

# Bayesian inversion for tomography through machine learning

Ozan Öktem

Department of Mathematics, KTH - Royal Institute of Technology, Stockholm  
The Alan Turing Institute, London

April 3, 2019

Imaging and machine learning  
Mathematics of Imaging Workshop #3  
Institut Henri Poincaré, Paris



# Overview

- Tomography
- Theory for Bayesian inversion
- Deep direct estimation
- Task adapted reconstruction
- Deep posterior sampling

Tomography



# Mathematical model

- Radiative transport equation: Equation for directional intensity of particles
  - Loses energy (absorption)
  - Gains energy (emission)
  - Redistributions energy (scattering)
- Ray transform: Ignore scattering and emission

$$\mathcal{A}(x)(\ell) = \int_{\ell} x(s) \, ds \quad \text{for } \ell = \text{line in } \mathbb{R}^n$$

- $x: \mathbb{R}^n \rightarrow \mathbb{R}$  is the linear attenuation coefficient.
- $\mathcal{A}: \{\text{functions on } \mathbb{R}^n\} \rightarrow \{\text{functions on lines in } \mathbb{R}^n\}$
- Tomographic imaging
  - Acquisition geometry: Set of lines  $\ell$  from which data is sampled.
  - Data:  $y(\ell) = \mathcal{A}(x)(\ell) + \text{'noise'}$  for lines  $\ell$  in acquisition geometry
  - Task: Recover  $x$  (model parameter)

# Mathematical model

- Radiative transport equation: Equation for directional intensity of particles
  - Loses energy (absorption)
  - Gains energy (emission)
  - Redistributions energy (scattering)
- Ray transform: Ignore scattering and emission

$$\mathcal{A}(x)(\ell) = \int_{\ell} x(s) \, ds \quad \text{for } \ell = \text{line in } \mathbb{R}^n$$

- $x: \mathbb{R}^n \rightarrow \mathbb{R}$  is the linear attenuation coefficient.
- $\mathcal{A}: \{\text{functions on } \mathbb{R}^n\} \rightarrow \{\text{functions on lines in } \mathbb{R}^n\}$
- Tomographic imaging
  - Acquisition geometry: Set of lines  $\ell$  from which data is sampled.
  - Data:  $y(\ell) = \mathcal{A}(x)(\ell) + \text{'noise'}$  for lines  $\ell$  in acquisition geometry
  - Task: Recover  $x$  (model parameter)

## Theory for Bayesian inversion

# Inverse problems

Functional analytic vs. statistical viewpoints

## Inverse problem (Frequentist/Functional analytic)

Model:  $\mathbf{y} = \mathcal{A}(x^*) + \mathbf{e}$

- Forward model:  $\mathcal{A}: X \rightarrow Y$  deterministic model for data.
- $X$  = possible model parameters,  $Y$  = possible data.
- Observation noise:  $\mathbf{e} \sim \pi_{\text{noise}}$  known distribution
- $x^* \in X$  unknown

Measured data:  $y \in Y$  single sample of  $\mathbf{y}$ .

Task: Given data  $y$ , recover the **model parameter**  $x^*$ .

# Inverse problems

Functional analytic vs. statistical viewpoints

## Inverse problem (Bayesian)

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

- Forward model:  $\mathcal{A}: X \rightarrow Y$  deterministic model for data.
- $X$  = possible model parameters,  $Y$  = possible data.
- Observation noise:  $\mathbf{e} \sim \pi_{\text{noise}}$  known distribution
- $x^* \in X$  unknown
- Prior:  $\mathbf{x} \sim \pi_{\text{prior}}$ .

Measured data:  $y \in Y$  single sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

Task: Given data  $y$ , recover the **posterior distribution** of  $(\mathbf{x} \mid \mathbf{y} = y)$ .

# Bayes' theorem

- Bayes' theorem:

$$\frac{d\pi_{\text{post}}}{d\pi_{\text{prior}}}(x | y) = \frac{1}{Z(y)} \pi_{\text{data}}(y | x)$$

- Posterior:  $\pi_{\text{post}}(x | y)$  = probability of model parameter  $x$  given data  $y$ .
- Data likelihood:  $\pi_{\text{data}}(y | x)$  = probability of data  $y$  given model parameter  $x$ .
- Prior:  $\pi_{\text{prior}}(x)$  = probability of model parameter  $x$ .
- Normalisation:  $Z(y)$  is the probability of data.

# Bayes' theorem

- Bayes' theorem:

$$\frac{d\pi_{\text{post}}}{d\pi_{\text{prior}}}(x | y) = \frac{1}{Z(y)} \pi_{\text{data}}(y | x)$$

- Posterior:  $\pi_{\text{post}}(x | y)$  = probability of model parameter  $x$  given data  $y$ .
- Data likelihood:  $\pi_{\text{data}}(y | x)$  = probability of data  $y$  given model parameter  $x$ .
- Prior:  $\pi_{\text{prior}}(x)$  = probability of model parameter  $x$ .
- Normalisation:  $Z(y)$  is the probability of data.

- Challenges:

- Posterior  $\Leftarrow$  solution to inverse problem

# Bayes' theorem

- Bayes' theorem:

$$\frac{d\pi_{\text{post}}}{d\pi_{\text{prior}}}(x | y) = \frac{1}{Z(y)} \pi_{\text{data}}(y | x)$$

- Posterior:  $\pi_{\text{post}}(x | y)$  = probability of model parameter  $x$  given data  $y$ .
- Data likelihood:  $\pi_{\text{data}}(y | x)$  = probability of data  $y$  given model parameter  $x$ .
- Prior:  $\pi_{\text{prior}}(x)$  = probability of model parameter  $x$ .
- Normalisation:  $Z(y)$  is the probability of data.

- Challenges:

- Posterior  $\Leftarrow$  solution to inverse problem
- Data likelihood  $\Leftarrow$  physics model for generation of data

$$\pi_{\text{data}}(y | x) = \pi_{\text{noise}}(\mathcal{A}(x) - y) \text{ if } \mathbf{e} \sim \pi_{\text{noise}} \text{ independent of } \mathbf{x}.$$

# Bayes' theorem

- Bayes' theorem:

$$\frac{d\pi_{\text{post}}}{d\pi_{\text{prior}}}(x | y) = \frac{1}{Z(y)} \pi_{\text{data}}(y | x)$$

- Posterior:  $\pi_{\text{post}}(x | y)$  = probability of model parameter  $x$  given data  $y$ .
- Data likelihood:  $\pi_{\text{data}}(y | x)$  = probability of data  $y$  given model parameter  $x$ .
- Prior:  $\pi_{\text{prior}}(x)$  = probability of model parameter  $x$ .
- Normalisation:  $Z(y)$  is the probability of data.

- Challenges:

- Posterior  $\Leftarrow$  solution to inverse problem
- Data likelihood  $\Leftarrow$  physics model for generation of data
- Prior  $\Leftarrow$  can be handcrafted or learned

# Bayes' theorem

- Bayes' theorem:

$$\frac{d\pi_{\text{post}}}{d\pi_{\text{prior}}}(x | y) = \frac{1}{Z(y)} \pi_{\text{data}}(y | x)$$

- Posterior:  $\pi_{\text{post}}(x | y)$  = probability of model parameter  $x$  given data  $y$ .
- Data likelihood:  $\pi_{\text{data}}(y | x)$  = probability of data  $y$  given model parameter  $x$ .
- Prior:  $\pi_{\text{prior}}(x)$  = probability of model parameter  $x$ .
- **Normalisation**:  $Z(y)$  is the probability of data.

- Challenges:

- Posterior  $\iff$  solution to inverse problem
- Data likelihood  $\iff$  physics model for generation of data
- Prior  $\iff$  can be handcrafted or learned
- **Normalisation**  $\iff$  major issue if  $Z(y)$  needs to be computed

## Regularisation theory for Bayesian inversion

**Bayesian inversion:** Recover (posterior) distribution for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $y$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$  with  $x^*$  unknown and  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$ .

- **Existence:** Regular conditional distributions for  $(\mathbf{x} \mid \mathbf{y} = y)$  and  $(\mathbf{y} \mid \mathbf{x} = x)$  exists whenever  $X$  and  $Y$  are complete and separable metric spaces.

# Regularisation theory for Bayesian inversion

**Bayesian inversion:** Recover (posterior) distribution for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $y$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$  with  $x^*$  unknown and  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$ .

- **Existence:** Regular conditional distributions for  $(\mathbf{x} \mid \mathbf{y} = y)$  and  $(\mathbf{y} \mid \mathbf{x} = x)$  exists whenever  $X$  and  $Y$  are complete and separable metric spaces.
- **Stability:** If negative log likelihood of data is locally Hölder continuous, then:
  - Posterior is Lipschitz in the data w.r.t. the Hellinger metric
  - Posterior moments, e.g., mean and covariance, are continuous w.r.t. data

⇒ Statistical formulation regularises.

# Regularisation theory for Bayesian inversion

**Bayesian inversion:** Recover (posterior) distribution for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $y$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$  with  $x^*$  unknown and  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$ .

- **Existence:** Regular conditional distributions for  $(\mathbf{x} \mid \mathbf{y} = y)$  and  $(\mathbf{y} \mid \mathbf{x} = x)$  exists whenever  $X$  and  $Y$  are complete and separable metric spaces.
- **Stability:** If negative log likelihood of data is locally Hölder continuous, then:
  - Posterior is Lipschitz in the data w.r.t. the Hellinger metric
  - Posterior moments, e.g., mean and covariance, are continuous w.r.t. data $\implies$  Statistical formulation regularises.
- **Convergence:**  $\pi_{\text{post}}(\mathbf{x} \mid y)$  concentrates to  $x^*$  as  $\mathbf{e} \rightarrow 0$  (posterior consistency).

# Regularisation theory for Bayesian inversion

**Bayesian inversion:** Recover (posterior) distribution for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $y$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$  with  $x^*$  unknown and  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$ .

- **Existence:** Regular conditional distributions for  $(\mathbf{x} \mid \mathbf{y} = y)$  and  $(\mathbf{y} \mid \mathbf{x} = x)$  exists whenever  $X$  and  $Y$  are complete and separable metric spaces.
- **Stability:** If negative log likelihood of data is locally Hölder continuous, then:
  - Posterior is Lipschitz in the data w.r.t. the Hellinger metric
  - Posterior moments, e.g., mean and covariance, are continuous w.r.t. data $\implies$  Statistical formulation regularises.
- **Convergence:**  $\pi_{\text{post}}(\mathbf{x} \mid y)$  concentrates to  $x^*$  as  $\mathbf{e} \rightarrow 0$  (posterior consistency).
- **Asymptotic analysis:** Posterior consistency holds for many choices of priors.
  - Contraction rates: Determine the rate at which  $\pi_{\text{post}}(\mathbf{x} \mid y)$  concentrates to  $x^*$ .
  - Microscopic fluctuations: Characterise  $\pi_{\text{post}}(\mathbf{x} \mid y)$  asymptotically near  $x^*$ .

Results represent analysis of properties independent of the choice of prior.

# Regularisation theory for Bayesian inversion

Functional analytic regularisation	Bayesian inversion
Existence	Existence
Stability	Stability
Convergence	Posterior consistency
Convergence rates	Contraction rates
??	Microscopic fluctuations
Stability/error estimates	??

## Example of point estimators I

- Maximum likelihood: Maximise the data likelihood, i.e. for given  $y$  solve

$$\arg \max_x \pi_{\text{data}}(y | x) \iff \arg \min_x \left[ -\log \pi_{\text{data}}(y | x) \right].$$

- + No high-dimensional integration
- + Computationally feasible for large-scale problems.
- Does not regularise, not suitable for ill-posed problems.

## Example of point estimators I

- **Maximum likelihood:** Maximise the data likelihood, i.e. for given  $y$  solve

$$\arg \max_x \pi_{\text{data}}(y | x) \iff \arg \min_x \left[ -\log \pi_{\text{data}}(y | x) \right].$$

- + No high-dimensional integration
- + Computationally feasible for large-scale problems.
- Does not regularise, not suitable for ill-posed problems.

- **Maximum a posteriori:** Maximise the posterior, for given  $y$  solve

$$\arg \max_x \pi_{\text{post}}(x | y) \iff \arg \min_x \left[ -\log \pi_{\text{data}}(y | x) - \log \pi_{\text{prior}}(x) \right].$$

- + No high-dimensional integration
- + Regularises in most cases, suitable for most ill-posed problems.
- Need to specify prior.
- Optimisation can be non-smooth  $\implies$  unfeasible for large-scale problems.

## Example of point estimators II

- Conditional mean: Model parameter  $\hat{x}$  is the conditional mean, i.e. for given  $y$

$$\hat{x} := \mathbb{E}[x | y = y] = \int_X x \pi_{\text{post}}(x | y) dx.$$

- + Regularises, suitable for ill-posed problems.
- Need to specify posterior  $\implies$  specify prior
- Involves high-dimensional integration  $\implies$  unfeasible for mid-scale problems.

## Example of point estimators II

- **Conditional mean:** Model parameter  $\hat{x}$  is the conditional mean, i.e. for given  $y$

$$\hat{x} := \mathbb{E}[\mathbf{x} \mid \mathbf{y} = y] = \int_X x \pi_{\text{post}}(x \mid y) dx.$$

- + Regularises, suitable for ill-posed problems.
- Need to specify posterior  $\Rightarrow$  specify prior
- Involves high-dimensional integration  $\Rightarrow$  unfeasible for mid-scale problems.

- **Bayes estimator:** Minimise expected loss w.r.t.  $\ell_X: X \times X \rightarrow \mathbb{R}$ , i.e. for given  $y$

$$\hat{x} := \widehat{\mathcal{R}}(y) \quad \text{where} \quad \widehat{\mathcal{R}} \in \arg \min_{\mathcal{R}: Y \rightarrow X} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mu} [\ell_X(\mathcal{R}(\mathbf{y}), \mathbf{x})].$$

- + Regularises, suitable for ill-posed problems.
- + Suitable for supervised learning.
- + Equivalent to conditional mean when loss is squared 2-norm.
- Need to specify joint distribution  $\Rightarrow$  specify prior and computing normalisation.
- Involves high-dimensional integration and optimisation  
 $\Rightarrow$  unfeasible for small-scale problems.

