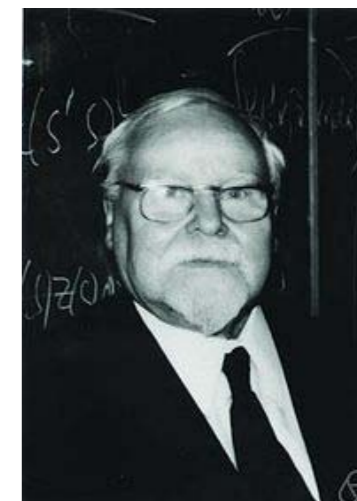
**We consider a constrained optimization problem**

$$\text{minimize } \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{Hf} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2$$

- The "regularizer" picks the solution which we <u>think</u> is best
- Allows us to infuse <u>prior knowledge</u> into the problem

# Regularization framework enables the selection of the most suitable solution among alternatives

**Consider a noisy linear system with noise of bounded norm**

$$\mathbf{y} = \mathbf{H}\mathbf{f} + \mathbf{e} \quad \text{such that} \quad \|\mathbf{y} - \mathbf{H}\mathbf{f}\|_{\ell_2}^2 \leq \sigma^2$$

$$\mathbf{y} \qquad \mathbf{H} \qquad \mathbf{f} \quad \mathbf{e}$$



**We consider a constrained optimization problem**

$$\text{minimize } \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2$$

Examples

when no noise

Elad, "Sparse and Redundant Representations," 2010

$\mathbf{H}\mathbf{f} = \mathbf{y}$

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{f}\|_{\ell_2}^2$$

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{f}\|_{\ell_1}$$

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{f}\|_{\ell_{1.5}}^{1.5}$$

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{f}\|_{\ell_{0.7}}^{0.7}$$

# Question: How to regularize in imaging?

# Question: How to regularize in imaging?
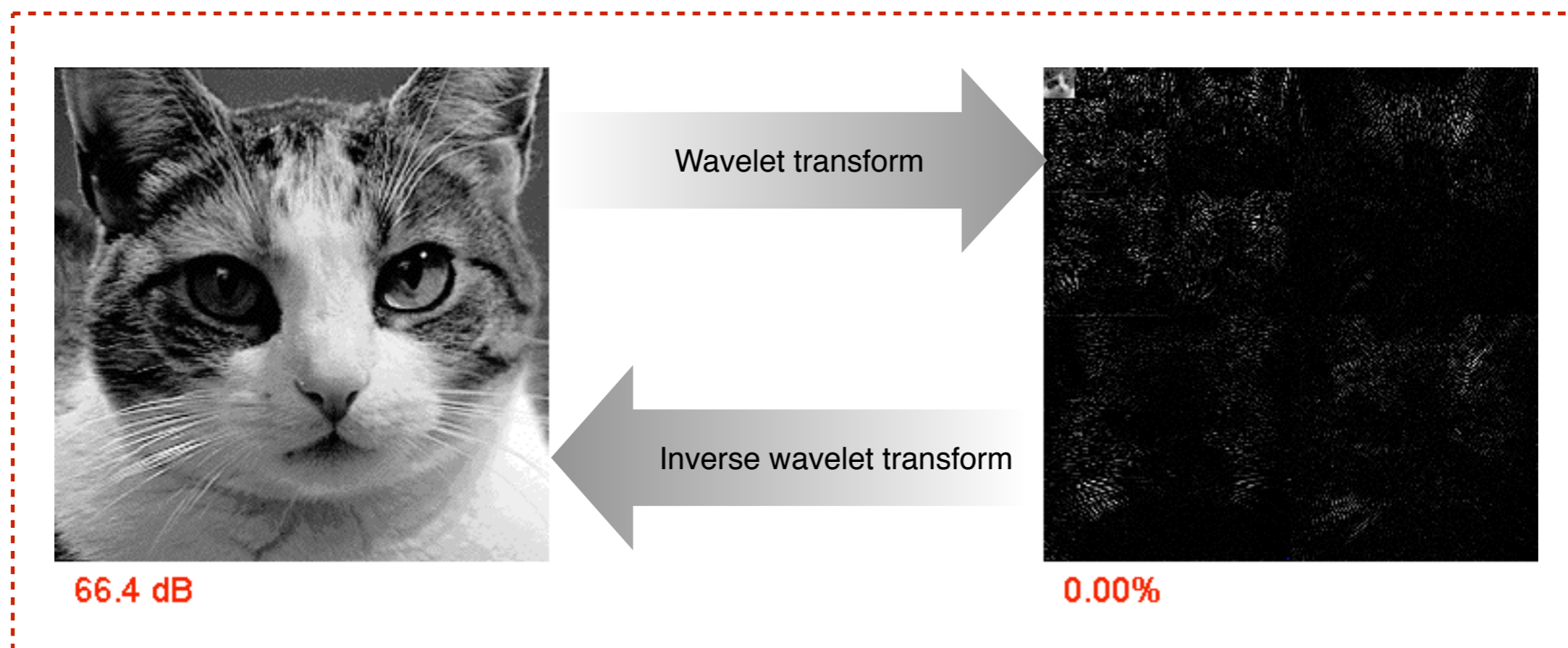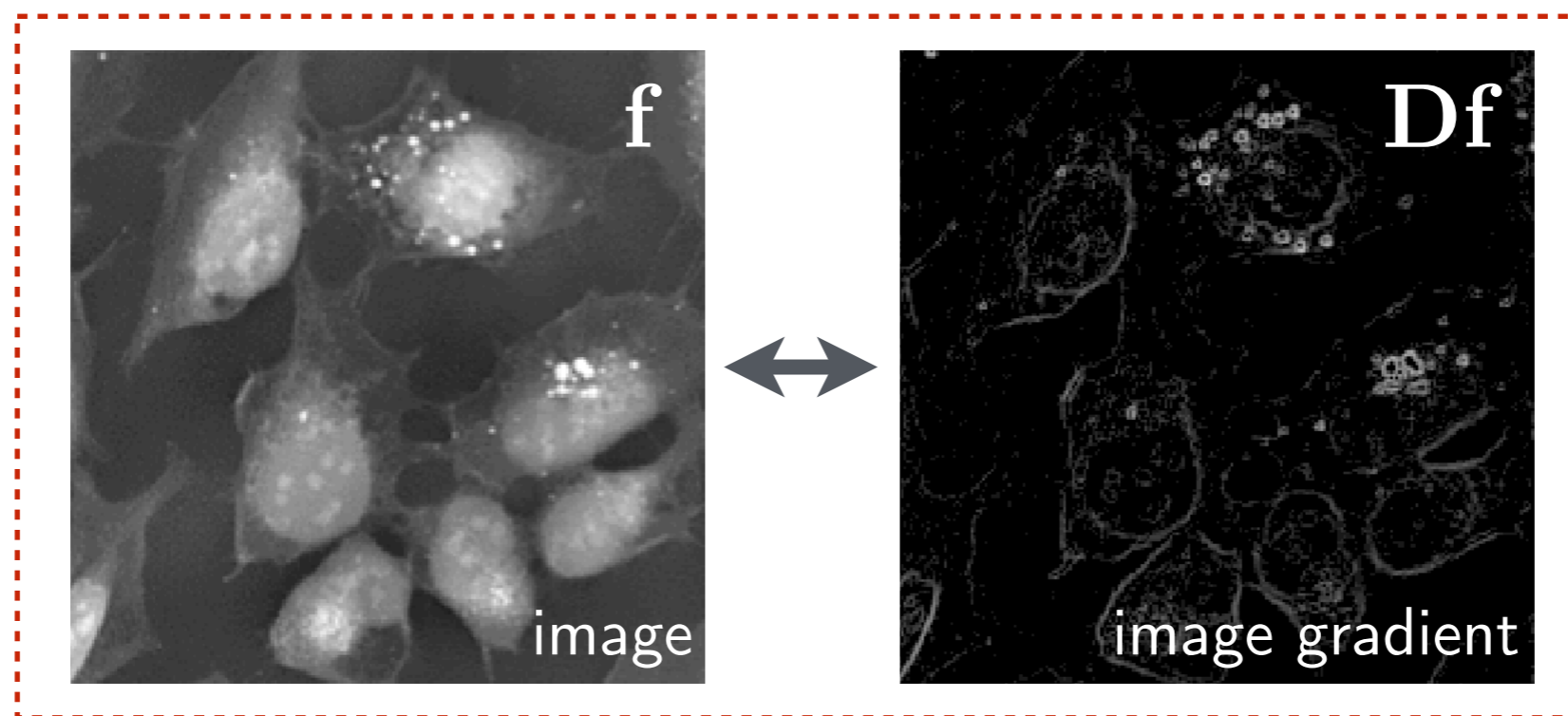
minimize $\mathcal{R}(\mathbf{f})$ subject to $\|\mathbf{Hf} - \mathbf{y}\|^2_{\ell_2} \leq \sigma^2$

## Question: How to regularize in imaging?

■ Dealing with **ill-posed problems**: Tikhonov **regularization**

**Classical approach: Tikhonov regularization**

$\mathcal{R}(s) = \|\mathbf{L}s\|_2^2$: regularization (or smoothness) functional

$\mathbf{L}$: regularization operator (i.e., Gradient)

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_2}^2$$

*Assumption:*

*image is smooth*

$$\min_{\mathbf{s}} \mathcal{R}(\mathbf{s}) \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2 \leq \sigma^2$$

*Andrey N. Tikhonov (1906-1993)*

■ Equivalent variational problem

$$\mathbf{s}^{\star} = \arg\min \underbrace{\|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2}_{\text{data consistency}} + \underbrace{\lambda\|\mathbf{L}\mathbf{s}\|_2^2}_{\text{regularization}}$$

Formal linear solution: $\quad \mathbf{s} = (\mathbf{H}^T\mathbf{H} + \lambda\mathbf{L}^T\mathbf{L})^{-1}\mathbf{H}^T\mathbf{y} = \mathbf{R}_\lambda \cdot \mathbf{y}$

minimize $\mathcal{R}(\mathbf{f})$ subject to $\|\mathbf{H}\mathbf{f} - \mathbf{y}\|_2^2 \leq \sigma^2$

**Statistical formulation (20th century)**

# Question: How to regularize in imaging?

**Classical approach: Tikhonov regularization**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \widehat{\mathbf{f}}_{\text{Tikh}} = (\mathbf{D}^{\mathsf{T}}\mathbf{D})^{-1}\mathbf{H}^{\mathsf{T}}\left[\mathbf{H}(\mathbf{D}^{\mathsf{T}}\mathbf{D})^{-1}\mathbf{H}^{\mathsf{T}}\right]^{-1}\mathbf{y}$$

<span style="color:red">unique closed-form solution</span>
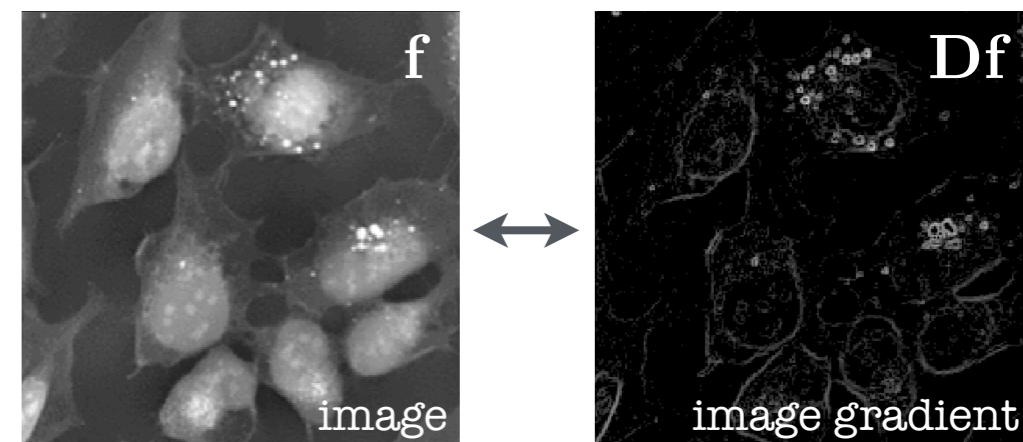
minimize $\mathcal{R}(\mathbf{f})$ subject to $\|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2$