## Additive Gaussian noise and Gibbs type of prior

- Inverse problem:  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$
- Data likelihood:  $\mathbf{e} \sim \text{Normal}(0, \sigma^2)$ 
  - $\pi_{\text{noise}}(\mathbf{e}) \propto \exp\left(-\frac{1}{\sigma^2} \|\mathbf{e}\|^2\right)$
  - $\pi_{\text{data}}(\mathbf{y} \mid \mathbf{x}) = \pi_{\text{noise}}(\mathcal{A}(\mathbf{x}) - \mathbf{y}) \implies -\log \pi_{\text{data}}(\mathbf{y} \mid \mathbf{x}) \propto \frac{1}{\sigma^2} \|\mathcal{A}(\mathbf{x}) - \mathbf{y}\|^2$

# Additive Gaussian noise and Gibbs type of prior

- Inverse problem:  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$
- Data likelihood:  $\mathbf{e} \sim \text{Normal}(0, \sigma^2)$ 
  - $\pi_{\text{noise}}(\mathbf{e}) \propto \exp\left(-\frac{1}{\sigma^2} \|\mathbf{e}\|^2\right)$
  - $\pi_{\text{data}}(\mathbf{y} | \mathbf{x}) = \pi_{\text{noise}}(\mathcal{A}(\mathbf{x}) - \mathbf{y}) \implies -\log \pi_{\text{data}}(\mathbf{y} | \mathbf{x}) \propto \frac{1}{\sigma^2} \|\mathcal{A}(\mathbf{x}) - \mathbf{y}\|^2$
- Maximum likelihood:

$$\arg \min_x \frac{1}{\sigma^2} \|\mathcal{A}(\mathbf{x}) - \mathbf{y}\|^2$$

# Additive Gaussian noise and Gibbs type of prior

- Inverse problem:  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$
- Data likelihood:  $\mathbf{e} \sim \text{Normal}(0, \sigma^2)$ 
  - $\pi_{\text{noise}}(\mathbf{e}) \propto \exp\left(-\frac{1}{\sigma^2} \|\mathbf{e}\|^2\right)$
  - $\pi_{\text{data}}(\mathbf{y} | \mathbf{x}) = \pi_{\text{noise}}(\mathcal{A}(\mathbf{x}) - \mathbf{y}) \implies -\log \pi_{\text{data}}(\mathbf{y} | \mathbf{x}) \propto \frac{1}{\sigma^2} \|\mathcal{A}(\mathbf{x}) - \mathbf{y}\|^2$
- Maximum likelihood:
$$\arg \min_{\mathbf{x}} \frac{1}{\sigma^2} \|\mathcal{A}(\mathbf{x}) - \mathbf{y}\|^2$$
- Prior: Gibbs type of prior  $\pi_{\text{prior}}(\mathbf{x}) \propto \exp(-S_\theta(\mathbf{x}))$  with  $S_\theta: X \rightarrow \mathbb{R}$  convex.
- Maximum a posteriori:

$$\arg \min_{\mathbf{x}} \left[ S_\theta(\mathbf{x}) + \frac{1}{\sigma^2} \|\mathcal{A}(\mathbf{x}) - \mathbf{y}\|^2 \right]$$

# Choosing the prior

- Selecting the prior:
  - Theory: Influence of prior on posterior diminishes as noise  $\rightarrow 0$ .
  - Prior acts primarily as a regulariser

# Choosing the prior

- Selecting the prior:
  - Theory: Influence of prior on posterior diminishes as noise  $\rightarrow 0$ .
  - Prior acts primarily as a regulariser
  - Choice of prior matters for finite data with non-zero observational noise!
  - Assign high probability to “natural” model parameters.

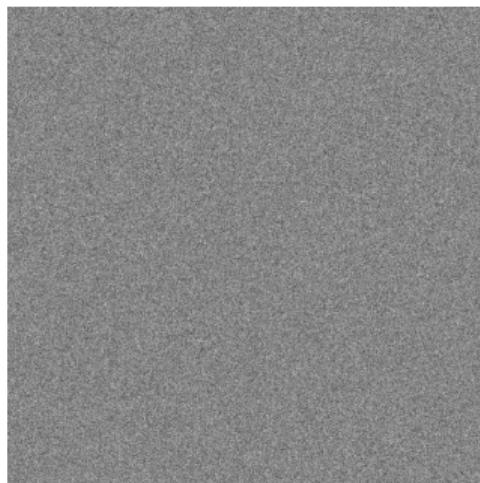
# Choosing the prior

- Selecting the prior:
  - Theory: Influence of prior on posterior diminishes as noise  $\rightarrow 0$ .
  - Prior acts primarily as a regulariser
  - Choice of prior matters for finite data with non-zero observational noise!
  - Assign high probability to “natural” model parameters.
- Gibbs type of prior:

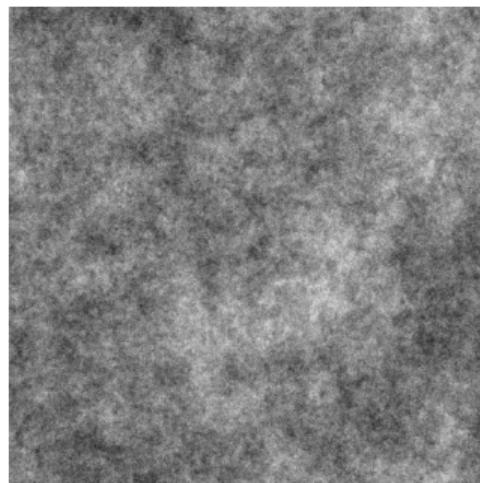
$$\pi_{\text{prior}}(x) \propto \exp(-S_\theta(x)) \quad \text{with } S_\theta: X \rightarrow \mathbb{R} \text{ convex.}$$

Often  $\theta$  is scalar and  $S_\theta(x) = \theta S(x)$  for fixed  $S: X \rightarrow \mathbb{R}$ .

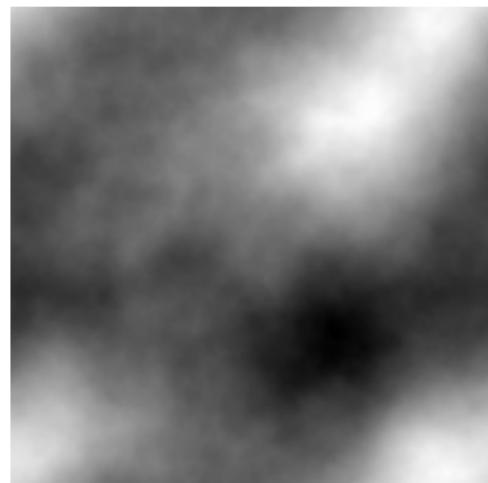
# Gibbs type of priors in imaging



$$S(x) = \|x\|_2^2$$

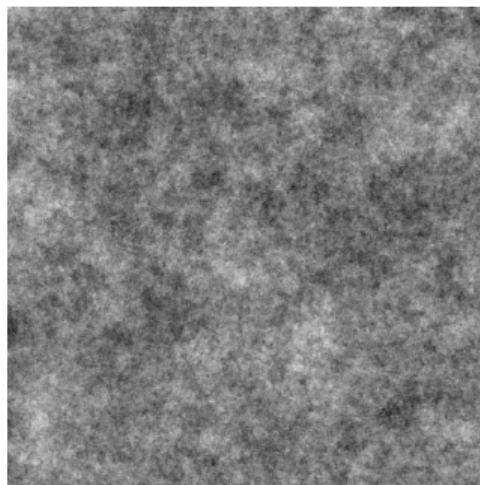


$$S(x) = \|\nabla x\|_2^2$$

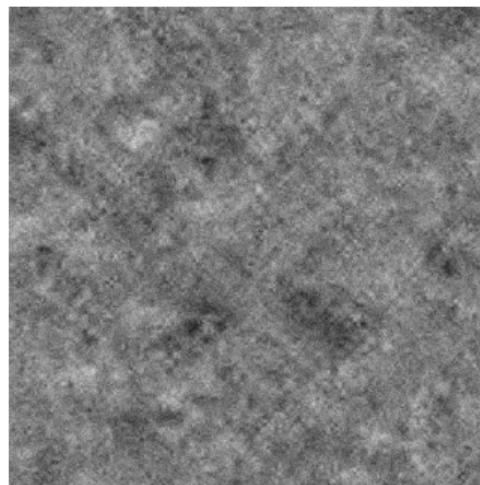


$$S(x) = \|\Delta x\|_2^2$$

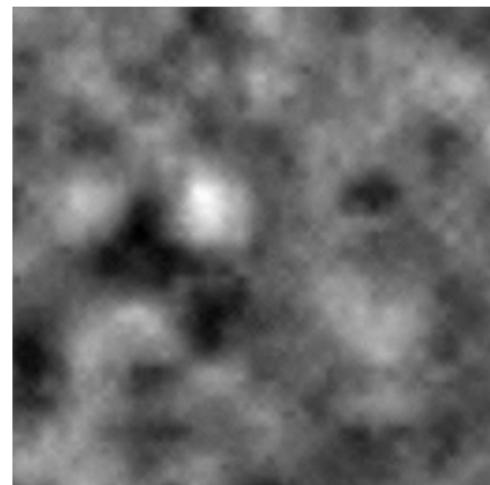
# Gibbs type of priors in imaging



$$S(x) = \|\nabla x\|_1$$

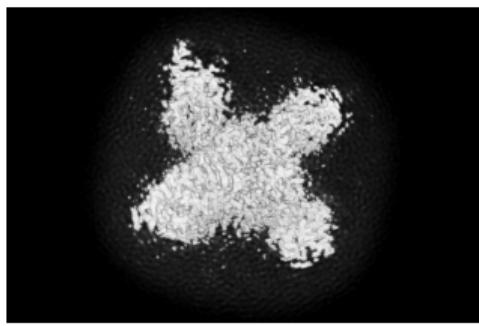
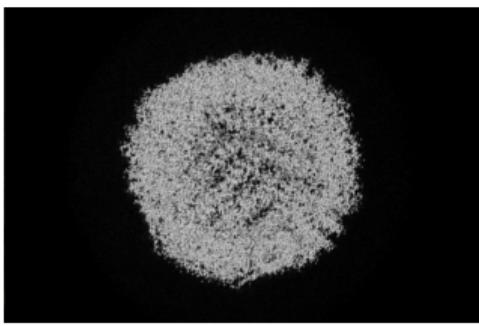
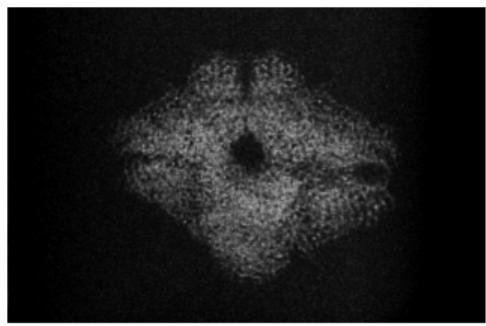
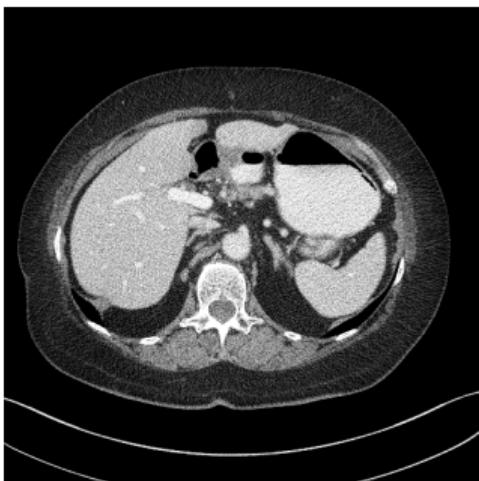


$$S(x) = \|x\|_{B_{1,1}^1}$$



$$S(x) = \|x\|_{B_{1,1}^2}$$

## Examples of natural images



# Conclusions about Bayesian inversion

- + General framework for solving inverse problems, plug-and-play specification of:
  - Forward model and statistical model of noise
  - Prior
- + Strong regularising properties
- + Parameter free
- + Uncertainty quantification
- Classical numerical methods unfeasible:
  - Require handcrafted prior
  - Computing many estimators, like conditional mean, computationally unfeasible.

# Conclusions about Bayesian inversion

- + General framework for solving inverse problems, plug-and-play specification of:
  - Forward model and statistical model of noise
  - Prior
- + Strong regularising properties
- + Parameter free
- + Uncertainty quantification
- Classical numerical methods unfeasible:
  - Require handcrafted prior
  - Computing many estimators, like conditional mean, computationally unfeasible.

Only MAP estimator currently computed in large-scale imaging applications.

# Conclusions about Bayesian inversion

- + General framework for solving inverse problems, plug-and-play specification of:
  - Forward model and statistical model of noise
  - Prior
- + Strong regularising properties
- + Parameter free
- + Uncertainty quantification
- Classical numerical methods unfeasible:
  - Require handcrafted prior
  - Computing many estimators, like conditional mean, computationally unfeasible.

Only MAP estimator currently computed in large-scale imaging applications.

- Can one use techniques from deep learning to address above disadvantages?

## Deep direct estimation

Joint work with Jonas Adler

Reference(s): (Adler & Öktem, 2017, 2018b, 2018a)

# Estimators in Bayesian inversion

Bayesian inversion: Recover the (posterior) distribution for  $(\mathbf{x} \mid \mathbf{y} = \mathbf{y})$ .

Here,  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $\mathbf{y}$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Challenge: Posterior too complex, explore it by:

- Computing suitable estimators
- Generating samples.



# Estimators in Bayesian inversion

Bayesian inversion: Recover a **suitable estimator** for  $(\mathbf{x} \mid \mathbf{y} = \mathbf{y})$ .

Here,  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $\mathbf{y}$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Challenge: Posterior too complex, explore it by:

- Computing suitable estimators
- Generating samples.



# Estimators in Bayesian inversion

Bayesian inversion: Recover a **suitable estimator** for  $(\mathbf{x} \mid \mathbf{y} = y)$ .

Here,  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$  with  $x^*$  unknown.