# Question: How to regularize in imaging?

**Classical approach: Tikhonov regularization**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \widehat{\mathbf{f}}_{\text{Tikh}} = (\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\left[\mathbf{H}(\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\right]^{-1}\mathbf{y}$$

Assumption:

image is smooth

Question: Is image smoothness a reasonable assumption?

minimize $\mathcal{R}(\mathbf{f})$ subject to $\|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2 \le \sigma^2$
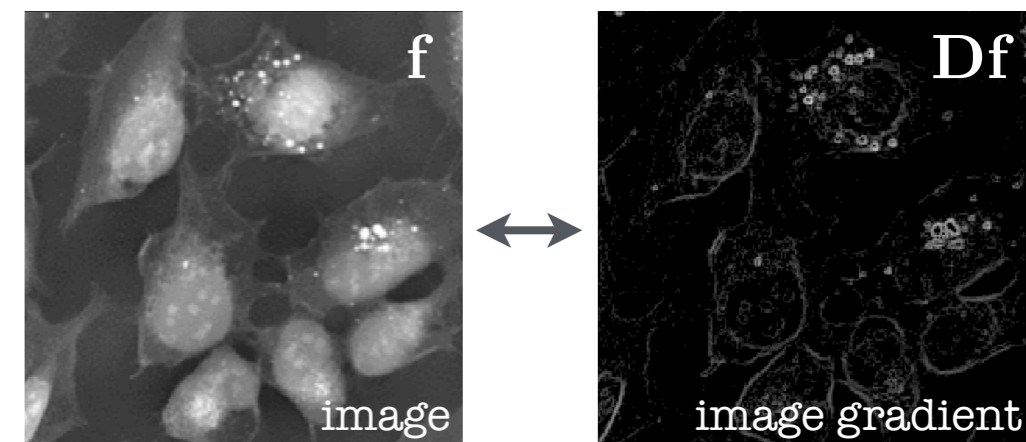
# Question: How to regularize in imaging?

**Classical approach: Tikhonov regularization**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{Df}\|_{\ell_2}^2 \quad \Rightarrow \quad \widehat{\mathbf{f}}_{\text{Tikh}} = (\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\left[\mathbf{H}(\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\right]^{-1}\mathbf{y}$$

**Modern approach: Transform-domain sparsity**



66.4 dB         Wavelet transform         Inverse wavelet transform         0.00%

minimize $\mathcal{R}(\mathbf{f})$ subject to $\|\mathbf{Hf} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2$
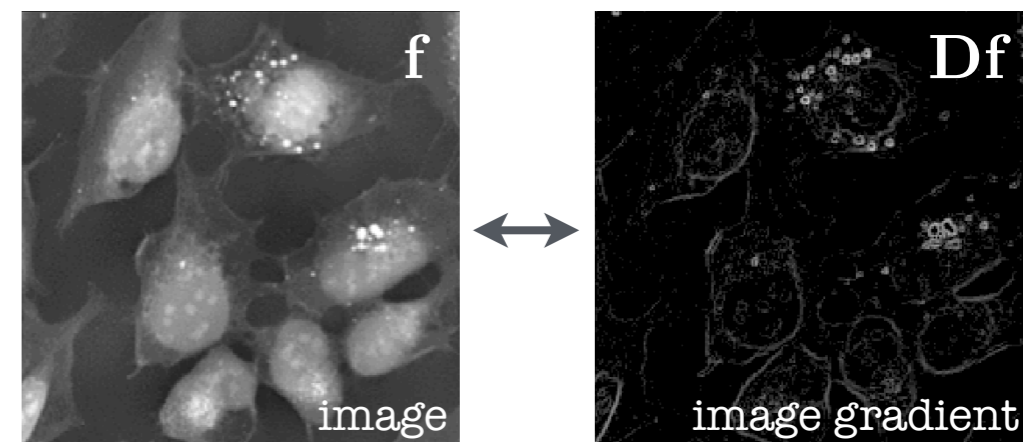
# Question: How to regularize in imaging?

**Classical approach: Tikhonov regularization**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{Df}\|_{\ell_2}^2 \quad \Rightarrow \quad \widehat{\mathbf{f}}_{\text{Tikh}} = (\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\left[\mathbf{H}(\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\right]^{-1}\mathbf{y}$$

**Modern approach: Transform-domain sparsity**



$$\text{minimize } \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{Hf} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2$$

# Question: How to regularize in imaging?

Classical approach: Tikhonov regularization

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{Df}\|_{\ell_2}^2 \quad \Rightarrow \quad \widehat{\mathbf{f}}_{\text{Tikh}} = (\mathbf{D}^\top \mathbf{D})^{-1}\mathbf{H}^\top \left[\mathbf{H}(\mathbf{D}^\top \mathbf{D})^{-1}\mathbf{H}^\top\right]^{-1}\mathbf{y}$$

Modern approach: Transform-domain sparsity



f

Df

image &harr; image gradient

# Question: How to regularize in imaging?

**Classical approach: Tikhonov regularization**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \widehat{\mathbf{f}}_{\mathsf{Tikh}} = (\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\left[\mathbf{H}(\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\right]^{-1}\mathbf{y}$$

**Modern approach: Transform-domain sparsity**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_0} = \#\{i : [\mathbf{D}\mathbf{f}]_i \neq 0\}$$

intractable nonconvex optimiazation



f · image

Df · image gradient

# Question: How to regularize in imaging?

**Classical approach: Tikhonov regularization**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \widehat{\mathbf{f}}_{\mathsf{Tikh}} = (\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\left[\mathbf{H}(\mathbf{D}^\mathsf{T}\mathbf{D})^{-1}\mathbf{H}^\mathsf{T}\right]^{-1}\mathbf{y}$$

**Modern approach: Transform-domain sparsity**

$$\mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_1}$$

- convex (but nondifferentiable)
- promotes sparsity



f      Df

image      image gradient

# To conclude "regularization"

Many imaging problems are ill-posed:
there are infinitely many solutions

$$y = Hf + e$$

Regularization is a strategy to select the
solution that "makes sense"

$$\text{minimize } \mathcal{R}(f) \text{ subject to } \|Hf - y\|_{\ell_2}^2 \leq \sigma^2$$

Classical image regularizers are linear,
but increasingly they are nonlinear

(20th) $\quad \mathcal{R}(f) = \|Df\|_{\ell_2}^2 \quad \Rightarrow \quad \mathcal{R}(f) = \|Df\|_{\ell_1} \quad$ (21st)

# Today we will talk about

◉ Forward models in imaging
Relating the unknowns to the measured data

◉ Notions of ill-posedness and regularization
When measurements are not enough

◉ **Optimization at large scales**
When analytical solutions are not enough

◉ Plug-and-Play Priors (PnP) at large scales
When traditional optimization is not enough

# Proximal operator corresponds to image denoising

# Proximal operator corresponds to image denoising

**A more convenient formulation**

$$\min \; \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2 \Leftrightarrow \min_{\mathbf{f}} \left\{ \frac{1}{2}\|\mathbf{y} - \mathbf{H}\mathbf{f}\|_{\ell_2}^2 + \lambda \mathcal{R}(\mathbf{f}) \right\}$$

constrained optimization          unconstrained optimization

# Proximal operator corresponds to image denoising

**A more convenient formulation**

$$\min \ \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{Hf} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2 \ \Leftrightarrow \ \min_{\mathbf{f}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{Hf}\|_{\ell_2}^2 + \lambda \mathcal{R}(\mathbf{f}) \right\}$$

**Image denoising corresponds to identity measurement matrix**

$$\min_{\mathbf{f}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{f}\|_{\ell_2}^2 + \lambda \mathcal{R}(\mathbf{f}) \right\}$$

Question: Can you comment on convexity?

# Proximal operator corresponds to image denoising

A more convenient formulation

$$\min \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2 \iff \min_{\mathbf{f}} \left\{ \frac{1}{2}\|\mathbf{y} - \mathbf{H}\mathbf{f}\|_{\ell_2}^2 + \lambda\mathcal{R}(\mathbf{f}) \right\}$$
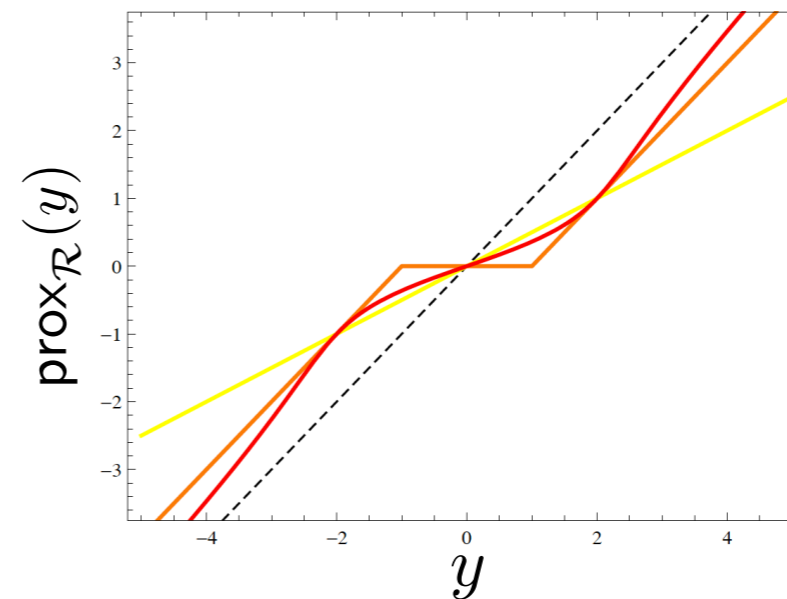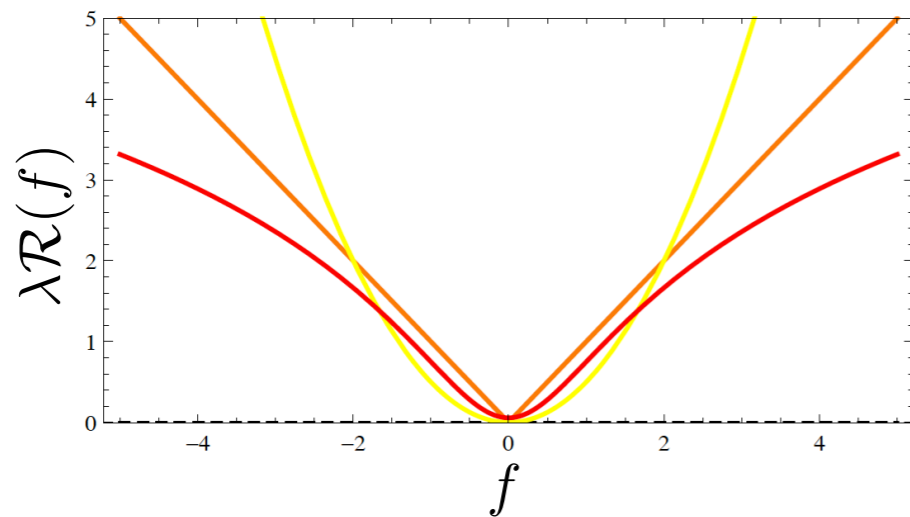
Image denoising corresponds to
identity measurement matrix

$$\min_{\mathbf{f}} \left\{ \frac{1}{2}\|\mathbf{y} - \mathbf{f}\|_{\ell_2}^2 + \lambda\mathcal{R}(\mathbf{f}) \right\}$$

For a convex regularizer,
the objective is strongly convex
=> there is a <u>unique</u> minimizer

# Proximal operator corresponds to image denoising

A more convenient formulation

$$\min \, \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{Hf} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2 \iff \min_{\mathbf{f}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{Hf}\|_{\ell_2}^2 + \lambda \mathcal{R}(\mathbf{f}) \right\}$$

Image denoising corresponds to
identity measurement matrix

$$\min_{\mathbf{f}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{f}\|_{\ell_2}^2 + \lambda \mathcal{R}(\mathbf{f}) \right\}$$

We can thus define the prox operator that
solves the denoising problem

$$\text{prox}_{\lambda \mathcal{R}}(\mathbf{y}) \triangleq \arg\min_{\mathbf{f}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{f}\|_{\ell_2}^2 + \lambda \mathcal{R}(\mathbf{f}) \right\}$$

# Proximal operator corresponds to image denoising

### Some examples of poitwise proximals



- $\blacksquare$ linear attenuation $\quad\quad\quad \ell_2$ minimization
- $\blacksquare$ soft-threshold $\quad\quad\quad\quad \ell_1$ minimization
- $\blacksquare$ shrinkage function $\quad\approx\quad \ell_p$ relaxation for $p \to 0$

Unser and Tafti, "An Introduction to Sparse Stochastic Processes," 2014

# FISTA and ADMM are two popular algorithms for large-scale and nonsmooth optimization

# FISTA and ADMM are two popular algorithms for large-scale and nonsmooth optimization

**Consider the objective function**

$$\boxed{\mathcal{C}(\mathbf{f}) = \mathcal{D}(\mathbf{f}) + \mathcal{R}(\mathbf{f})} \quad \text{where} \quad \mathcal{D}(\mathbf{f}) \triangleq \frac{1}{2}\|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2$$

data fit + regularizer

# FISTA and ADMM are two popular algorithms for large-scale and nonsmooth optimization

**Consider the objective function**

$$\mathcal{C}(\mathbf{f}) = \mathcal{D}(\mathbf{f}) + \mathcal{R}(\mathbf{f}) \quad \text{where} \quad \mathcal{D}(\mathbf{f}) \triangleq \frac{1}{2}\|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2$$

**Fast iterative shrinkage/thresholding algorithm (FISTA) vs. Alternating direction method of multipliers (ADMM)**

$$\mathbf{z}^k \leftarrow \mathbf{s}^{k-1} - \gamma\nabla\mathcal{D}(\mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \text{prox}_{\gamma\mathcal{R}}(\mathbf{z}^k)$$

$$\mathbf{s}^k \leftarrow \mathbf{f}^k + ((q_{k-1} - 1)/q_k)(\mathbf{f}^k - \mathbf{f}^{k-1})$$

$$\mathbf{z}^k \leftarrow \text{prox}_{\gamma\mathcal{D}}(\mathbf{f}^{k-1} - \mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \text{prox}_{\gamma\mathcal{R}}(\mathbf{z}^k + \mathbf{s}^{k-1})$$

$$\mathbf{s}^k \leftarrow \mathbf{s}^{k-1} + (\mathbf{z}^k - \mathbf{f}^k)$$

ISTA: $q_k = 1 \Rightarrow O(1/t)$

FISTA: specific $q_k \Rightarrow O(1/t^2)$

ADMM fast practical convergence

# FISTA and ADMM are two popular algorithms for large-scale and nonsmooth optimization

Consider the objective function

$$\mathcal{C}(\mathbf{f}) = \mathcal{D}(\mathbf{f}) + \mathcal{R}(\mathbf{f}) \quad \text{where} \quad \mathcal{D}(\mathbf{f}) \triangleq \frac{1}{2}\|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2$$

Fast iterative shrinkage/thresholding algorithm (FISTA) vs. Alternating direction method of multipliers (ADMM)

$$\mathbf{z}^k \leftarrow \mathbf{s}^{k-1} - \gamma \nabla \mathcal{D}(\mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \text{prox}_{\gamma\mathcal{R}}(\mathbf{z}^k)$$

$$\mathbf{s}^k \leftarrow \mathbf{f}^k + ((q_{k-1} - 1)/q_k)(\mathbf{f}^k - \mathbf{f}^{k-1})$$

$$\mathbf{z}^k \leftarrow \text{prox}_{\gamma\mathcal{D}}(\mathbf{f}^{k-1} - \mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \text{prox}_{\gamma\mathcal{R}}(\mathbf{z}^k + \mathbf{s}^{k-1})$$

$$\mathbf{s}^k \leftarrow \mathbf{s}^{k-1} + (\mathbf{z}^k - \mathbf{f}^k)$$

Question: Which one is computationally more efficient?

# FISTA and ADMM are two popular algorithms for large-scale and nonsmooth optimization

**Consider the objective function**

$$\mathcal{C}(\mathbf{f}) = \mathcal{D}(\mathbf{f}) + \mathcal{R}(\mathbf{f}) \quad \text{where} \quad \mathcal{D}(\mathbf{f}) \triangleq \frac{1}{2}\|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2$$

**Fast iterative shrinkage/thresholding algorithm (FISTA) vs. Alternating direction method of multipliers (ADMM)**

$$\mathbf{z}^k \leftarrow \mathbf{s}^{k-1} - \gamma\nabla\mathcal{D}(\mathbf{s}^{k-1})$$
$$\mathbf{f}^k \leftarrow \text{prox}_{\gamma\mathcal{R}}(\mathbf{z}^k)$$
$$\mathbf{s}^k \leftarrow \mathbf{f}^k + ((q_{k-1} - 1)/q_k)(\mathbf{f}^k - \mathbf{f}^{k-1})$$

$$\mathbf{z}^k \leftarrow \text{prox}_{\gamma\mathcal{D}}(\mathbf{f}^{k-1} - \mathbf{s}^{k-1})$$
$$\mathbf{f}^k \leftarrow \text{prox}_{\gamma\mathcal{R}}(\mathbf{z}^k + \mathbf{s}^{k-1})$$
$$\mathbf{s}^k \leftarrow \mathbf{s}^{k-1} + (\mathbf{z}^k - \mathbf{f}^k)$$

**Per-iteration complexity of ADMM is generally higher**

$$\nabla\mathcal{D}(\mathbf{f}) = \mathbf{H}^\mathsf{T}(\mathbf{H}\mathbf{f} - \mathbf{y})$$

$$\text{prox}_{\gamma\mathcal{D}}(\mathbf{f}) = [\mathbf{I} + \gamma\mathbf{H}^\mathsf{T}\mathbf{H}]^{-1}(\mathbf{f} + \gamma\mathbf{H}^\mathsf{T}\mathbf{y})$$

requires matrix inversion

# To conclude "optimization"

Many imaging problems are ill-posed:
there are infinitely many solutions

$$\mathbf{y} = \mathbf{H}\mathbf{f} + \mathbf{e}$$

Regularization is a strategy to select the
solution that "makes sense"

$$\text{minimize } \mathcal{R}(\mathbf{f}) \text{ subject to } \|\mathbf{H}\mathbf{f} - \mathbf{y}\|_{\ell_2}^2 \leq \sigma^2$$

Classical image regularizers are linear,
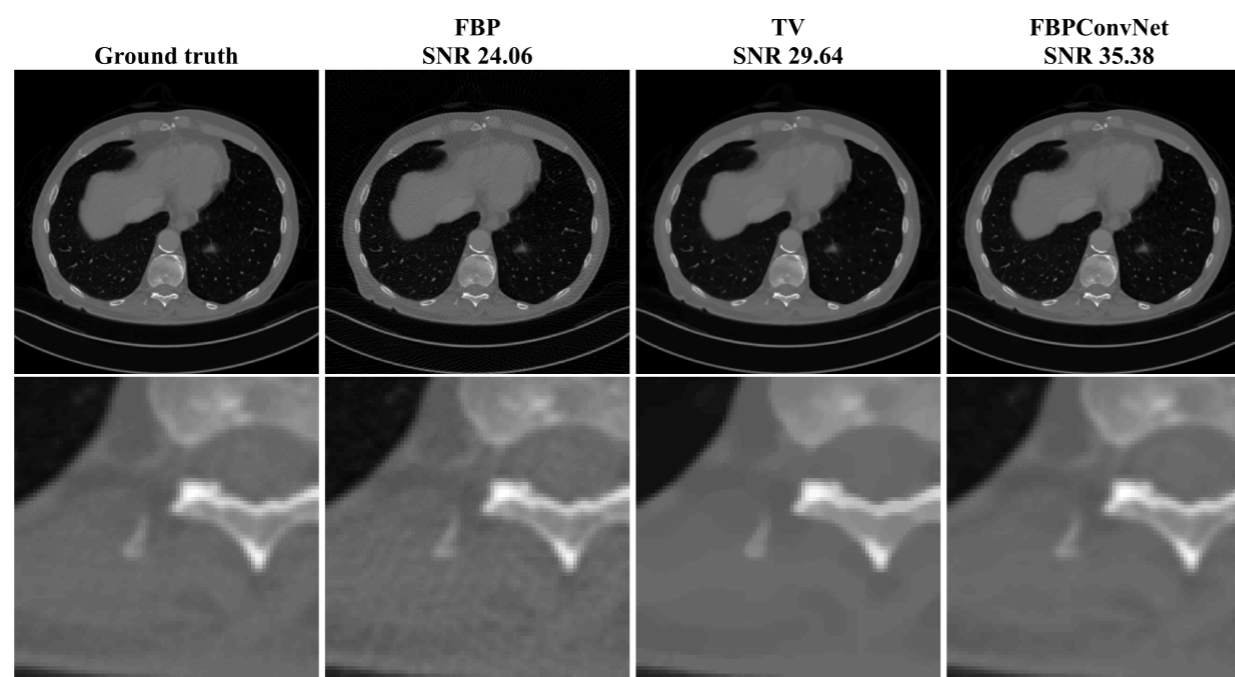but increasingly they are nonlinear

(20th) $\quad \mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \mathcal{R}(\mathbf{f}) = \|\mathbf{D}\mathbf{f}\|_{\ell_1} \quad$ (21st)

# Today we will talk about

◉ Forward models in imaging
Relating the unknowns to the measured data

◉ Notions of ill-posedness and regularization
When measurements are not enough

◉ Optimization at large scales
When analytical solutions are not enough

◉ **Plug-and-Play Priors (PnP) at large scales**
When traditional optimization is not enough

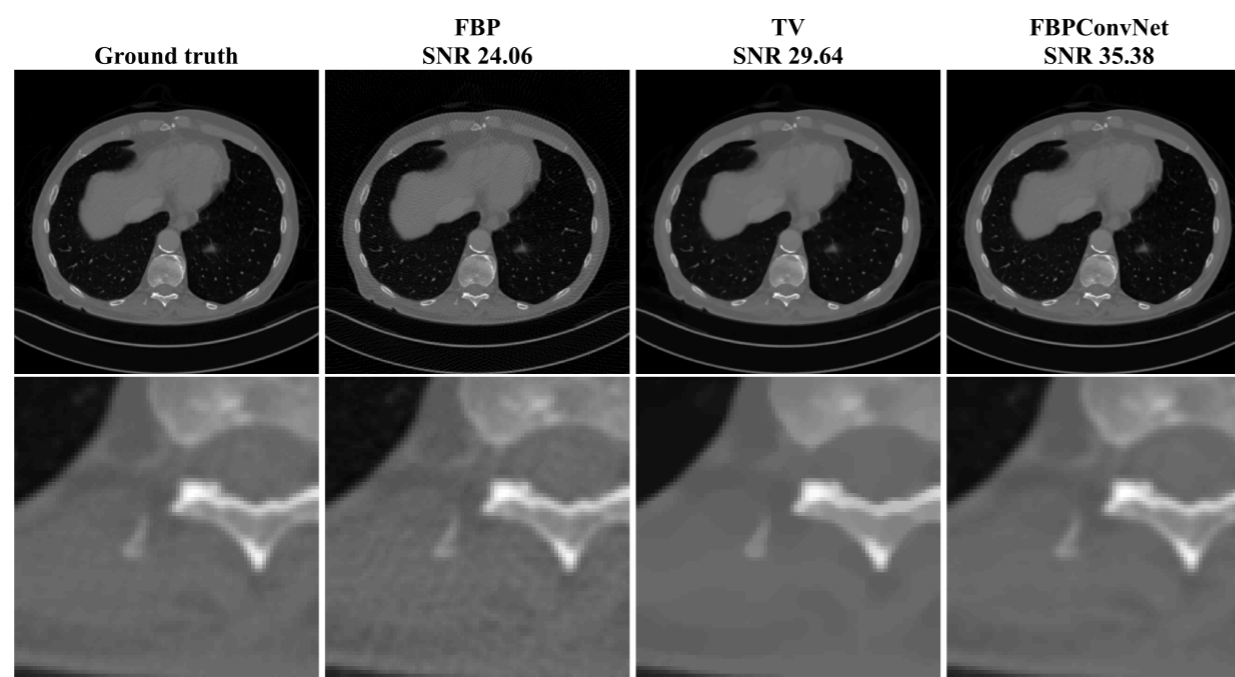# Deep learning is currently getting the best performance for image reconstruction