## Example estimators

- Mean (reconstruction):

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- Point-wise variance (uncertainty):

$$\mathbb{E}\left[ (\mathbf{x} - \mathbb{E}[\mathbf{x} \mid \mathbf{y} = y])^2 \mid \mathbf{y} = y \right]$$

- Point-wise covariance:

$$\mathbb{E}\left[ (\mathbf{x}_1 - \mathbb{E}[\mathbf{x}_1 \mid \mathbf{y} = y])(\mathbf{x}_2 - \mathbb{E}[\mathbf{x}_2 \mid \mathbf{y} = y]) \mid \mathbf{y} = y \right]$$

- Bayesian hypothesis testing (decision making):

$$\mathbb{P}(\mathbf{x} \in X_0 \mid \mathbf{y} = y) = \mathbb{E}[\mathbb{1}_{X_0}(\mathbf{x}) \mid \mathbf{y} = y]$$



# Estimators in Bayesian inversion

Bayesian inversion: Recover a suitable estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$ .

Here,  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$  with  $x^*$  unknown.

## Example estimators

- Mean (reconstruction):

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- Point-wise variance (uncertainty):

$$\mathbb{E}\left[ (\mathbf{x} - \mathbb{E}[\mathbf{x} \mid \mathbf{y} = y])^2 \mid \mathbf{y} = y \right]$$

- Point-wise covariance:

$$\mathbb{E}\left[ (\mathbf{x}_1 - \mathbb{E}[\mathbf{x}_1 \mid \mathbf{y} = y])(\mathbf{x}_2 - \mathbb{E}[\mathbf{x}_2 \mid \mathbf{y} = y]) \mid \mathbf{y} = y \right]$$

- Bayesian hypothesis testing (decision making):

$$\mathbb{P}(\mathbf{x} \in X_0 \mid \mathbf{y} = y) = \mathbb{E}[\mathbb{1}_{X_0}(\mathbf{x}) \mid \mathbf{y} = y]$$



# Estimators in Bayesian inversion

Bayesian inversion: Recover a suitable estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$ .

Here,  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  is a single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

## Example estimators

- $\mathbb{E}[\mathbf{w} \mid \mathbf{y} = y]$  for suitable choice of  $\mathbf{w}$ 
  - Mean:  $\mathbf{w} := \mathbf{x}$
  - Point-wise variance:  $\mathbf{w} := (\mathbf{x} - \mathbb{E}[\mathbf{x} \mid \mathbf{y} = y])^2$
  - Bayesian hypothesis testing:  $\mathbf{w} := \mathbb{1}_{X_0}(\mathbf{x})$
- Deep direct estimation: Use machine learning to directly compute such estimators from supervised training data.



# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(y) := \mathbb{E}[\mathbf{w} | \mathbf{y} = y]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).  
Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} | \mathbf{x} = \mathbf{x}^*)$  with  $x^*$  unknown.

# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(y) := \mathbb{E}[\mathbf{w} | \mathbf{y} = y]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).

Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} | \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Approach: Learn  $y \mapsto \hat{\mathcal{R}}(y)$  from  $(x_i, y_i) \in X \times Y$  i.i.d. samples of  $(\mathbf{x}, \mathbf{y})$ .

# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(y) := \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).

Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Approach: Learn  $y \mapsto \hat{\mathcal{R}}(y)$  from  $(x_i, y_i) \in X \times Y$  i.i.d. samples of  $(\mathbf{x}, \mathbf{y})$ .

- Let  $Y$  be a measurable space,  $W$  a measurable Hilbert space, and  $\mathbf{y}$  and  $\mathbf{w}$  are  $Y$ - and  $W$ -valued random variables, respectively. Then

$$\hat{\mathcal{R}} \in \arg \min_{\mathcal{R}: Y \rightarrow W} \mathbb{E}_{(\mathbf{w}, \mathbf{y})} \left[ \|\mathcal{R}(\mathbf{y}) - \mathbf{w}\|_W^2 \right] \implies \hat{\mathcal{R}}(y) = \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y].$$

# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(y) := \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).

Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Approach: Learn  $y \mapsto \hat{\mathcal{R}}(y)$  from  $(x_i, y_i) \in X \times Y$  i.i.d. samples of  $(\mathbf{x}, \mathbf{y})$ .

- Let  $Y$  be a measurable space,  $W$  a measurable Hilbert space, and  $\mathbf{y}$  and  $\mathbf{w}$  are  $Y$ - and  $W$ -valued random variables, respectively. Then

$$\hat{\mathcal{R}} \in \arg \min_{\mathcal{R}: Y \rightarrow W} \mathbb{E}_{(\mathbf{w}, \mathbf{y})} \left[ \|\mathcal{R}(\mathbf{y}) - \mathbf{w}\|_W^2 \right] \implies \hat{\mathcal{R}}(y) = \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y].$$

- Supervised learning: By law of large numbers

$$\frac{1}{m} \sum_{i=1}^m \|\mathcal{R}(y_i) - w_i\|_W^2 \longrightarrow \mathbb{E}_{(\mathbf{w}, \mathbf{y})} \left[ \|\mathcal{R}(\mathbf{y}) - \mathbf{w}\|_W^2 \right] \text{ a.s. as } m \rightarrow \infty$$

$(w_i, y_i)$  i.i.d. samples of  $(\mathbf{w}, \mathbf{y})$  where  $w_i$  computed from  $(x_i, y_i)$ .

# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(\mathbf{y}) := \mathbb{E}[\mathbf{w} \mid \mathbf{y} = \mathbf{y}]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).

Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Approach: Learn  $y \mapsto \hat{\mathcal{R}}(y)$  from  $(x_i, y_i) \in X \times Y$  i.i.d. samples of  $(\mathbf{x}, \mathbf{y})$ .

- Let  $Y$  be a measurable space,  $W$  a measurable Hilbert space, and  $\mathbf{y}$  and  $\mathbf{w}$  are  $Y$ - and  $W$ -valued random variables, respectively. Then

$$\hat{\mathcal{R}} \in \arg \min_{\mathcal{R}: Y \rightarrow W} \mathbb{E}_{(\mathbf{w}, \mathbf{y})} \left[ \|\mathcal{R}(\mathbf{y}) - \mathbf{w}\|_W^2 \right] \implies \hat{\mathcal{R}}(y) = \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y].$$

- Supervised learning:

$$\hat{\mathcal{R}} \approx \arg \min_{\mathcal{R}: Y \rightarrow W} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{R}(y_i) - w_i\|_W^2 \right].$$

# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(\mathbf{y}) := \mathbb{E}[\mathbf{w} \mid \mathbf{y} = \mathbf{y}]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).

Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Approach: Learn  $y \mapsto \hat{\mathcal{R}}(y)$  from  $(x_i, y_i) \in X \times Y$  i.i.d. samples of  $(\mathbf{x}, \mathbf{y})$ .

- Let  $Y$  be a measurable space,  $W$  a measurable Hilbert space, and  $\mathbf{y}$  and  $\mathbf{w}$  are  $Y$ - and  $W$ -valued random variables, respectively. Then

$$\hat{\mathcal{R}} \in \arg \min_{\mathcal{R}: Y \rightarrow W} \mathbb{E}_{(\mathbf{w}, \mathbf{y})} \left[ \|\mathcal{R}(\mathbf{y}) - \mathbf{w}\|_W^2 \right] \implies \hat{\mathcal{R}}(y) = \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y].$$

- Supervised learning:

$$\hat{\mathcal{R}} \approx \arg \min_{\mathcal{R}: Y \rightarrow W} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{R}(y_i) - w_i\|_W^2 \right].$$

Issue: Computationally unfeasible + ill-posed (overfitting)

# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(\mathbf{y}) := \mathbb{E}[\mathbf{w} \mid \mathbf{y} = \mathbf{y}]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).

Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Approach: Learn  $y \mapsto \hat{\mathcal{R}}(y)$  from  $(x_i, y_i) \in X \times Y$  i.i.d. samples of  $(\mathbf{x}, \mathbf{y})$ .

- Let  $Y$  be a measurable space,  $W$  a measurable Hilbert space, and  $\mathbf{y}$  and  $\mathbf{w}$  are  $Y$ - and  $W$ -valued random variables, respectively. Then

$$\hat{\mathcal{R}} \in \arg \min_{\mathcal{R}: Y \rightarrow W} \mathbb{E}_{(\mathbf{w}, \mathbf{y})} \left[ \|\mathcal{R}(\mathbf{y}) - \mathbf{w}\|_W^2 \right] \implies \hat{\mathcal{R}}(y) = \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y].$$

- Empirical risk minimisation:

$$\hat{\mathcal{R}} \approx \mathcal{R}_{\hat{\theta}} \quad \text{with} \quad \hat{\theta} \in \arg \min_{\theta} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{R}_\theta(y_i) - w_i\|_W^2 \right].$$

Issue: Computationally unfeasible + ill-posed (overfitting)

$\implies$  Parametrised family  $\mathcal{R}_\theta: Y \rightarrow W$  (deep neural network)

# Deep direct estimation

Case with supervised data

Task: Compute estimator  $\hat{\mathcal{R}}(y) := \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y]$  for suitable  $\mathbf{w}$  (function of  $\mathbf{x}$  and  $\mathbf{y}$ ).

Here  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  single sample of  $(\mathbf{y} \mid \mathbf{x} = \mathbf{x}^*)$  with  $\mathbf{x}^*$  unknown.

Approach: Learn  $y \mapsto \hat{\mathcal{R}}(y)$  from  $(x_i, y_i) \in X \times Y$  i.i.d. samples of  $(\mathbf{x}, \mathbf{y})$ .

- Let  $Y$  be a measurable space,  $W$  a measurable Hilbert space, and  $\mathbf{y}$  and  $\mathbf{w}$  are  $Y$ - and  $W$ -valued random variables, respectively. Then

$$\hat{\mathcal{R}} \in \arg \min_{\mathcal{R}: Y \rightarrow W} \mathbb{E}_{(\mathbf{w}, \mathbf{y})} \left[ \|\mathcal{R}(\mathbf{y}) - \mathbf{w}\|_W^2 \right] \implies \hat{\mathcal{R}}(y) = \mathbb{E}[\mathbf{w} \mid \mathbf{y} = y].$$

- Empirical risk minimisation:

$$\hat{\mathcal{R}} \approx \mathcal{R}_{\hat{\theta}} \quad \text{with} \quad \hat{\theta} \in \arg \min_{\theta} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{R}_\theta(y_i) - w_i\|_W^2 \right].$$

Architecture: Parametrisation of  $\mathcal{R}_\theta$  may encode **relation** between  $\mathbf{x}$  and  $\mathbf{y}$ :

$$\mathcal{A}: X \rightarrow Y \text{ and } [\partial \mathcal{A}(y)]^*: Y \rightarrow X.$$

# Deep direct estimation

Conditional mean: Unrolling a gradient decent scheme

**Unrolling:** Deep neural network architecture adapted to approximate an operator defined by an iterative scheme.

- Solution operators to PDEs (Hsieh et al., 2019; Rizzuti et al., 2019)
- Optimisation solvers (Gregor & LeCun, 2010; Banert et al., 2018)
- Regularised inverse (Putzky & Welling, 2017; Adler & Öktem, 2017)

# Deep direct estimation

Conditional mean: Unrolling a gradient decent scheme

**Unrolling:** Deep neural network architecture adapted to approximate an operator defined by an iterative scheme.

- Solution operators to PDEs (Hsieh et al., 2019; Rizzuti et al., 2019)
- Optimisation solvers (Gregor & LeCun, 2010; Banert et al., 2018)
- **Regularised inverse** (Putzky & Welling, 2017; Adler & Öktem, 2017)

---

## Algorithm Gradient descent

---

```
1: for  $i = 1, \dots$  do
2:    $x_{i+1} \leftarrow x_i - \alpha [\partial \mathcal{A}(x_i)]^* (\mathcal{A}(x_i) - y)$ 
```

---

Solution operator to variational problem

$$\min_x \|\mathcal{A}(x) - y\|_Y^2$$

# Deep direct estimation

Conditional mean: Unrolling a gradient decent scheme

**Unrolling:** Deep neural network architecture adapted to approximate an operator defined by an iterative scheme.

- Solution operators to PDEs (Hsieh et al., 2019; Rizzuti et al., 2019)
- Optimisation solvers (Gregor & LeCun, 2010; Banert et al., 2018)
- Regularised inverse (Putzky & Welling, 2017; Adler & Öktem, 2017)

---

## Algorithm Gradient descent

---

```
1: for  $i = 1, \dots$  do
2:    $x_{i+1} \leftarrow x_i - \alpha [\partial \mathcal{A}(x_i)]^* (\mathcal{A}(x_i) - y)$ 
```

---

## Solution operator to variational problem

$$\min_x \|\mathcal{A}(x) - y\|_Y^2$$

---

## Algorithm Learned gradient descent

---

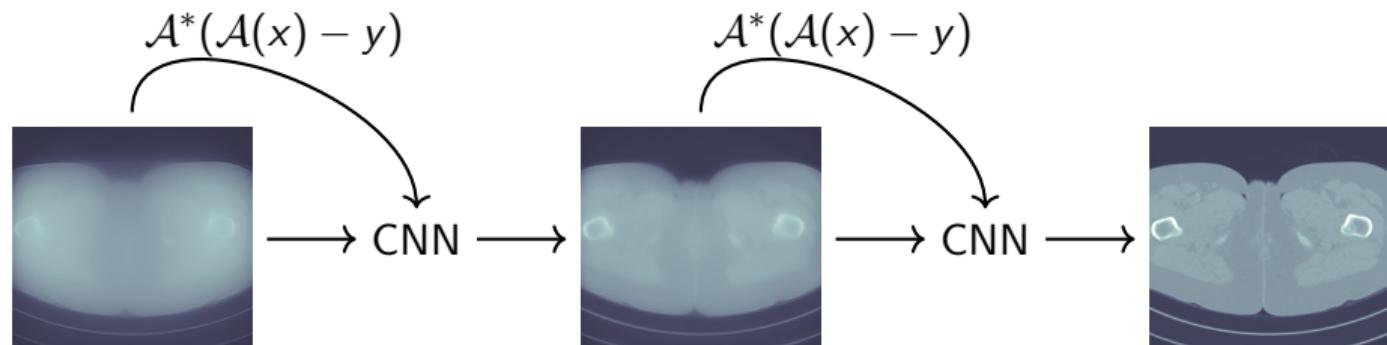
```
1: for  $i = 1, \dots, N$  do
2:    $x_{i+1} \leftarrow \Lambda_\theta(x_i, [\partial \mathcal{A}(x_i)]^* (\mathcal{A}(x_i) - y))$ 
3:    $\mathcal{R}_\theta(y) \leftarrow x_N$ 
```

---

- Finite number of iterates  $N$ .
- $\Lambda_\theta$  given by a convolutional neural network (CNN), learn parameter  $\theta$  from training data in  $X \times Y$ .