# Deep learning is currently getting the best performance for image reconstruction



X-Ray CT                    Jin *et al.*, 2016

# Deep learning is currently getting the best performance for image reconstruction



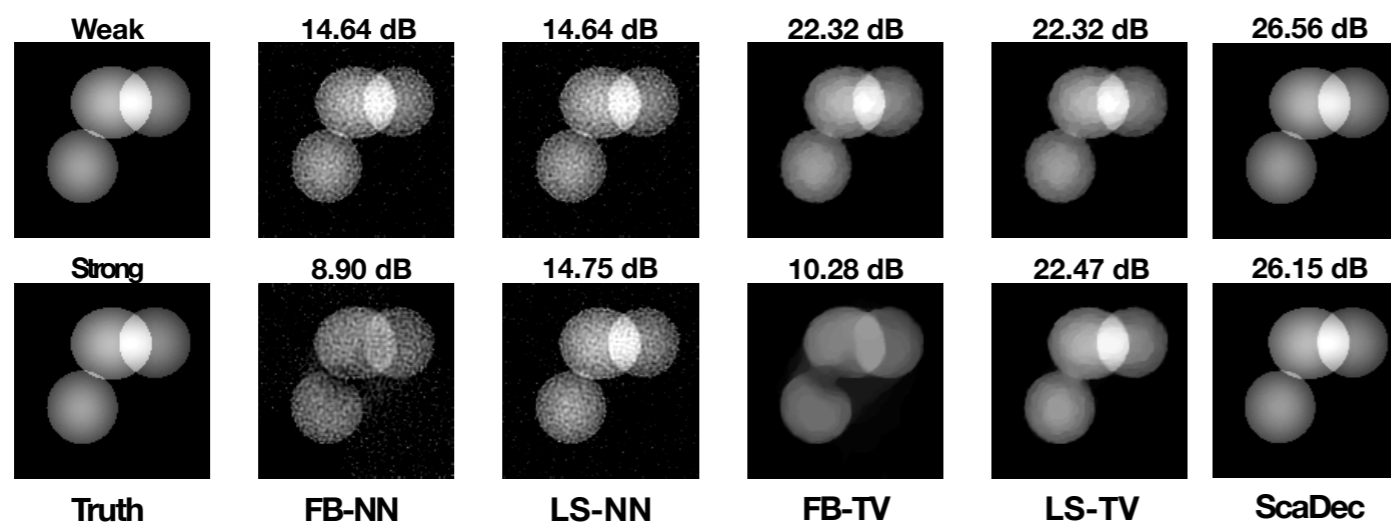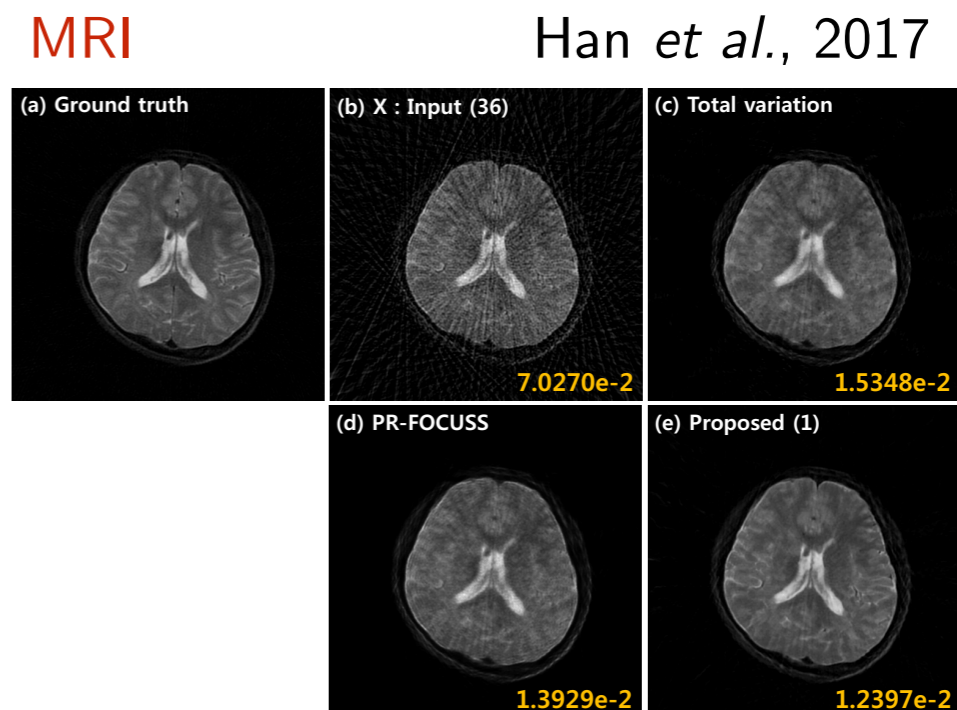X-Ray CT

Jin *et al.*, 2016
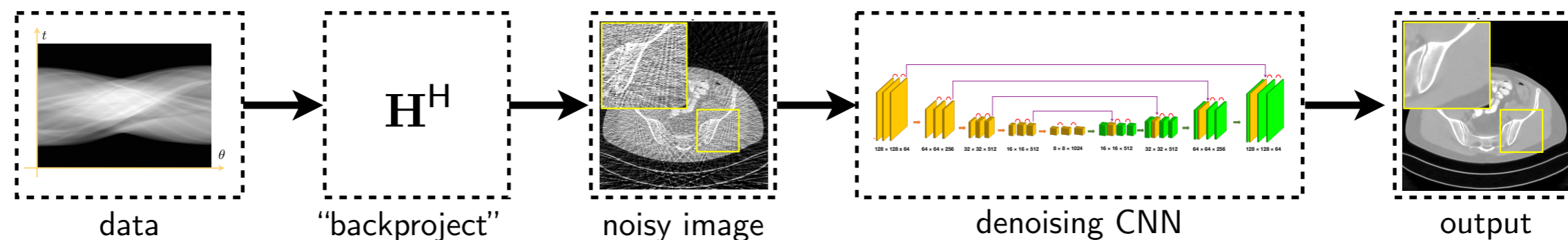
MRI

Han *et al.*, 2017

# Deep learning is currently getting the best performance for image reconstruction



Ground truth | FBP SNR 24.06 | TV SNR 29.64 | FBPConvNet SNR 35.38

X-Ray CT          Jin *et al.*, 2016

MRI          Han *et al.*, 2017

(a) Ground truth | (b) X : Input (36) 7.0270e-2 | (c) Total variation 1.5348e-2
(d) PR-FOCUSS 1.3929e-2 | (e) Proposed (1) 1.2397e-2

Weak    14.64 dB    14.64 dB    22.32 dB    26.56 dB

Strong          26.15 dB

Truth

Diffraction Tomography          Sun *et al.*, 2018

# A well established deep learning pipeline:
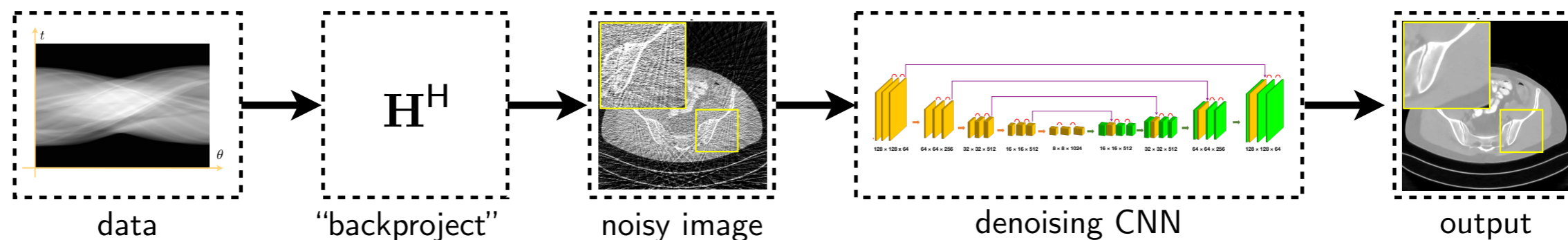# first backproject then denoise with a ConvNet

# A well established deep learning pipeline:
# first backproject then denoise with a ConvNet



Data processing pipeline
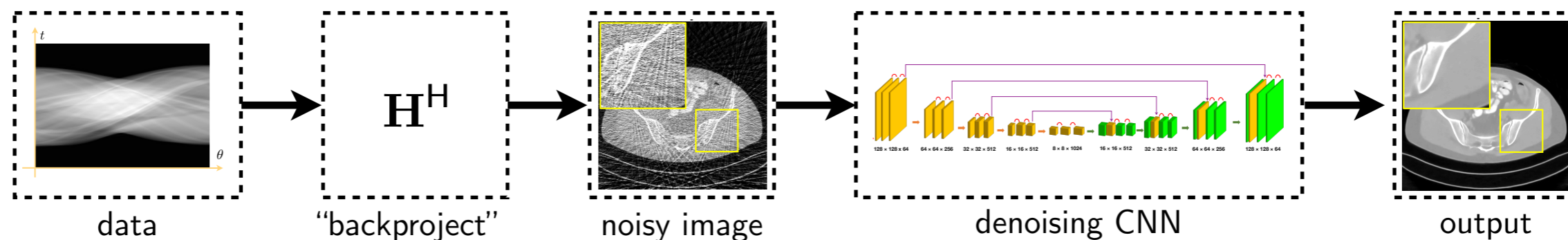
$H^H$

data     "backproject"     noisy image     denoising CNN     output

Ground truth    X : Input    Single-scale learning    Ground truth    Proposed    X : Input    Single-scale learning    Proposed

48 view

96 view

# A well established deep learning pipeline: first backproject then denoise with a ConvNet

**Data processing pipeline**

$H^H$

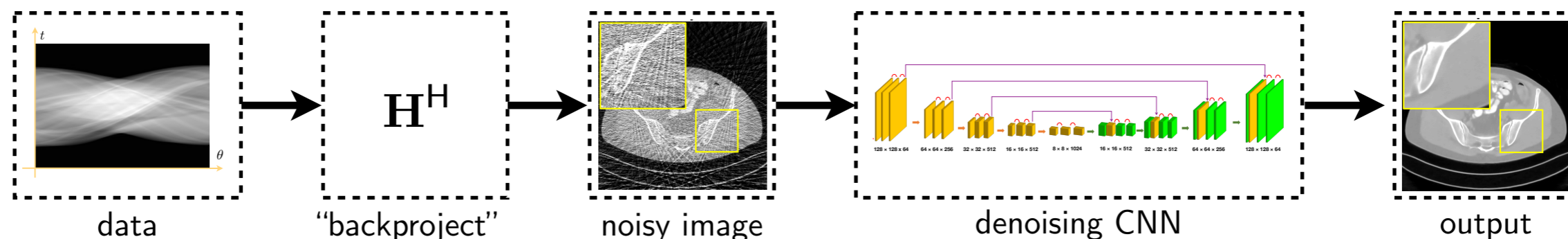data — "backproject" — noisy image — denoising CNN — output

Question: What are some of the key limitations of this approach?

| | Ground truth | X : Input | Single-scale learning | Proposed | Ground truth | X : Input | Single-scale learning | Proposed |

# A well established deep learning pipeline:
# first ~~backproject then denoise~~ with a ConvNet



Data processing pipeline

$H^H$

data     "backproject"     noisy image     denoising CNN     output

**1) Implicit dependance of CNN on the forward model**

*Hard to decouple the individual contributions of D and R*

| Ground truth | X : Input | Single-scale learning | Proposed | Ground truth | X : Input | Single-scale learning | Proposed |

48 view

96 view

# A well established deep learning pipeline:
# first backproject then denoise with a ConvNet

Data processing pipeline



$H^H$

data          "backproject"          noisy image          denoising CNN          output
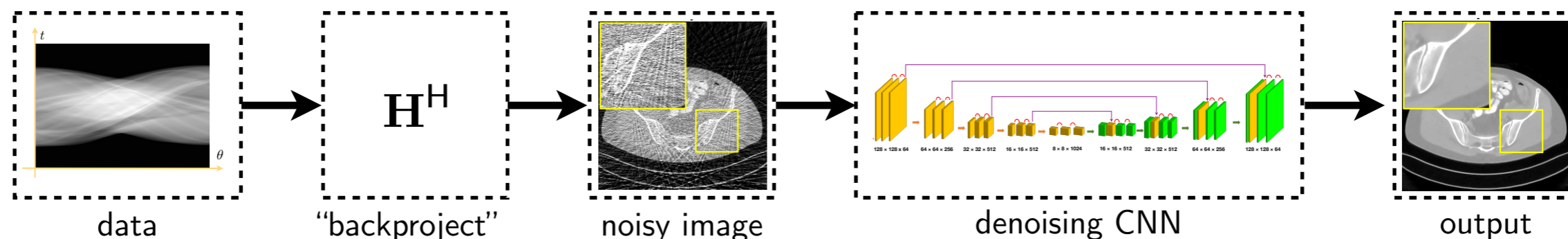
1) Implicit dependance of CNN on the forward model

2) Consistency with the measured data is unclear

No explicit measure of the deviation from the data

| Ground truth | X : Input | Single-scale learning | Proposed | Ground truth | X : Input | Single-scale learning | Proposed |

48 view

96 view

# A well established deep learning pipeline:
# first backproject then denoise with a ConvNet

**Data processing pipeline**



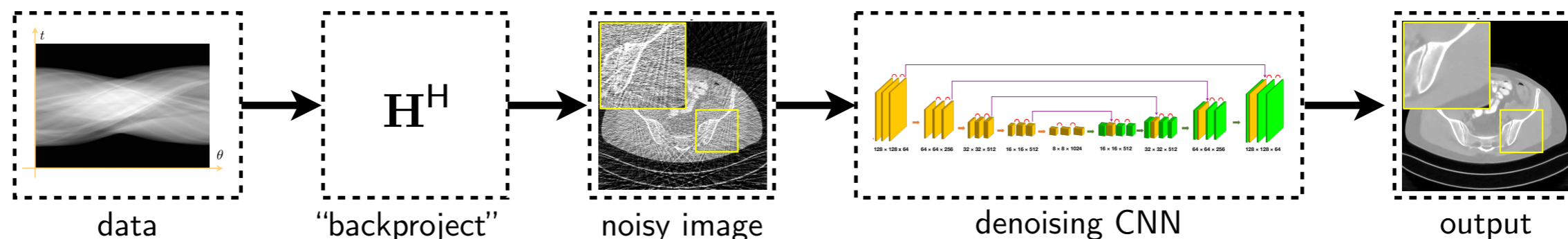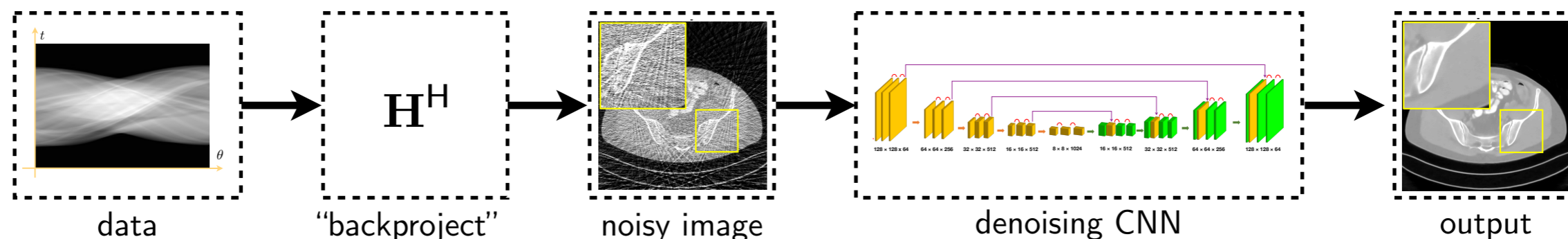data     "backproject"     noisy image     denoising CNN     output

$H^H$

1) Implicit dependance of CNN on the forward model

2) Consistency with the measured data is unclear

3) Difficult to impose non-trivial hard constraints on the image

Example: We absolutely need the image gradient to be smaller than epsilon

Ground truth    X : Input    Single-scale learning    Proposed    Ground truth    X : Input    Single-scale learning    Proposed

48 view

96 view

48 view

96 view

# A well established deep learning pipeline:
## first backproject then denoise with a ConvNet

**Data processing pipeline**



data     "backproject"     noisy image     denoising CNN     output

$H^H$

1) Implicit dependance of CNN on the forward model

2) Consistency with the measured data is unclear

3) Difficult to impose non-trivial hard constraints on the image

4) Not principled: how to select the right architecture?

Variations in the problem are not explicitly linked to model parameters

# A well established deep learning pipeline: first backproject then denoise with a ConvNet
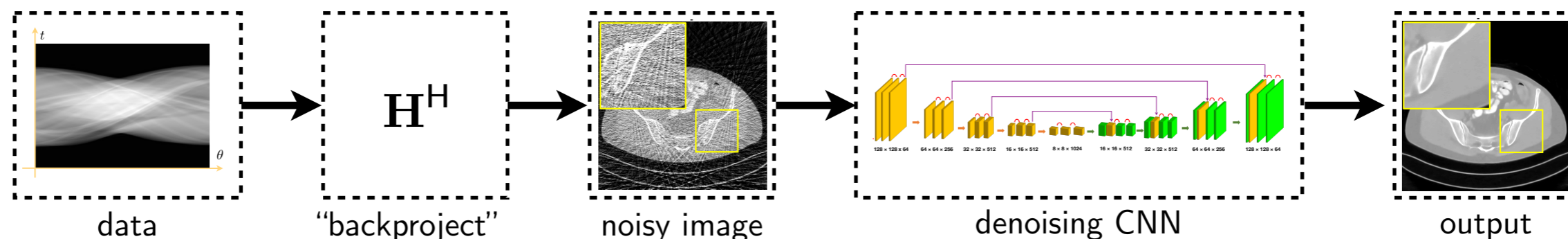
Data processing pipeline

$H^H$

data          "backproject"          noisy image          denoising CNN          output
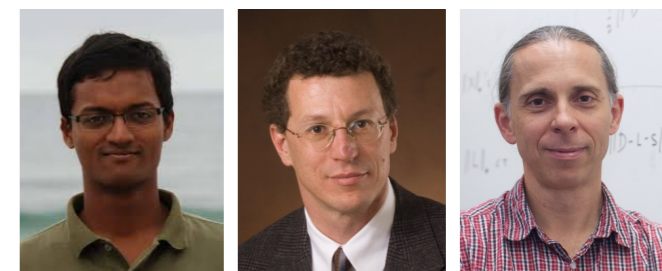
1) Implicit dependance of CNN on the forward model

2) Consistency with the measured data is unclear

3) Difficult to impose non-trivial hard constraints on the image

4) Not principled: how to select the right architecture?
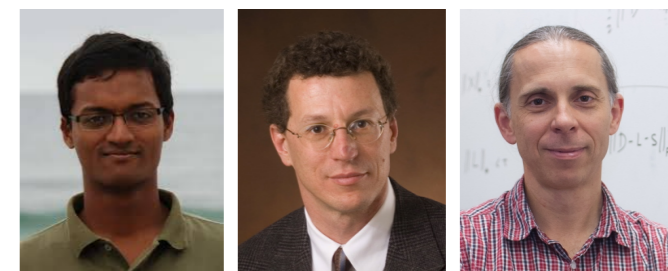
5) Difficult to generalize to nonlinear forward models

What happens if there is no backprojection?

Ground truth | X : Input | Single-scale learning | Proposed

48 view

96 view

# A well established deep learning pipeline:
# first backproject then denoise with a ConvNet

**Data processing pipeline**



$H^H$

data          "backproject"          noisy image          denoising CNN          output

1) Implicit dependance of CNN on the forward model

2) Consistency with the measured data is unclear

3) Difficult to impose non trivial hard constraints on the image

4) Not principled: how to select the right architecture?
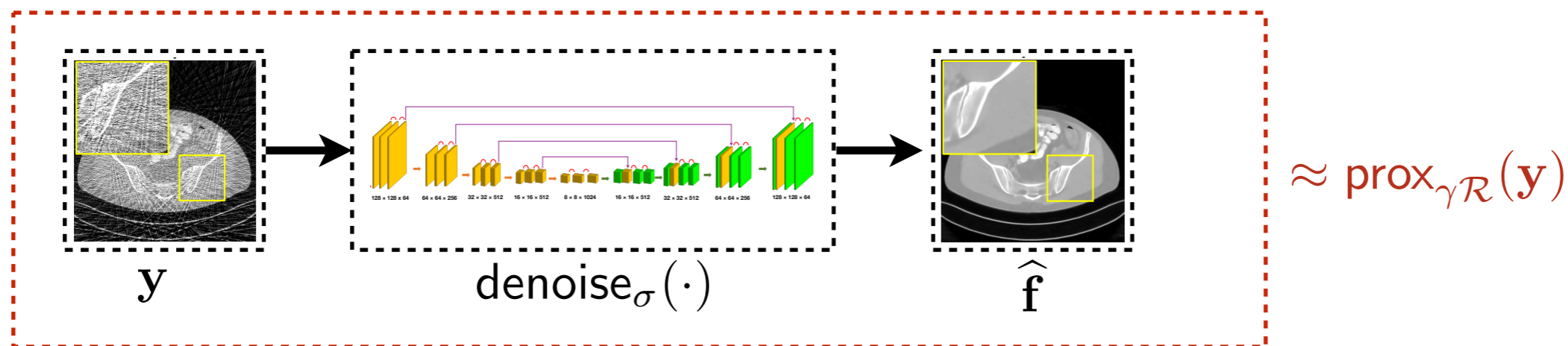
5) Difficult to generalize to nonlinear forward models



| Ground truth | X : Input | Single-scale learning | Proposed | Ground truth | X : Input | Single-scale learning | Proposed |

48 view

96 view

96 view

# Treating the denoising CNN as a proximal operator allows to separate the prior from the forward model

# Treating the denoising CNN as a proximal operator allows to separate the prior from the forward model
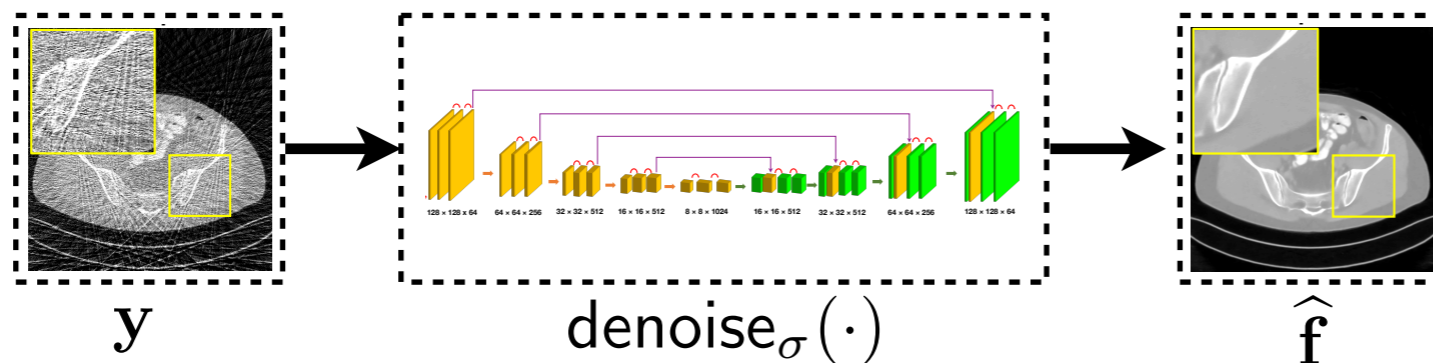
Venkatakrishnan *et al.*, "Plug-and-Play Priors for Model Based Reconstruction," 2013.