# Unrolling

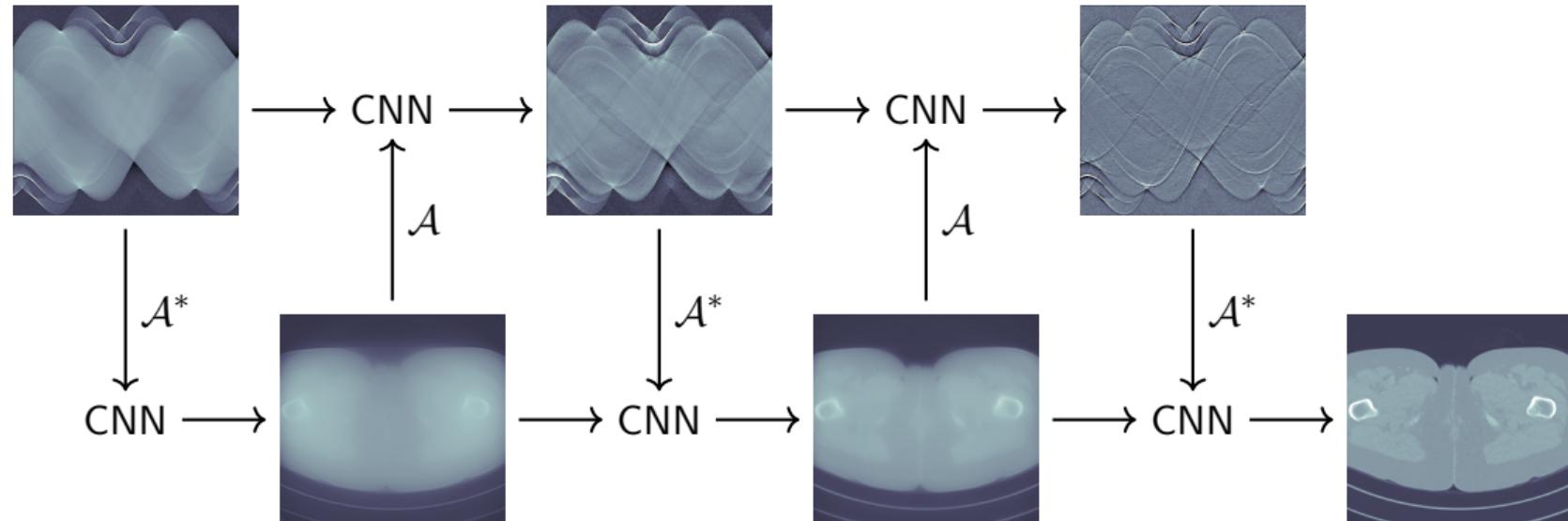
Gradient descent scheme



- $\mathcal{A}: X \rightarrow Y$  and  $\mathcal{A}^*: Y \rightarrow X$  handcrafted, not learned
- CNN: Convolutional neural network, learned

# Unrolling

Primal dual scheme



- $\mathcal{A}: X \rightarrow Y$  and  $\mathcal{A}^*: Y \rightarrow X$  handcrafted, not learned
- CNN: Convolutional neural network, learned

# Deep direct estimation

Conditional mean in 2D tomography of human phantom

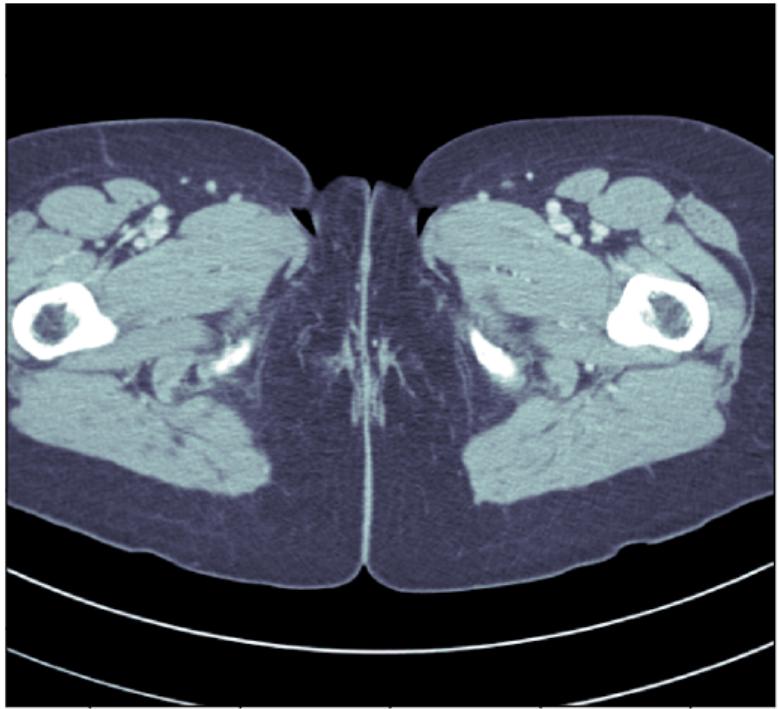
Inverse problem: Recover attenuation coefficient from tomographic data (sinogram)

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

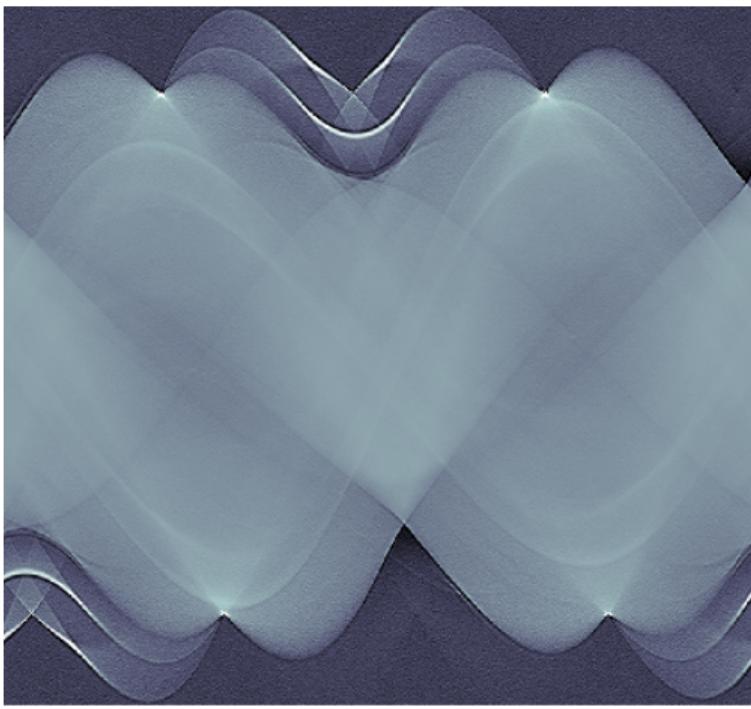
- Forward operator: 2D ray transform
- Geometry: Fan beam, 1000 lines/angle, 1000 angles
- Noise: Poisson noise,  $10^3$  incident photons/detector element and angle
- Image:  $512 \times 512$  pixel
- Training data: About 2 000 pairs  $(x_i, y_i)$  from 9 patients

Reconstruction methods:

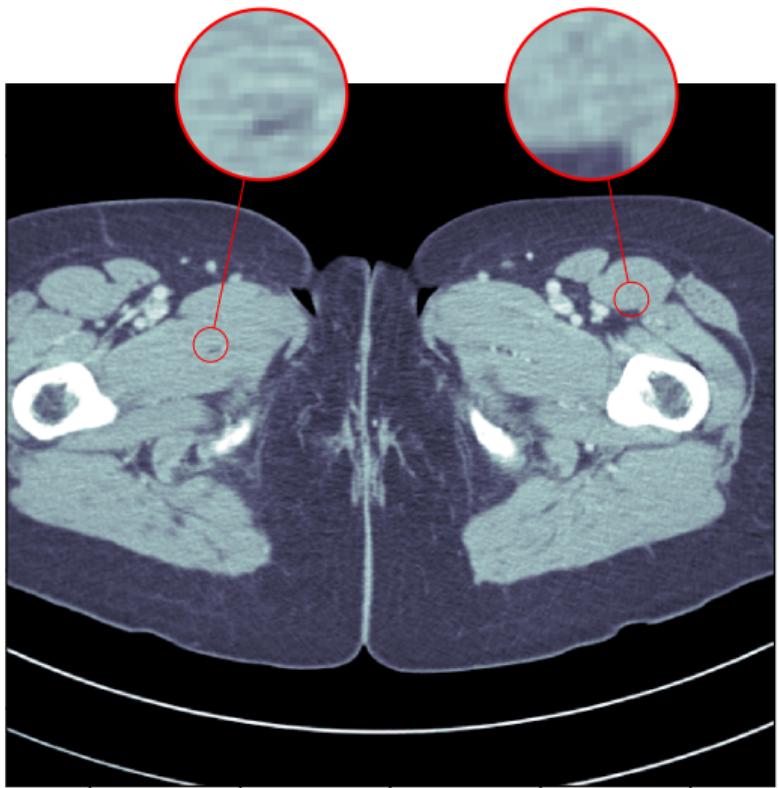
- Filtered backprojection (FBP)
- MAP estimator with total variation prior (MAP with TV)
- FBP reconstruction de-noised by U-Net (FBP + U-Net)
- Learned primal-dual reconstruction



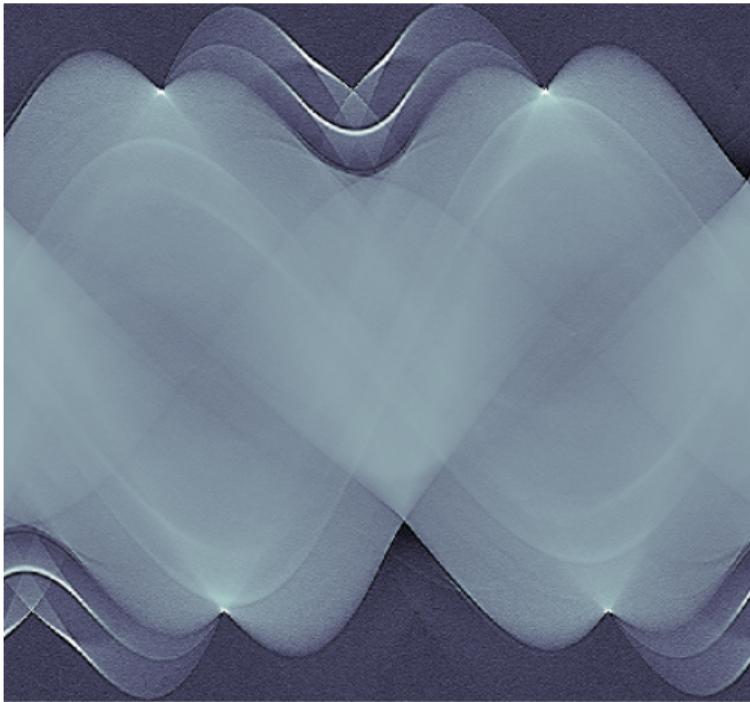
Ground truth.



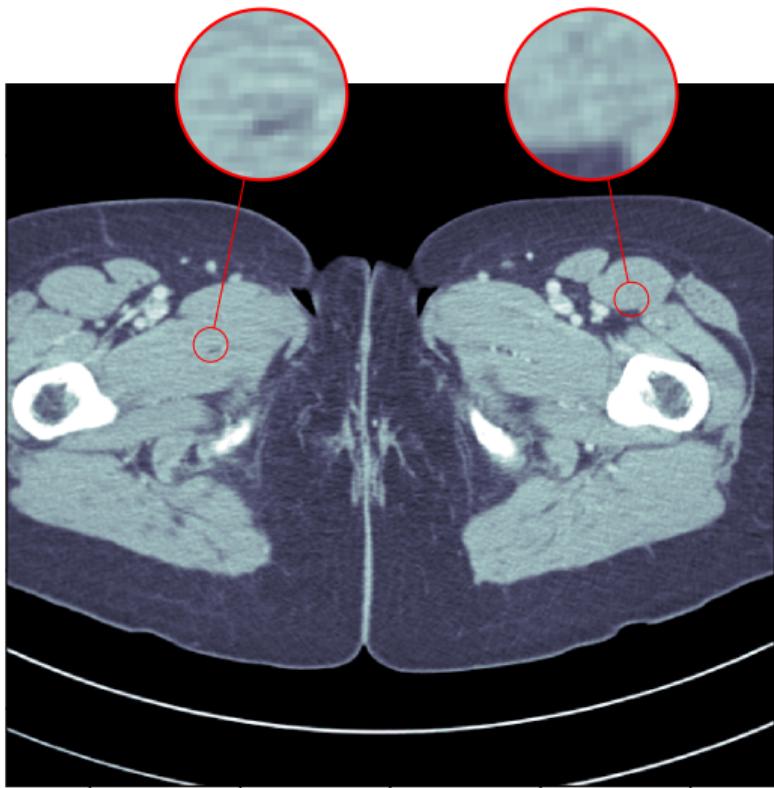
Data.



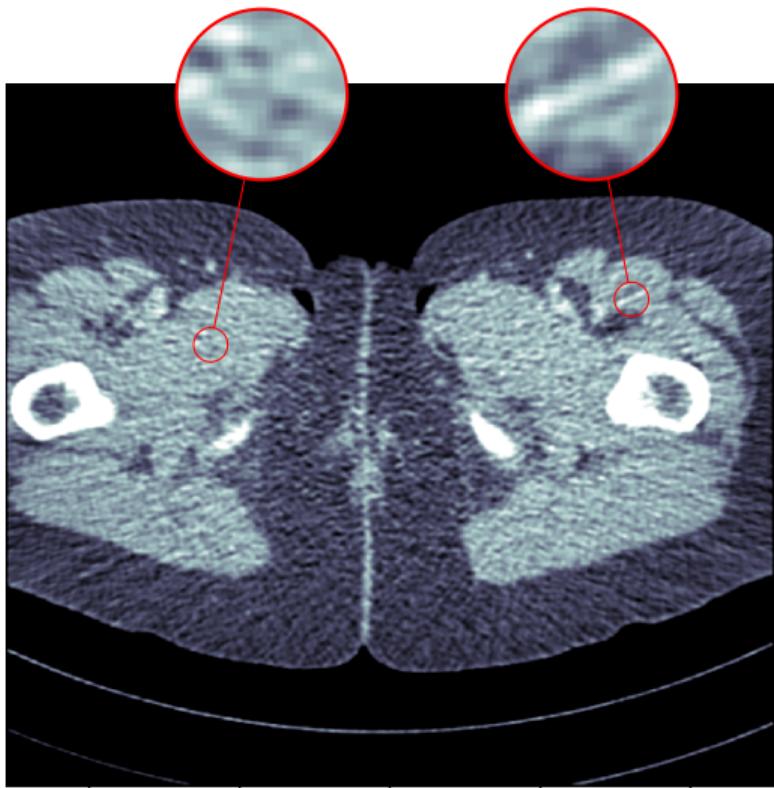
Ground truth.