# Treating the denoising CNN as a proximal operator allows to separate the prior from the forward model

**Train a CNN to denoise for various noise levels**



$$\approx \mathrm{prox}_{\gamma\mathcal{R}}(\mathbf{y})$$

$\mathbf{y}$  denoise$_{\sigma}(\cdot)$  $\widehat{\mathbf{f}}$



Venkatakrishnan *et al.*, "Plug-and-Play Priors for Model Based Reconstruction," 2013.

# Treating the denoising CNN as a proximal operator allows to separate the prior from the forward model

Train a CNN to denoise for various noise levels



$$\mathbf{y} \qquad \text{denoise}_\sigma(\cdot) \qquad \hat{\mathbf{f}}$$

## Use the trained CNN as a Plug-and-Play Prior (PnP)

$$\mathbf{z}^k \leftarrow \mathbf{s}^{k-1} - \gamma \nabla \mathcal{D}(\mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \text{denoise}_\sigma(\mathbf{z}^k)$$

$$\mathbf{z}^k \leftarrow \text{prox}_{\gamma\mathcal{D}}(\mathbf{f}^{k-1} - \mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \text{denoise}_\sigma(\mathbf{z}^k + \mathbf{s}^{k-1})$$

$$\mathbf{s}^k \leftarrow \mathbf{s}^{k-1} + (\mathbf{z}^k - \mathbf{f}^k)$$

PnP-FISTA

PnP-ADMM

Ground truth   X : Input   Single-scale learning   Proposed   Ground truth   X : Input   Single-scale learning   Proposed

48 view

96 view



Venkatakrishnan *et al.*, "Plug-and-Play Priors for Model Based Reconstruction," 2013.

# Plug-and-Play Priors (PnP) approach has been shown to yield state-of-the-art results

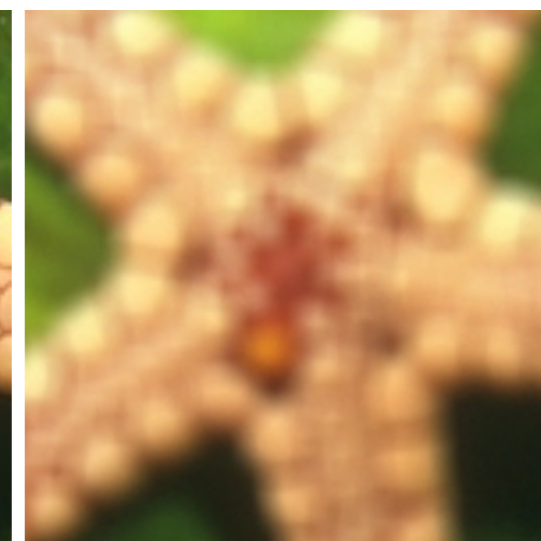# Plug-and-Play Priors (PnP) approach has been shown to yield state-of-the-art results

| Method | Average PSNR (dB) over 10 images |
|--------|-----------------------------------|
| TV | 29.22 |
| IDD-BM3D | 30.92 |
| ASDS-Reg | 30.11 |
| NCSR | 31.09 |
| PnP | **31.33** |

Romano *et al.*, "The Little Engine That Could: Regularization by Denoising," 2017
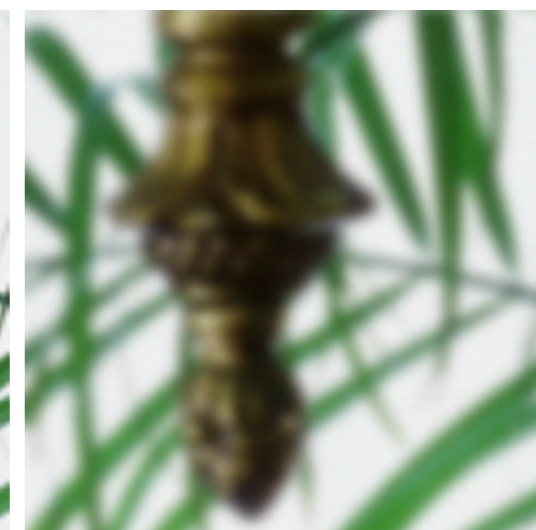
# Plug-and-Play Priors (PnP) approach has been shown to yield state-of-the-art results

| Method | Average PSNR (dB) over 10 images |
|--------|----------------------------------|
| TV | 29.22 |
| IDD-BM3D | 30.92 |
| ASDS-Reg | 30.11 |
| NCSR | 31.09 |
| PnP | **31.33** |



(a) Ground Truth      (b) Input 20.83dB

(d) NCSR 28.39dB      (e) $P^3$-TNRD 28.43dB

Romano *et al.*, "The Little Engine That Could: Regularization by Denoising," 2017

# Plug-and-Play Priors (PnP) approach has been shown to yield state-of-the-art results

| Method | Average PSNR (dB) over 10 images |
|--------|----------------------------------|
| TV | 29.22 |
| IDD-BM3D | 30.92 |
| ASDS-Reg | 30.11 |
| NCSR | 31.09 |
| PnP | **31.33** |



(a) Ground Truth (b) Input 21.40dB

(d) NCSR 30.03dB (e) $P^3$-TNRD 30.36dB

Romano *et al.*, "The Little Engine That Could: Regularization by Denoising," 2017

# Can we say anything about convergence?

# Can we say anything about convergence?

Sreehari *et al.*, "Plug-and-Play Priors for Bright Field Electron Tomography and Sparse Interpolation," 2016

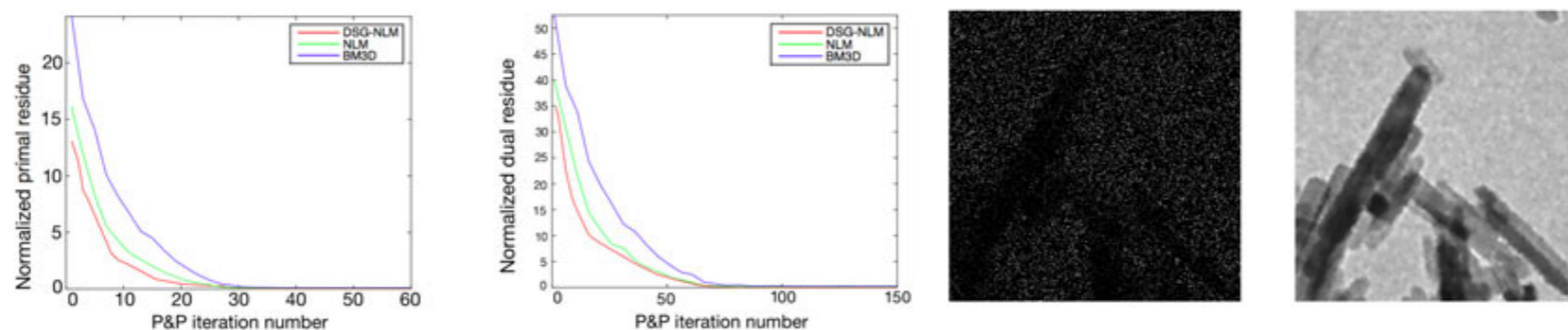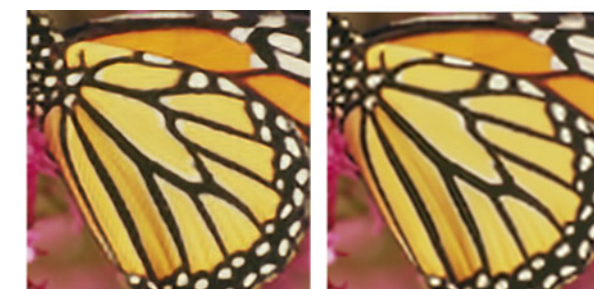Chan *et al.*, "Plug-and-Play ADMM for Image Restoration: Fixed-Point Convergence and Applications," 2016

# Can we say anything about convergence?

Result #1: When $\mathcal{D}(\cdot)$ is convex and $\nabla\text{denoise}_\sigma(\cdot)$ is a symmetric matrix with eigenvalues in $[0, 1]$, then $\text{denoise}_\sigma(\cdot)$ is a proximal operator.

Result #2: When both $\nabla\mathcal{D}(\cdot)$ and $\text{denoise}_\sigma(\cdot)$ are bounded operators, PnP-ADMM with damping converges to a fixed point.

Sreehari *et al.*, "Plug-and-Play Priors for Bright Field Electron Tomography and Sparse Interpolation," 2016

Chan *et al.*, "Plug-and-Play ADMM for Image Restoration: Fixed-Point Convergence and Applications," 2016

Result #1: When $\mathcal{D}(\cdot)$ is convex and $\nabla\text{denoise}_\sigma(\cdot)$ is a symmetric matrix with eigenvalues in $[0, 1]$, then $\text{denoise}_\sigma(\cdot)$ is a proximal operator.



Result #2: When both $\nabla\mathcal{D}(\cdot)$ and $\text{denoise}_\sigma(\cdot)$ are bounded operators, PnP-ADMM with damping converges to a fixed point.

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| DCNN [9] | 20.72 | 21.30 | 18.91 | 21.68 | 16.10 | 23.39 | 22.33 | 22.99 | 22.46 | 20.23 | 21.01 |
| SR [12] | 20.67 | 21.30 | 18.86 | 21.51 | 16.37 | 23.15 | 22.19 | 22.85 | 22.26 | 20.33 | 20.95 |
| SPSR [10] | 20.85 | 21.58 | 19.18 | 21.85 | 16.59 | 23.52 | 22.42 | 23.05 | 22.53 | 20.50 | 21.21 |
| TSE [52] | 20.59 | 21.24 | 18.80 | 21.49 | 16.40 | 23.14 | 22.21 | 22.78 | 22.21 | 20.30 | 20.92 |
| GPR [11] | 21.55 | 22.68 | 19.90 | 22.77 | 17.70 | 24.57 | 23.51 | 24.37 | 23.63 | 21.35 | 22.20 |
| Ours - M | **23.62** | **25.75** | **23.06** | **25.30** | **24.48** | **27.17** | **29.14** | **29.42** | **26.86** | **26.86** | **26.17** |

DCNN          PnP-ADMM

Sreehari *et al.*, "Plug-and-Play Priors for Bright Field Electron Tomography and Sparse Interpolation," 2016

Chan *et al.*, "Plug-and-Play ADMM for Image Restoration: Fixed-Point Convergence and Applications," 2016

# Can we say anything about convergence?

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for
Regularized Image Reconstruction," 2018
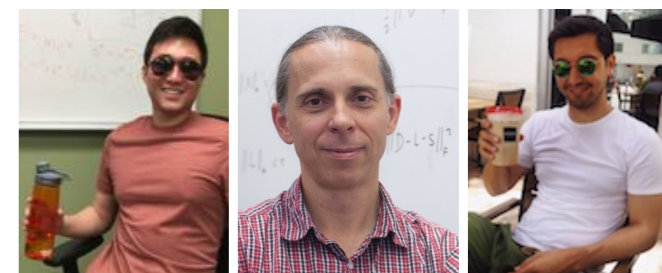
# Can we say anything about convergence?

**Useful definitions**

$$P(\mathbf{f}) \triangleq \mathrm{denoise}_\sigma(\mathbf{f} - \gamma \nabla \mathcal{D}(\mathbf{f}))$$

gradient-denoiser operator

$$\mathrm{fix}(P) \triangleq \{\mathbf{f} \in \mathbb{R}^n : \mathbf{f} = P(\mathbf{f})\}$$

its of fixed points

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# Can we say anything about convergence?

**Useful definitions**

$$\mathsf{P}(\mathbf{f}) \triangleq \mathsf{denoise}_\sigma(\mathbf{f} - \gamma \nabla \mathcal{D}(\mathbf{f})) \qquad \mathsf{fix}(\mathsf{P}) \triangleq \{\mathbf{f} \in \mathbb{R}^n : \mathbf{f} = \mathsf{P}(\mathbf{f})\}$$

#1: Let $\mathsf{denoise}_\sigma(\cdot) = \mathsf{prox}_{\gamma \mathcal{R}}(\cdot)$. Then, $\mathbf{f}^* \in \mathsf{fix}(\mathsf{P})$ iff it minimizes $\mathcal{C} = \mathcal{D} + \mathcal{R}$

#2: Run PnP-ISTA with a nonexpansive denoiser for $t \geq 1$ iterations. Then

$$\min_{k \in \{1, \dots, t\}} \left\{ \|\mathbf{f}^{k-1} - \mathsf{P}(\mathbf{f}^{k-1})\|_{\ell_2}^2 \right\} = O(1/t)$$

#3: For nonexpansive denoisers, fixed points of PnP-ADMM coincide with $\mathsf{fix}(\mathsf{P})$

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# PnP-SGD is an online extension useful when dealing with a large number of measurements

# PnP-SGD is an online extension useful when dealing with a large number of measurements

**Consider the following data-fidelity term**

$$\mathcal{D}(\mathbf{f}) = \frac{1}{2I} \sum_{i=1}^{I} \|\mathbf{y}_i - \mathbf{H}_i \mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \nabla \mathcal{D}(\mathbf{f}) = \frac{1}{I} \sum_{i=1}^{I} \mathbf{H}_i^\mathsf{T} (\mathbf{H}_i \mathbf{f} - \mathbf{y})$$

cost of computing the gradient is
liner in the number of measurements

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# PnP-SGD is an online extension useful when dealing with a large number of measurements

**Consider the following data-fidelity term**

$$\mathcal{D}(\mathbf{f}) = \frac{1}{2I} \sum_{i=1}^{I} \|\mathbf{y}_i - \mathbf{H}_i \mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \nabla\mathcal{D}(\mathbf{f}) = \frac{1}{I} \sum_{i=1}^{I} \mathbf{H}_i^\mathsf{T}(\mathbf{H}_i \mathbf{f} - \mathbf{y})$$

**PnP-SGD can accelerate imaging by parallelizing the processing of each data item**

$$\hat{\nabla}\mathcal{D}(\mathbf{s}^{k-1}) \leftarrow \mathsf{minibatchGradient}(\mathbf{s}^{k-1}, B)$$

$$\mathbf{z}^k \leftarrow \mathbf{s}^{k-1} - \gamma\hat{\nabla}\mathcal{D}(\mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \mathsf{denoise}_\sigma(\mathbf{z}^k)$$

$$\mathbf{s}^k \leftarrow \mathbf{f}^k + ((q_{k-1} - 1)/q_k)(\mathbf{f}^k - \mathbf{f}^{k-1})$$

use only *B* measurements per iteration instead of *I*

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# PnP-SGD is an online extension useful when dealing with a large number of measurements

**Consider the following data-fidelity term**

$$\mathcal{D}(\mathbf{f}) = \frac{1}{2I} \sum_{i=1}^{I} \|\mathbf{y}_i - \mathbf{H}_i\mathbf{f}\|_{\ell_2}^2 \quad \Rightarrow \quad \nabla\mathcal{D}(\mathbf{f}) = \frac{1}{I} \sum_{i=1}^{I} \mathbf{H}_i^\mathsf{T}(\mathbf{H}_i\mathbf{f} - \mathbf{y})$$

**PnP-SGD can accelerate imaging by parallelizing the processing of each data item**

$$\hat{\nabla}\mathcal{D}(\mathbf{s}^{k-1}) \leftarrow \mathsf{minibatchGradient}(\mathbf{s}^{k-1}, B)$$

$$\mathbf{z}^k \leftarrow \mathbf{s}^{k-1} - \gamma\hat{\nabla}\mathcal{D}(\mathbf{s}^{k-1})$$

$$\mathbf{f}^k \leftarrow \mathsf{denoise}_\sigma(\mathbf{z}^k)$$

$$\mathbf{s}^k \leftarrow \mathbf{f}^k + ((q_{k-1} - 1)/q_k)(\mathbf{f}^k - \mathbf{f}^{k-1})$$

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# PnP-SGD converges to the same set of fixed points as batch PnP algorithms

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# PnP-SGD converges to the same set of fixed points as batch PnP algorithms

#4: Run PnP-SGD for $t \geq 1$ iterations under some mild assumptions. Then
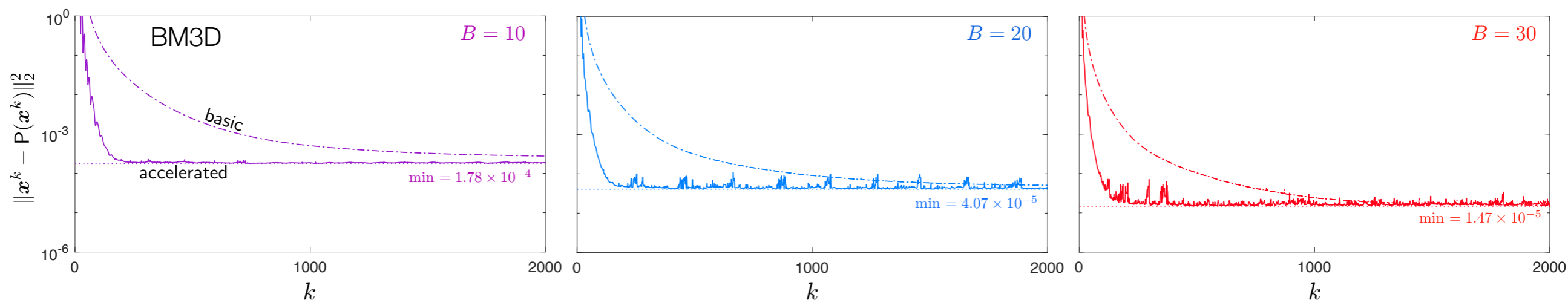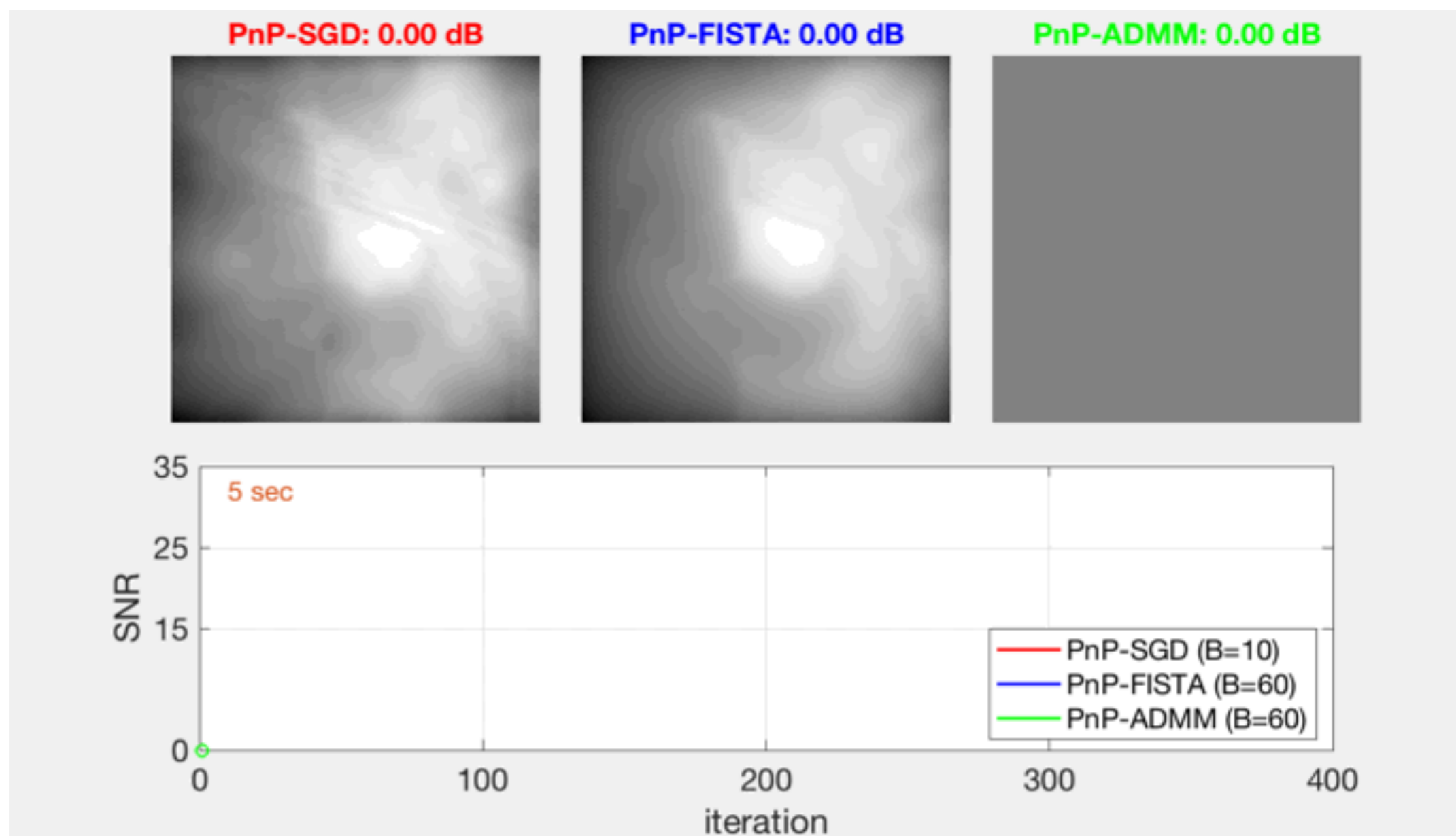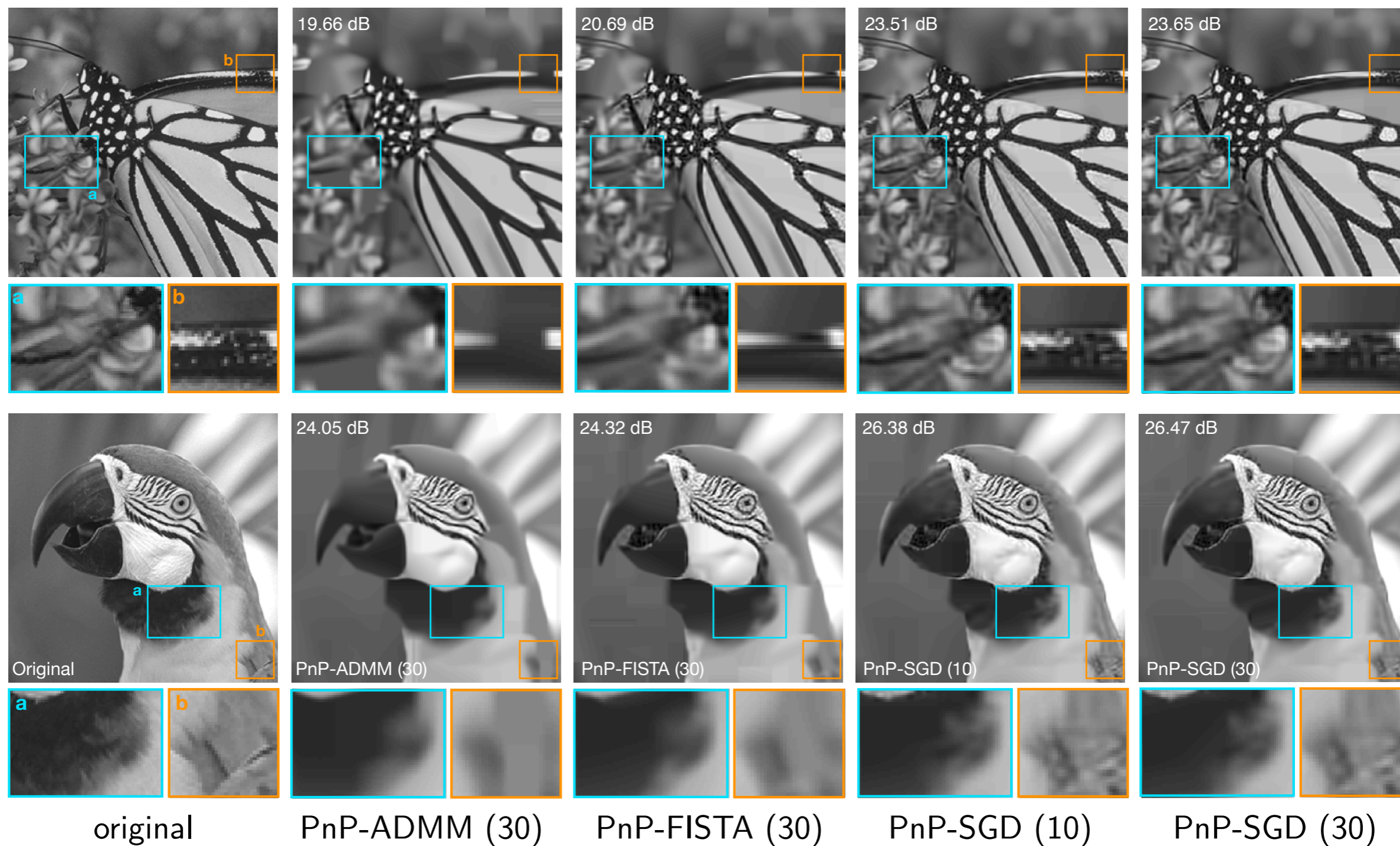
$$\mathbb{E}\left[\min_{k \in \{1,\ldots,t\}} \left\{ \|\mathbf{f}^{k-1} - \mathsf{P}(\mathbf{f}^{k-1})\|_{\ell_2}^2 \right\}\right] \leq C\left[\frac{\gamma^2 \nu^2}{B} + \frac{2\gamma\nu}{\sqrt{B}}\|\mathbf{f}^0 - \mathbf{f}^*\|_{\ell_2} + \frac{\|\mathbf{f}^0 - \mathbf{f}^*\|_{\ell_2}^2}{t}\right]$$