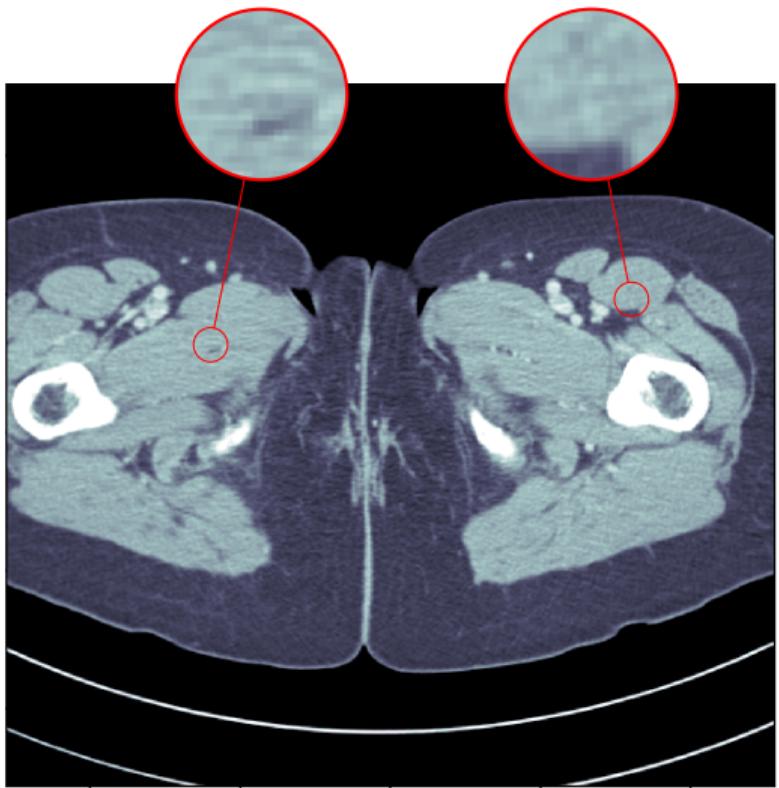
Data.



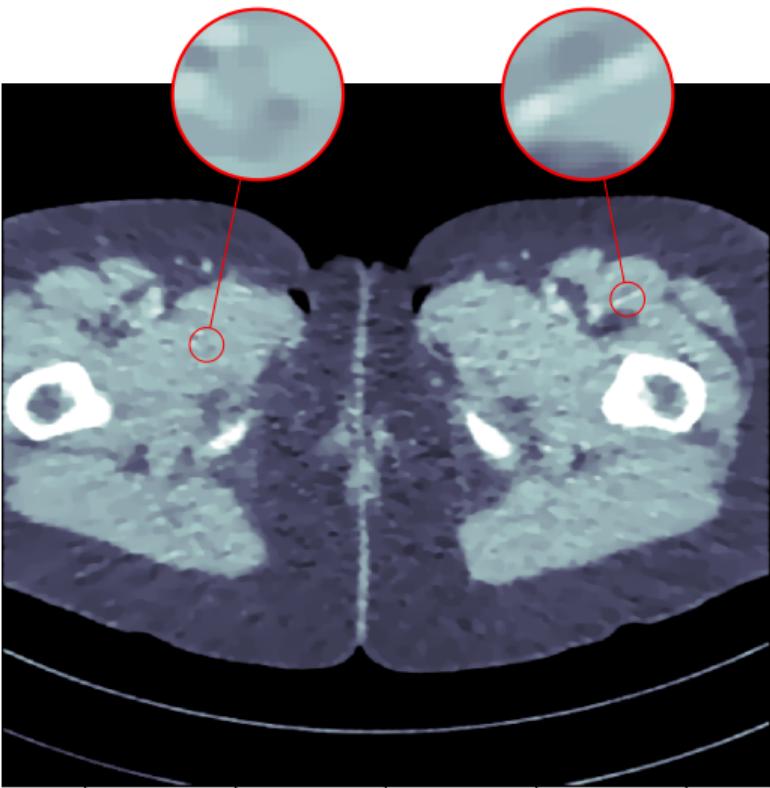
Ground truth.



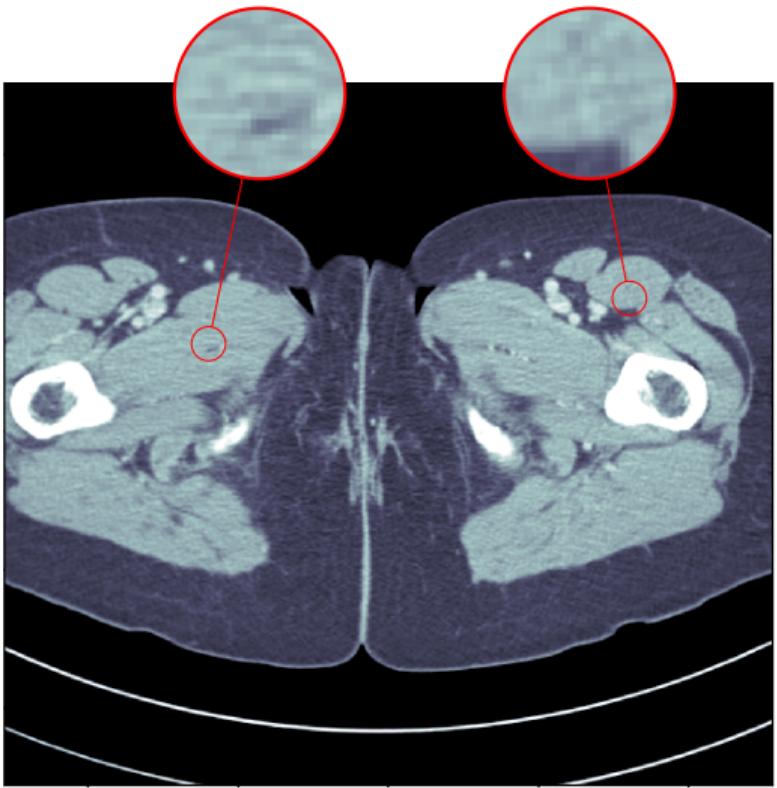
FBP.



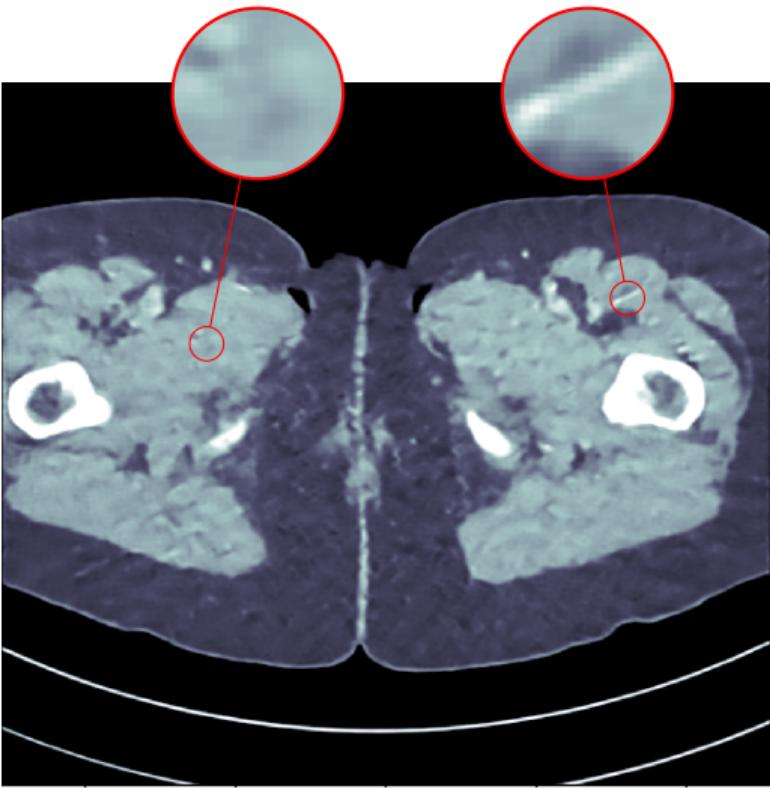
Ground truth.



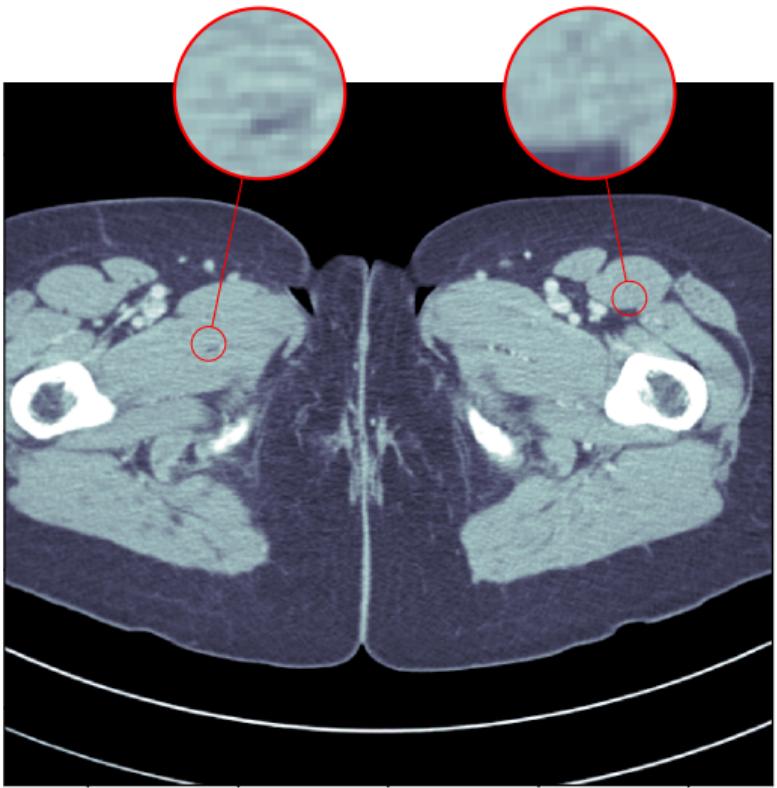
MAP with TV.



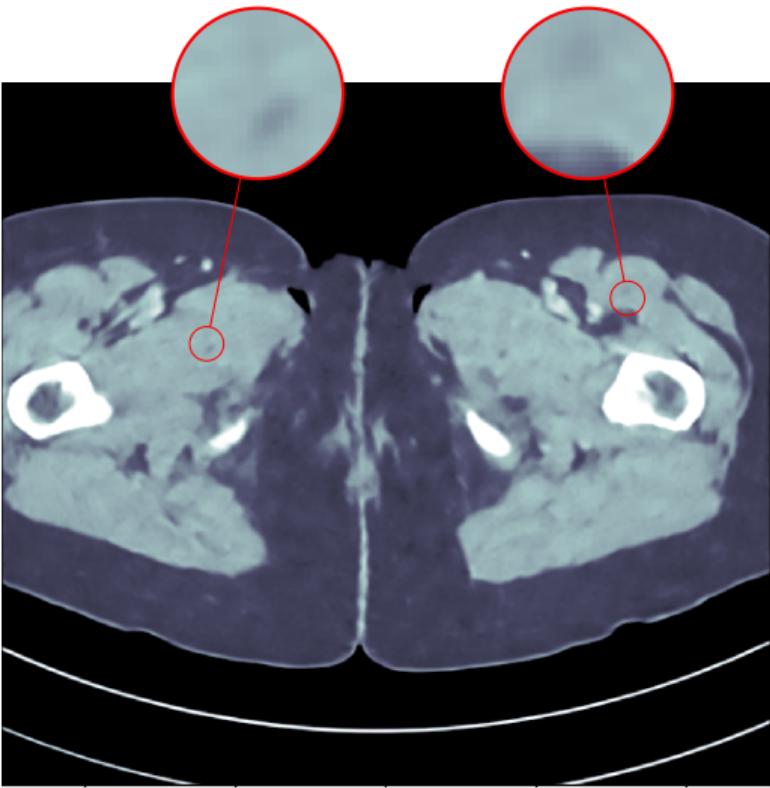
Ground truth.



FBP + U-Net denoising.



Ground truth.



Learned primal-dual.

# Learned iterative reconstruction

2D tomography of human phantom

Quantitative comparison (SSIM = structural similarity index, 1 = perfect match)

Method	PSNR (dB)	SSIM	Runtime (ms)	Parameters
FBP	33.65	0.829	423	1
MAP with TV	37.48	0.946	64 371	1
FBP + U-Net	41.92	0.941	463	$10^7$
Learned primal-dual	44.11	0.969	620	$2.4 \cdot 10^5$

## Comments

- Improved reconstruction quality against state-of-the-art
  - Very large improvement in PSNR
  - Significant clinically relevant improvement
- No need to manually set ‘obscure’ parameters
- Execution time allows for clinical implementation
- Very modest requirements on amount of supervised training data

# Learned iterative reconstruction

2D tomography of human phantom

Quantitative comparison (SSIM = structural similarity index, 1 = perfect match)

Method	PSNR (dB)	SSIM	Runtime (ms)	Parameters
FBP	33.65	0.829	423	1
MAP with TV	37.48	0.946	64 371	1
FBP + U-Net	41.92	0.941	463	$10^7$
Learned primal-dual	44.11	0.969	620	$2.4 \cdot 10^5$

## Comments

- Improved reconstruction quality against state-of-the-art
  - Very large improvement in PSNR
  - Significant clinically relevant improvement
- No need to manually set ‘obscure’ parameters
- Execution time allows for clinical implementation
- Very modest requirements on amount of supervised training data

# Integrating machine learning with inverse problems

- Characterise statistical properties of training data.
- Statistical model consistent with training data + loss function  
     $\Rightarrow$  type of estimator.
- Choose suitable neural network architecture for estimator:
  - Fully learned: Architecture entirely data driven.
  - Learned iterative: Architecture includes elements that make up the data likelihood,  
        e.g., the forward operator and the adjoint of its derivative.  
     $\Rightarrow$  parametrisation of estimator.
- Optimality criteria for neural network parameter + training data  
     $\Rightarrow$  learning/training problem.
- Optimal neural network parameter  $\Rightarrow$  learned estimator.