Convergence in expectation. C is a constant. Note the case when B = t

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

$$D(\boldsymbol{x}^k)$$



Sun, Wohlberg, Kamilov, "An Online Plug-ar
Regularized Image Reconstruction," 2018

# For many measurements PnP-SGD converges faster than batch algorithms

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# For the same measurement budget, PnP-SGD gets much higher quality results



|  | original | PnP-ADMM (30) | PnP-FISTA (30) | PnP-SGD (10) | PnP-SGD (30) |

Sun, Wohlberg, Kamilov, "An Online Plug-and-Play Algorithm for Regularized Image Reconstruction," 2018

# Conclusion

Image reconstruction is a fascinating research area that brings together physics, signal processing, nonlinear optimization, and machine learning

We are increasingly reliant on implicit regularization using nonlinear operators, such as deep neural networks or nonlinear filters

Plug-In SGD is a theoretically sound algorithm that can regularize at large-scales using nonlinear operators
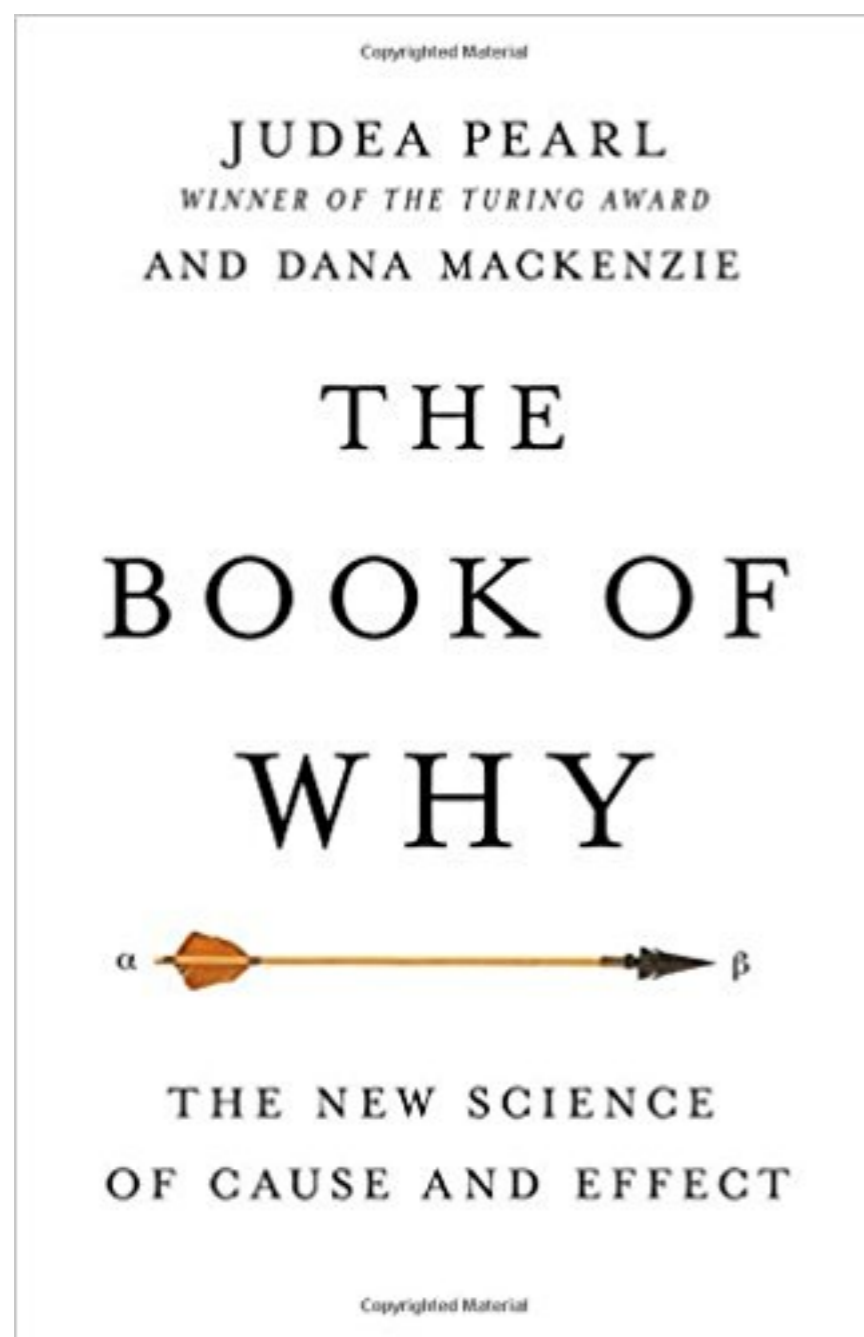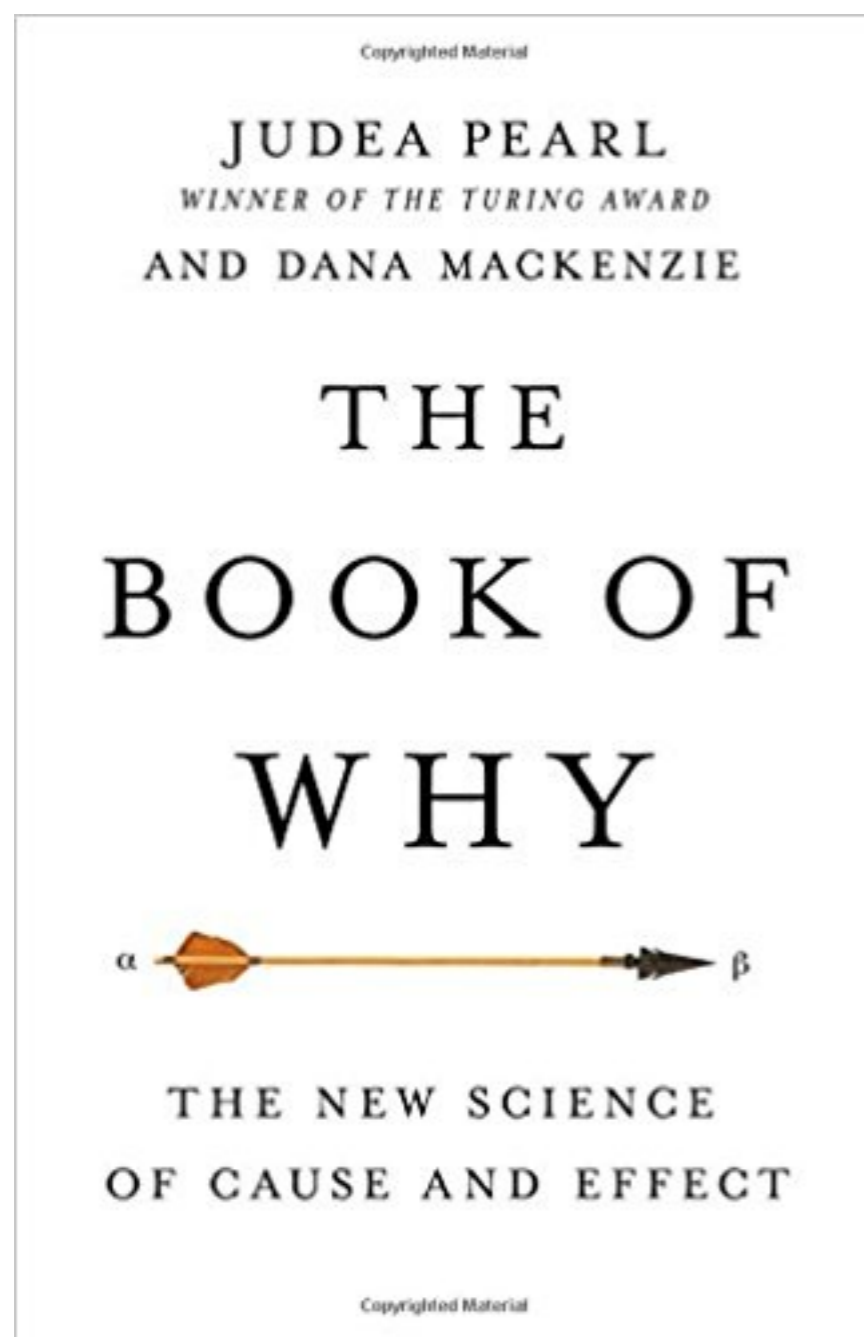
**CONTACT INFO**

Prof. Ulugbek Kamilov
Computational Imaging Group (CIG)
Washington University in St. Louis
Email: kamilov@wustl.edu
Web: http://cigroup.wustl.edu
Twitter: @wustlcig

# Judea Pearl won the Turing Award in 2011 for fundamental contributions to artificial intelligence

# Judea Pearl won the Turing Award in 2011 for fundamental contributions to artificial intelligence





Judea Pearl

# Judea Pearl won the Turing Award in 2011 for fundamental contributions to artificial intelligence



We live in an era that presumes Big Data to be the solution to all our problems (…) But I hope with this book to convince you that data are profoundly dumb. Data can tell you that the people who took a medicine recovered faster than those who did not take it, but they can't tell you why.



Judea Pearl

# Judea Pearl won the Turing Award in 2011 for fundamental contributions to artificial intelligence

The belief that data can tell the full story is a misconception. To produce truly useful insights, data must be combined with models that infuse what we know about the problem.

*Judea Pearl*