# Integrating machine learning with inverse problems

Estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

- Supervised learning

- Statistical model: Joint law  $(\mathbf{x}, \mathbf{y}) \sim \mu$  known.
- Training data:  $(x_i, y_i) \in X \times Y \implies \hat{\mu}$ .
- Estimator:  $\mathcal{R}_{\hat{\theta}}: Y \rightarrow X$  given as Bayes estimator, so  $\hat{\theta}$  solves

$$\hat{\theta} \in \arg \min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{\mu}} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})].$$

- Example references: (Putzky & Welling, 2017; Adler & Öktem, 2017, 2018b; Hammernik et al., 2018; Zhu et al., 2018)

# Integrating machine learning with inverse problems

Estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

- Learned prior/regulariser
  - Statistical model: Marginal distribution  $\mu_{\mathbf{x}}$  of  $(\mathbf{x}, \mathbf{y}) \sim \mu$  known.
  - Training data:  $x_i \in X \implies \hat{\mu}_{\mathbf{x}}$ .
  - Estimator:  $\mathcal{R}_{\hat{\theta}}: Y \rightarrow X$  given as MAP estimator with learned prior  $x \mapsto \hat{\mu}_{\mathbf{x}}$ , alternatively a variational method with learned regulariser  $x \mapsto -\log(\hat{\mu}_{\mathbf{x}}(x))$ .
  - Example references: (Romano, Elad, & Milanfar, 2017; Lunz et al., 2018; H. Li et al., 2018)

# Integrating machine learning with inverse problems

Estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

- **Unsupervised learning**

- **Statistical model:** Marginal distribution  $\mu_y$  of  $(\mathbf{x}, \mathbf{y}) \sim \mu$  known.
- **Training data:**  $y_i \in Y \implies \hat{\mu}_y$ .
- **Estimator:**  $\mathcal{R}_{\hat{\theta}}: Y \rightarrow X$  given as MAP estimator where  $\hat{\theta}$  solves

$$\hat{\theta} := \arg \min_{\theta} \mathbb{E}_{\mathbf{y} \sim \hat{\mu}_y} [\mathcal{L}(\mathcal{A}(\mathcal{R}_{\theta}(\mathbf{y})), \mathbf{y}) + \mathcal{S}(\mathcal{R}_{\theta}(\mathbf{y}))]$$

- **Example references:** (Gregor & LeCun, 2010; Oymak & Soltanolkotabi, 2017; Giryes et al., 2017)

# Integrating machine learning with inverse problems

Estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

- Unpaired data

- Statistical model: Marginal distributions  $\mu_{\mathbf{x}}$  and  $\mu_{\mathbf{y}}$  of  $(\mathbf{x}, \mathbf{y}) \sim \mu$  known.
- Training data:  $x_i \sim X$  and  $y_i \in Y \implies \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}$ .
- Estimator:  $\mathcal{R}_{\hat{\theta}}: Y \rightarrow X$  where  $\hat{\theta}$  solves

$$\begin{aligned}\hat{\theta} \in \arg \min_{\theta} \left\{ \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}} \left[ \ell_Y \left( \mathcal{A}(\mathcal{R}_{\theta}(\mathbf{y})), \mathbf{y} \right) + \ell_X \left( \mathcal{R}_{\theta}(\mathbf{y}), \mathbf{x} \right) \right] \right. \\ \left. + \lambda D \left( (\mathcal{R}_{\theta})_{\#}(\hat{\mu}_{\mathbf{y}}), \hat{\mu}_{\mathbf{x}} \right) \right\}.\end{aligned}$$

- Example references: (Mao et al., 2016; Mardani et al., 2019; Schwab et al., 2018)

# Integrating machine learning with inverse problems

Estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

- Unpaired data

- Statistical model: Marginal distributions  $\mu_{\mathbf{x}}$  and  $\mu_{\mathbf{y}}$  of  $(\mathbf{x}, \mathbf{y}) \sim \mu$  known.
- Training data:  $x_i \sim X$  and  $y_i \in Y \implies \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}$ .
- Estimator:  $\mathcal{R}_{\hat{\theta}}: Y \rightarrow X$  where  $\hat{\theta}$  solves

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}} \left[ \ell_Y \left( \mathcal{A}(\mathcal{R}_{\theta}(\mathbf{y})), \mathbf{y} \right) + \ell_X \left( \mathcal{R}_{\theta}(\mathbf{y}), \mathbf{x} \right) \right] + \lambda D \left( (\mathcal{R}_{\theta})_{\#}(\hat{\mu}_{\mathbf{y}}), \hat{\mu}_{\mathbf{x}} \right) \right\}.$$

- $\ell_X$  and  $\ell_Y$  loss functions on  $X$  and  $Y$ .

- Example references: (Mao et al., 2016; Mardani et al., 2019; Schwab et al., 2018)

# Integrating machine learning with inverse problems

Estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

- Unpaired data

- Statistical model: Marginal distributions  $\mu_{\mathbf{x}}$  and  $\mu_{\mathbf{y}}$  of  $(\mathbf{x}, \mathbf{y}) \sim \mu$  known.
- Training data:  $x_i \sim X$  and  $y_i \in Y \implies \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}$ .
- Estimator:  $\mathcal{R}_{\hat{\theta}}: Y \rightarrow X$  where  $\hat{\theta}$  solves

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}} \left[ \ell_Y \left( \mathcal{A}(\mathcal{R}_{\theta}(\mathbf{y})), \mathbf{y} \right) + \ell_X \left( \mathcal{R}_{\theta}(\mathbf{y}), \mathbf{x} \right) \right] + \lambda D \left( (\mathcal{R}_{\theta})_{\#}(\hat{\mu}_{\mathbf{y}}), \hat{\mu}_{\mathbf{x}} \right) \right\}.$$

- $\ell_X$  and  $\ell_Y$  loss functions on  $X$  and  $Y$ .
- $(\mathcal{R}_{\theta})_{\#}(\hat{\mu}_{\mathbf{y}})$  = push-forward of  $\hat{\mu}_{\mathbf{y}}$  by  $\mathcal{R}_{\theta}: Y \rightarrow X$

- Example references: (Mao et al., 2016; Mardani et al., 2019; Schwab et al., 2018)

# Integrating machine learning with inverse problems

Estimator for  $(\mathbf{x} \mid \mathbf{y} = y)$  where  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$  and  $y$  sample of  $(\mathbf{y} \mid \mathbf{x} = x^*)$ .

- Unpaired data

- Statistical model: Marginal distributions  $\mu_{\mathbf{x}}$  and  $\mu_{\mathbf{y}}$  of  $(\mathbf{x}, \mathbf{y}) \sim \mu$  known.
- Training data:  $x_i \sim X$  and  $y_i \in Y \implies \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}$ .
- Estimator:  $\mathcal{R}_{\hat{\theta}}: Y \rightarrow X$  where  $\hat{\theta}$  solves

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{\mu}_{\mathbf{x}} \otimes \hat{\mu}_{\mathbf{y}}} \left[ \ell_Y \left( \mathcal{A}(\mathcal{R}_{\theta}(\mathbf{y})), \mathbf{y} \right) + \ell_X \left( \mathcal{R}_{\theta}(\mathbf{y}), \mathbf{x} \right) \right] + \lambda \textcolor{red}{D} \left( (\mathcal{R}_{\theta})_{\#}(\hat{\mu}_{\mathbf{y}}), \hat{\mu}_{\mathbf{x}} \right) \right\}.$$

- $\ell_X$  and  $\ell_Y$  loss functions on  $X$  and  $Y$ .
- $(\mathcal{R}_{\theta})_{\#}(\hat{\mu}_{\mathbf{y}})$  = push-forward of  $\hat{\mu}_{\mathbf{y}}$  by  $\mathcal{R}_{\theta}: Y \rightarrow X$
- $\textcolor{red}{D}$  = distance between probability distributions on  $X$ , computed using GANs, introduces a separate deep neural network (discriminator/critic).
- Example references: (Mao et al., 2016; Mardani et al., 2019; Schwab et al., 2018)

## Task adapted reconstruction

Joint work with Jonas Adler, Olivier Verdier, Sebastian Lunz, and Carola Schönlieb

Reference(s): (Adler et al., 2018a, 2018b)

# Task adapted reconstruction

## Motivation

- Medical imaging is not done for fun, we want to solve a **task**!
  - Images typically summarised, either by an expert or by using specific descriptors, in an analysis step.
  - Summaries used as input for decision making.
- Task adapted reconstruction: Methods that jointly perform reconstruction and task (segmentation, classification, radiomics, . . . ).

# Task adapted reconstruction

## Motivation

- Medical imaging is not done for fun, we want to solve a **task**!
  - Images typically summarised, either by an expert or by using specific descriptors, in an analysis step.
  - Summaries used as input for decision making.
- **Task adapted reconstruction:** Methods that jointly perform reconstruction and task (segmentation, classification, radiomics, . . . ).

# Task Adapted Reconstruction

## Role of loss functions

- Bayes estimator:  $\mathcal{R}_\theta: Y \rightarrow X$  where  $\theta$  minimises Bayes risk:

$$L(\theta) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mu} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})] \quad \text{given loss } \ell_X: X \times X \rightarrow \mathbb{R}.$$

- Traditional losses:  $\|x_1 - x_2\|_2^2$  or  $\|x_1 - x_2\|_1$
- Fancy losses
  - Adversarial: Uses a trained neural network to quantify if the result is 'good'.
  - Perceptual: Quantifies image similarity by comparing features extracted using a previously trained neural network.
- Adapt loss to the task.

# Task Adapted Reconstruction

## Role of loss functions

- Bayes estimator:  $\mathcal{R}_\theta: Y \rightarrow X$  where  $\theta$  minimises Bayes risk:

$$L(\theta) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mu} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})] \quad \text{given loss } \ell_X: X \times X \rightarrow \mathbb{R}.$$

- Traditional losses:  $\|x_1 - x_2\|_2^2$  or  $\|x_1 - x_2\|_1$
- Fancy losses
  - Adversarial: Uses a trained neural network to quantify if the result is 'good'.
  - Perceptual: Quantifies image similarity by comparing features extracted using a previously trained neural network.
- Adapt loss to the task.

# Task Adapted Reconstruction

## Role of loss functions

- Bayes estimator:  $\mathcal{R}_\theta: Y \rightarrow X$  where  $\theta$  minimises Bayes risk:

$$L(\theta) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mu} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})] \quad \text{given loss } \ell_X: X \times X \rightarrow \mathbb{R}.$$

- Traditional losses:  $\|x_1 - x_2\|_2^2$  or  $\|x_1 - x_2\|_1$
- Fancy losses
  - **Adversarial:** Uses a trained neural network to quantify if the result is ‘good’.
  - **Perceptual:** Quantifies image similarity by comparing features extracted using a previously trained neural network.
- Adapt loss to the task.

# Task Adapted Reconstruction

## Role of loss functions

- Bayes estimator:  $\mathcal{R}_\theta: Y \rightarrow X$  where  $\theta$  minimises Bayes risk:

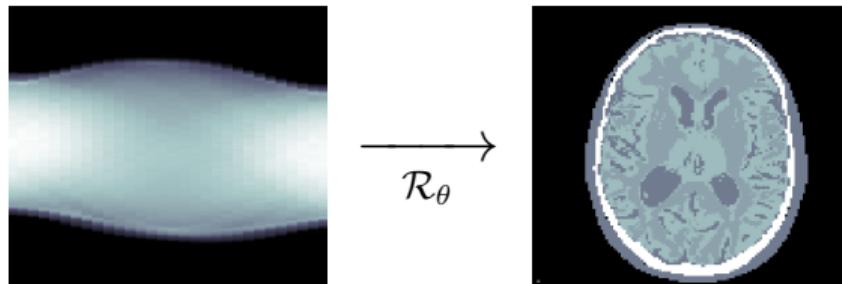
$$L(\theta) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mu} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})] \quad \text{given loss } \ell_X: X \times X \rightarrow \mathbb{R}.$$

- Traditional losses:  $\|x_1 - x_2\|_2^2$  or  $\|x_1 - x_2\|_1$
- Fancy losses
  - **Adversarial:** Uses a trained neural network to quantify if the result is ‘good’.
  - **Perceptual:** Quantifies image similarity by comparing features extracted using a previously trained neural network.
- Adapt loss to the task.

# Task Adapted Reconstruction

## Overview

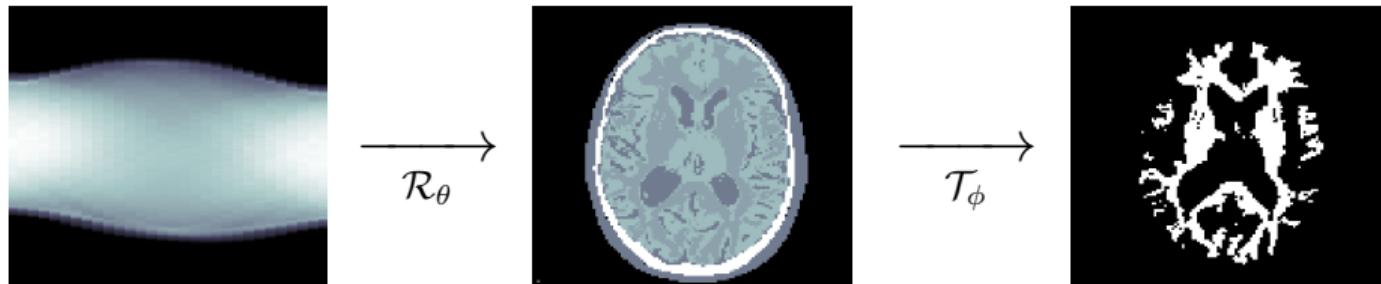
- We can learn to go from data to reconstruction
- Combine with learned task operator
- End-to-end differentiable training!



# Task Adapted Reconstruction

## Overview

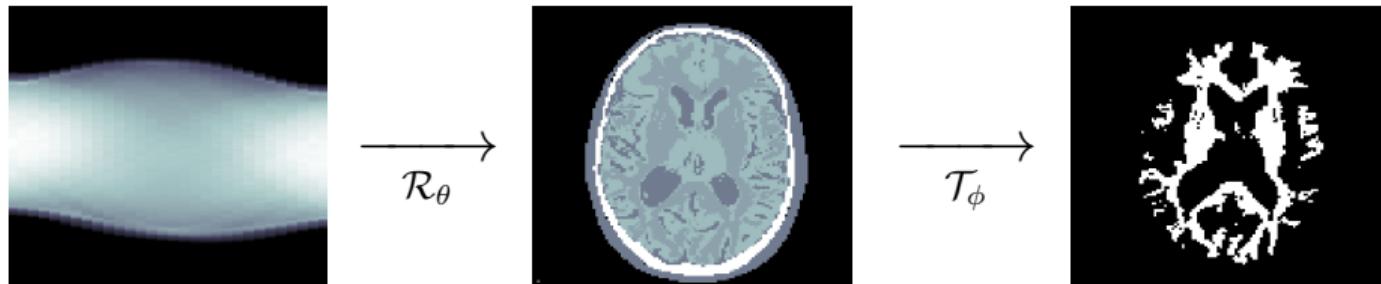
- We can learn to go from data to reconstruction
- Combine with learned task operator
- End-to-end differentiable training!



# Task Adapted Reconstruction

## Overview

- We can learn to go from data to reconstruction
- Combine with learned task operator
- End-to-end differentiable training!



# Task Adapted Reconstruction

## Approaches

- Sequential training: First train a reconstruction, then train the task

$$L_{\text{rec}}(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})]$$

$$L_{\text{task}}(\phi) = \mathbb{E}_{\mathbf{y}, \mathbf{d}} [\ell_D(\mathcal{T}_\phi \circ \mathcal{R}_{\hat{\theta}}(\mathbf{y}), \mathbf{d})] \quad \text{where } \hat{\theta} \text{ minimises } \theta \mapsto L_{\text{rec}}(\theta).$$

Training data: Samples  $(x_i, y_i)$  generated by  $(\mathbf{x}, \mathbf{y})$  and  $(y_i, d_i)$  generated by  $(\mathbf{y}, \mathbf{d})$

- End-to-end training: Straight from data to task

$$L(\theta, \phi) = \mathbb{E}_{\mathbf{y}, \mathbf{d}} [\ell_D(\mathcal{T}_\phi \circ \mathcal{R}_\theta(\mathbf{y}), \mathbf{d})].$$

Training data: Samples  $(y_i, d_i)$  generated by  $(\mathbf{y}, \mathbf{d})$ .

- Task adapted training: Anything in between

$$L(\theta, \phi) = \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{d}} [C \ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x}) + (1 - C) \ell_D(\mathcal{T}_\phi \circ \mathcal{R}_\theta(\mathbf{y}), \mathbf{d})].$$

Training data: Samples  $(x_i, y_i, d_i)$  generated by  $(\mathbf{x}, \mathbf{y}, \mathbf{d})$ .

# Task Adapted Reconstruction

## Approaches

- Sequential training: First train a reconstruction, then train the task

$$L_{\text{rec}}(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})]$$

$$L_{\text{task}}(\phi) = \mathbb{E}_{\mathbf{y}, \mathbf{d}} [\ell_D(\mathcal{T}_\phi \circ \mathcal{R}_{\hat{\theta}}(\mathbf{y}), \mathbf{d})] \quad \text{where } \hat{\theta} \text{ minimises } \theta \mapsto L_{\text{rec}}(\theta).$$

Training data: Samples  $(x_i, y_i)$  generated by  $(\mathbf{x}, \mathbf{y})$  and  $(y_i, d_i)$  generated by  $(\mathbf{y}, \mathbf{d})$

- End-to-end training: Straight from data to task

$$L(\theta, \phi) = \mathbb{E}_{\mathbf{y}, \mathbf{d}} [\ell_D(\mathcal{T}_\phi \circ \mathcal{R}_\theta(\mathbf{y}), \mathbf{d})].$$

Training data: Samples  $(y_i, d_i)$  generated by  $(\mathbf{y}, \mathbf{d})$ .

- Task adapted training: Anything in between

$$L(\theta, \phi) = \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{d}} [C \ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x}) + (1 - C) \ell_D(\mathcal{T}_\phi \circ \mathcal{R}_\theta(\mathbf{y}), \mathbf{d})].$$

Training data: Samples  $(x_i, y_i, d_i)$  generated by  $(\mathbf{x}, \mathbf{y}, \mathbf{d})$ .

# Task Adapted Reconstruction

## Approaches

- Sequential training: First train a reconstruction, then train the task

$$L_{\text{rec}}(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x})]$$

$$L_{\text{task}}(\phi) = \mathbb{E}_{\mathbf{y}, \mathbf{d}} [\ell_D(\mathcal{T}_\phi \circ \mathcal{R}_{\hat{\theta}}(\mathbf{y}), \mathbf{d})] \quad \text{where } \hat{\theta} \text{ minimises } \theta \mapsto L_{\text{rec}}(\theta).$$

Training data: Samples  $(x_i, y_i)$  generated by  $(\mathbf{x}, \mathbf{y})$  and  $(y_i, d_i)$  generated by  $(\mathbf{y}, \mathbf{d})$

- End-to-end training: Straight from data to task

$$L(\theta, \phi) = \mathbb{E}_{\mathbf{y}, \mathbf{d}} [\ell_D(\mathcal{T}_\phi \circ \mathcal{R}_\theta(\mathbf{y}), \mathbf{d})].$$

Training data: Samples  $(y_i, d_i)$  generated by  $(\mathbf{y}, \mathbf{d})$ .

- Task adapted training: Anything in between

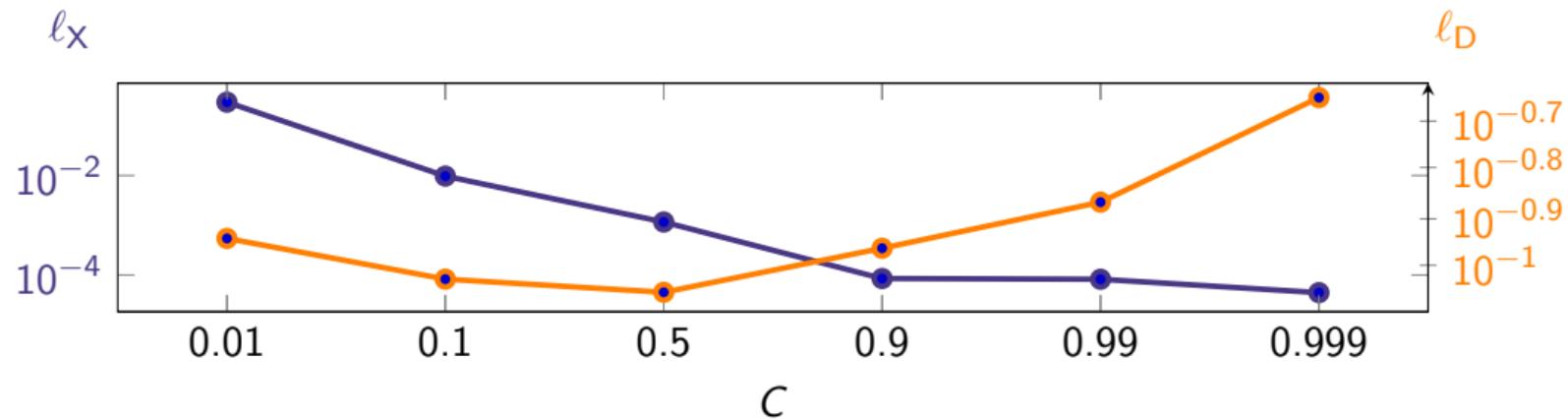
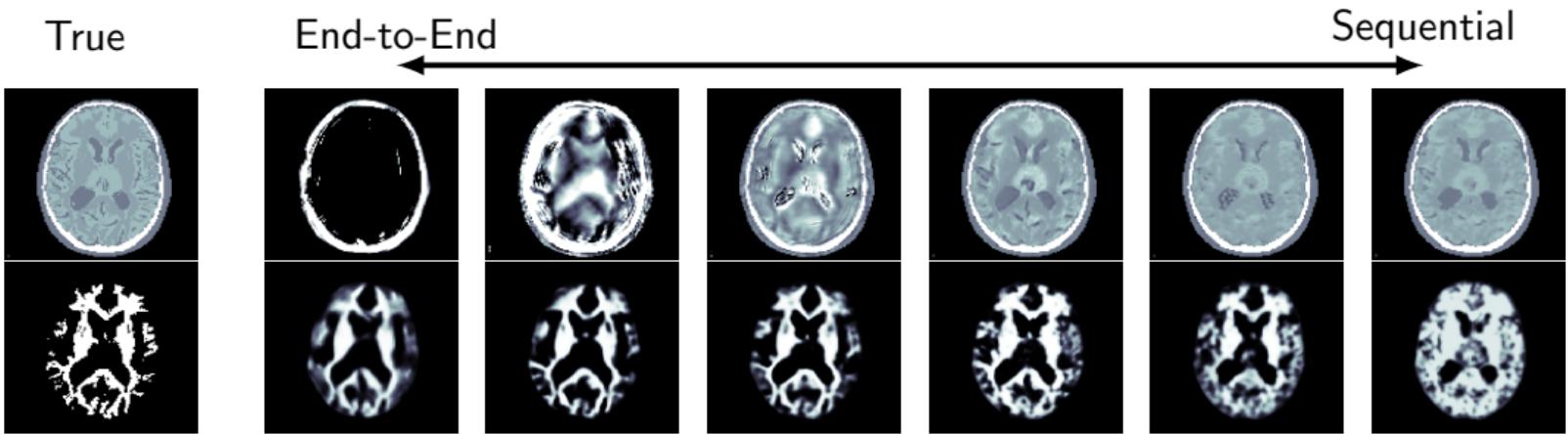
$$L(\theta, \phi) = \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{d}} [C \ell_X(\mathcal{R}_\theta(\mathbf{y}), \mathbf{x}) + (1 - C) \ell_D(\mathcal{T}_\phi \circ \mathcal{R}_\theta(\mathbf{y}), \mathbf{d})].$$

Training data: Samples  $(x_i, y_i, d_i)$  generated by  $(\mathbf{x}, \mathbf{y}, \mathbf{d})$ .

# Task Adapted Reconstruction

Joint reconstruction and segmentation

- 7 CT brain scans
  - Segmented semi-manually
  - Simulated low-dose data
- Task: Segment white matter given CT sinogram
- Reconstruction operator:  $\mathcal{R}_\theta: Y \rightarrow X$  (learned primal-dual)
- Task operator:  $\mathcal{T}_\phi: Y \rightarrow D$  (U-Net)



# Task Adapted Reconstruction

## Tasks

Segmentation only an example, can perform reconstruction jointly with **any** task given by a trainable differentiable neural network.

- Semantic segmentation (Thoma, 2016; Guo et al., 2018).
- Caption generation (Karpathy & Fei-Fei, 2017; C. Y. Li et al., 2018).
- Inpainting (Xie et al., 2012).
- Depixelization/super-resolution (Romano, Isidoro, & Milanfar, 2017).
- Demosaicing (Syu et al., 2018).
- Image translation (Wolterink et al., 2017).
- Object recognition (Sermanet et al., 2013; He et al., 2016; Farabet et al., 2013).
- Non-rigid image registration (Ghosal & Ray, 2017; Dalca et al., 2018).

...

## Deep posterior sampling

Joint work with Jonas Adler

Reference(s): (Adler & Öktem, 2018a)

# Sampling methods

- **Input:** Training data  $\mathbf{x}_i$  generated by  $\mathbf{x} \sim \mu$ .
- **Task:** Sample from unknown distribution  $\mu$ .
- **Approaches:**
  - MCMC: Proximal MCMC, Hamiltonian Monte-Carlo, ...
  - Variational Inference: Variational Bayes, Variational Auto-Encoders (deep learning variant of Variational Bayes), ...
  - Plug and Play Generative Networks
  - Pixel Recurrent Models
  - Generative Adversarial Networks

# Sampling methods

- Input: Training data  $\mathbf{x}_i$  generated by  $\mathbf{x} \sim \mu$ .
- Task: Sample from unknown distribution  $\mu$ .
- Approaches:
  - MCMC: Proximal MCMC, Hamiltonian Monte-Carlo, ...
  - Variational Inference: Variational Bayes, Variational Auto-Encoders (deep learning variant of Variational Bayes), ...
  - Plug and Play Generative Networks
  - Pixel Recurrent Models
  - Generative Adversarial Networks

# Generative Adversarial Networks

- **Main idea:** Train two networks, generator and discriminator.
- Generator tries to generate 'true' samples, discriminator tries to say 'good/bad'.



# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \mathcal{W}(\mathcal{G}_{\theta}, \mu)$$

- $\mathcal{W}$  is the Wasserstein 1-distance, quantifies similarity of probability measures.
- Generator:  $\mathcal{G}_{\theta}$  is a probability distribution on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \mathcal{W}(\mathcal{G}_{\theta}, \mu)$$

- $\mathcal{W}$  is the Wasserstein 1-distance, quantifies similarity of probability measures.
- Generator:  $\mathcal{G}_{\theta}$  is a probability distribution on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .

Issue:  $\mu$  unknown  $\implies$  not possible to evaluate  $\mathcal{W}(\mathcal{G}_{\theta}, \mu)$ .

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{D \in \text{Lip}(X)} \left\{ \mathbb{E}_{x \sim \mu} [D(x)] - \mathbb{E}_{v \sim \mathcal{G}_\theta} [D(v)] \right\} \right\}$$

- Discriminator:  $D: X \rightarrow \mathbb{R}$  is a 1-Lipschitz function.
- Generator:  $\mathcal{G}_\theta$  is a probability distribution on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .

Issue:  $\mu$  unknown  $\implies$  not possible to evaluate  $\mathcal{W}(\mathcal{G}_\theta, \mu)$ .

$\implies$  Re-write using the Kantorovich-Rubinstein dual characterization of  $\mathcal{W}$ .

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{D \in \text{Lip}(X)} \left\{ \mathbb{E}_{x \sim \mu} [D(x)] - \mathbb{E}_{v \sim \mathcal{G}_\theta} [D(v)] \right\} \right\}$$

- Discriminator:  $D: X \rightarrow \mathbb{R}$  is a **1-Lipschitz** function.
- Generator:  $\mathcal{G}_\theta$  is a probability distribution on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .

Issue: How to ensure  $D$  is 1-Lipschitz?

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{D: X \rightarrow \mathbb{R}} \left\{ \mathbb{E}_{\mathbf{x} \sim \mu} [D(\mathbf{x})] - \mathbb{E}_{\mathbf{v} \sim \mathcal{G}_\theta} [D(\mathbf{v})] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mu}} \left[ \left( \|\partial D(\tilde{\mathbf{x}})\| - 1 \right)^2 \right] \right\} \right\}$$

- **Discriminator:**  $D: X \rightarrow \mathbb{R}$  is a measurable function.
- **Generator:**  $\mathcal{G}_\theta$  is a probability distribution on  $X$ .
- **Desired property:**  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) \mathbf{v}$  with  $\epsilon \sim U(0, 1)$ .

Issue: How to ensure  $D$  is 1-Lipschitz?

⇒ softly enforce this condition by adding a penalty term.

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{D: X \rightarrow \mathbb{R}} \left\{ \mathbb{E}_{\mathbf{x} \sim \mu} [D(\mathbf{x})] - \mathbb{E}_{\mathbf{v} \sim \mathcal{G}_\theta} [D(\mathbf{v})] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mu}} \left[ \left( \| \partial D(\tilde{\mathbf{x}}) \| - 1 \right)^2 \right] \right\} \right\}$$

- Discriminator:  $D: X \rightarrow \mathbb{R}$  is a **measurable function**.
- Generator:  $\mathcal{G}_\theta$  is a probability distribution on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) \mathbf{v}$  with  $\epsilon \sim U(0, 1)$ .

Issue: Unfeasible to maximise over all measurable mappings

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \mathbb{E}_{\mathbf{x} \sim \mu} [D_{\phi}(\mathbf{x})] - \mathbb{E}_{\mathbf{v} \sim \mathcal{G}_{\theta}} [D_{\phi}(\mathbf{v})] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mu}} \left[ \left( \|\partial D_{\phi}(\tilde{\mathbf{x}})\| - 1 \right)^2 \right] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \rightarrow \mathbb{R}$  a function parametrised by **neural network**.
- Generator:  $\mathcal{G}_{\theta}$  is a probability distribution on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) \mathbf{v}$  with  $\epsilon \sim U(0, 1)$ .

Issue: Unfeasible to maximise over all measurable mappings

⇒ Parametrise discriminator by a neural network.

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \mathbb{E}_{\mathbf{x} \sim \mu} [D_{\phi}(\mathbf{x})] - \mathbb{E}_{\mathbf{v} \sim \mathcal{G}_{\theta}} [D_{\phi}(\mathbf{v})] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mu}} \left[ \left( \|\partial D_{\phi}(\tilde{\mathbf{x}})\| - 1 \right)^2 \right] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \rightarrow \mathbb{R}$  a function parametrised by neural network.
- Generator:  $\mathcal{G}_{\theta}$  is a **probability distribution** on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}} \approx \mu$ .
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) \mathbf{v}$  with  $\epsilon \sim U(0, 1)$ .

Issue: Difficult to parametrise family of probability distributions on  $X$

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \mathbb{E}_{\mathbf{x} \sim \mu} [D_{\phi}(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim \sigma} [D_{\phi}(G_{\theta}(\mathbf{z}))] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mu}} \left[ \left( \|\partial D_{\phi}(\tilde{\mathbf{x}})\| - 1 \right)^2 \right] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \rightarrow \mathbb{R}$  a function parametrised by neural network.
- Generator:  $G_{\theta}: Z \rightarrow X$  a **function** parametrised by **neural network**.
- Desired property:  $G_{\hat{\theta}}(\mathbf{z})$  is approximately  $\mu$ -distributed when  $\mathbf{z} \sim \sigma$  ( $\sigma$  known).
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) G_{\theta}(\mathbf{z})$  with  $\epsilon \sim U(0, 1)$ .

Issue: Difficult to parametrise family of probability distributions on  $X$

⇒ Write generator as deterministic function with random input  $\mathbf{z} \sim \sigma$  with  $\sigma$  known.

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \mathbb{E}_{\mathbf{x} \sim \mu} [D_{\phi}(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim \sigma} [D_{\phi}(G_{\theta}(\mathbf{z}))] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mu}} \left[ \left( \|\partial D_{\phi}(\tilde{\mathbf{x}})\| - 1 \right)^2 \right] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \rightarrow \mathbb{R}$  a function parametrised by neural network.
- Generator:  $G_{\theta}: Z \rightarrow X$  a function parametrised by neural network.
- Desired property:  $G_{\hat{\theta}}(\mathbf{z})$  is approximately  $\mu$ -distributed when  $\mathbf{z} \sim \sigma$  ( $\sigma$  known).
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) G_{\theta}(\mathbf{z})$  with  $\epsilon \sim U(0, 1)$ .

Issue:  $\mu$  unknown  $\implies$  not possible to compute  $\mu$ - and  $\tilde{\mu}$ -expectations

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \frac{1}{m} \sum_{i=1}^m [D_{\phi}(x_i)] - \mathbb{E}_{z \sim \sigma} [D_{\phi}(G_{\theta}(z))] + \lambda \frac{1}{m} \sum_{i=1}^m [\left( \|\partial D_{\phi}(\tilde{x}_i)\| - 1 \right)^2] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \rightarrow \mathbb{R}$  a function parametrised by neural network.
- Generator:  $G_{\theta}: Z \rightarrow X$  a function parametrised by neural network.
- Desired property:  $G_{\hat{\theta}}(\mathbf{z})$  is approximately  $\mu$ -distributed when  $\mathbf{z} \sim \sigma$  ( $\sigma$  known).
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) G_{\theta}(\mathbf{z})$  with  $\epsilon \sim U(0, 1)$ .

Issue:  $\mu$  unknown  $\implies$  not possible to compute  $\mu$ - and  $\tilde{\mu}$ -expectations

$\implies$  Use empirical distributions given by  $x_i$  and  $\tilde{x}_i = \epsilon_i x_i + (1 - \epsilon_i) G_{\theta}(z_i)$ .

# Wasserstein GAN

Input: Unsupervised data  $x_i \in X$  generated by  $\mathbf{x} \sim \mu$ .

Task: Sample from  $\mathbf{x} \sim \mu$  where  $\mu$  is unknown.

Method: Learn how to sample from  $\mu$  by solving (Arjovsky et al., 2017)

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \frac{1}{m} \sum_{i=1}^m [D_{\phi}(x_i)] - \mathbb{E}_{z \sim \sigma} [D_{\phi}(G_{\theta}(z))] + \lambda \frac{1}{m} \sum_{i=1}^m [\left( \| \partial D_{\phi}(x_i) \| - 1 \right)^2] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \rightarrow \mathbb{R}$  a function parametrised by neural network.
- Generator:  $G_{\theta}: Z \rightarrow X$  a function parametrised by neural network.
- Desired property:  $G_{\hat{\theta}}(\mathbf{z})$  is approximately  $\mu$ -distributed when  $\mathbf{z} \sim \sigma$  ( $\sigma$  known).
- $\tilde{\mu}$  probability measure generated by  $\epsilon \mathbf{x} + (1 - \epsilon) G_{\theta}(\mathbf{z})$  with  $\epsilon \sim U(0, 1)$ .

Conclusion: Above expression is suitable for deep learning.

# Conditional Wasserstein GAN

Input: Supervised data  $(x_i, y_i) \in X \times Y$  generated by  $(\mathbf{x}, \mathbf{y}) \sim \mu$ .

Task: Sample from  $(\mathbf{x} | \mathbf{y} = y) \sim \pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$  where posterior is unknown.

Method: Learn how to sample from posterior by solving

$$\hat{\theta} \in \arg \min_{\theta} \mathbb{E}_{\mathbf{y}} \left[ \mathcal{W}(\mathcal{G}_{\theta}(\mathbf{y}), \pi_{\text{post}}(\mathbf{x} | \mathbf{y})) \right]$$

- $\mathcal{W}$  quantifies similarity between  $\mathcal{G}_{\theta}(y)$  and  $\pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$ .
- Generator:  $\mathcal{G}_{\theta}(y)$  is a probability distribution on  $X$ .
- Desired property:  $\mathcal{G}_{\hat{\theta}}(y) \approx \pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$ .

# Conditional Wasserstein GAN

**Input:** Supervised data  $(x_i, y_i) \in X \times Y$  generated by  $(\mathbf{x}, \mathbf{y}) \sim \mu$ .

**Task:** Sample from  $(\mathbf{x} | \mathbf{y} = y) \sim \pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$  where posterior is unknown.

**Method:** Learn how to sample from posterior by solving

$$\hat{\theta} \in \arg \min_{\theta} \mathbb{E}_{\mathbf{y}} \left[ \mathcal{W}(\mathcal{G}_{\theta}(\mathbf{y}), \pi_{\text{post}}(\mathbf{x} | \mathbf{y})) \right]$$

- $\mathcal{W}$  quantifies similarity between  $\mathcal{G}_{\theta}(y)$  and  $\pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$ .
- **Generator:**  $\mathcal{G}_{\theta}(y)$  is a probability distribution on  $X$ .
- **Desired property:**  $\mathcal{G}_{\hat{\theta}}(y) \approx \pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$ .

Repeating the same procedure as for Wasserstein GAN . . .

# Conditional Wasserstein GAN

Input: Supervised data  $(x_i, y_i) \in X \times Y$  generated by  $(\mathbf{x}, \mathbf{y}) \sim \mu$ .

Task: Sample from  $(\mathbf{x} | \mathbf{y} = y) \sim \pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$  where posterior is unknown.

Method: Learn how to sample from posterior by solving

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \frac{1}{m} \sum_{i=1}^m \left[ D_{\phi}(x_i, y_i) - \mathbb{E}_{\mathbf{z} \sim \sigma} [D_{\phi}(G_{\theta}(\mathbf{z}, y_i), y_i)] \right] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \times Y \rightarrow \mathbb{R}$ .
- Generator:  $G_{\theta}: Z \times Y \rightarrow X$ .
- Desired property:  $G_{\theta}(\mathbf{z}, y)$  approximately  $\pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$ -distributed when  $\mathbf{z} \sim \sigma$ .

# Conditional Wasserstein GAN

Input: Supervised data  $(x_i, y_i) \in X \times Y$  generated by  $(\mathbf{x}, \mathbf{y}) \sim \mu$ .

Task: Sample from  $(\mathbf{x} | \mathbf{y} = y) \sim \pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$  where posterior is unknown.

Method: Learn how to sample from posterior by solving

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \frac{1}{m} \sum_{i=1}^m \left[ D_{\phi}(x_i, y_i) - \mathbb{E}_{\mathbf{z} \sim \sigma} [D_{\phi}(G_{\theta}(\mathbf{z}, y_i), y_i)] \right] \right\} \right\}$$

- Discriminator:  $D_{\phi}: X \times Y \rightarrow \mathbb{R}$ .
- Generator:  $G_{\theta}: Z \times Y \rightarrow X$ .
- Desired property:  $G_{\theta}(\mathbf{z}, y)$  approximately  $\pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$ -distributed when  $\mathbf{z} \sim \sigma$ .
- Issue: Only single sample  $x_i \in X$  for each  $y_i \in Y$  in training data  
 $\implies G_{\hat{\theta}}(\mathbf{z}, y_i) \approx \delta_{x_i}$  independent of  $\mathbf{z} \sim \sigma$  (mode collapse).

# Conditional Wasserstein GAN

**Input:** Supervised data  $(x_i, y_i) \in X \times Y$  generated by  $(\mathbf{x}, \mathbf{y}) \sim \mu$ .

**Task:** Sample from  $(\mathbf{x} | \mathbf{y} = y) \sim \pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$  where posterior is unknown.

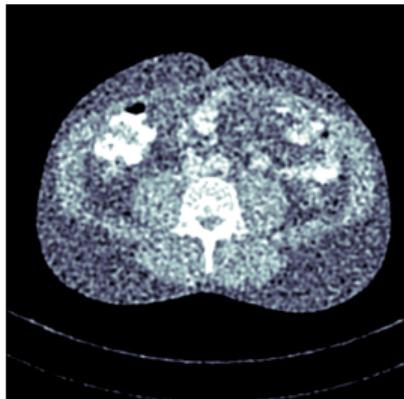
**Method:** Learn how to sample from posterior by solving

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \max_{\phi} \left\{ \frac{1}{m} \sum_{i=1}^m \left[ D_{\phi}(x_i, y_i) - \mathbb{E}_{\mathbf{z} \sim \sigma} [D_{\phi}(G_{\theta}(\mathbf{z}, y_i), y_i)] \right] \right\} \right\}$$

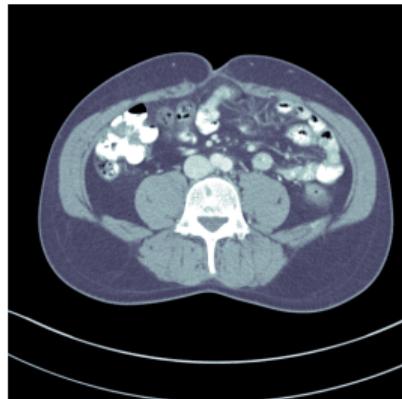
- **Discriminator:**  $D_{\phi}: X \times Y \rightarrow \mathbb{R}$ .
- **Generator:**  $G_{\theta}: Z \times Y \rightarrow X$ .
- **Desired property:**  $G_{\theta}(\mathbf{z}, y)$  approximately  $\pi_{\text{post}}(\mathbf{x} | \mathbf{y} = y)$ -distributed when  $\mathbf{z} \sim \sigma$ .
- **Issue:** Only single sample  $x_i \in X$  for each  $y_i \in Y$  in training data  
 $\implies G_{\hat{\theta}}(\mathbf{z}, y_i) \approx \delta_{x_i}$  independent of  $\mathbf{z} \sim \sigma$  (mode collapse).

**Solution:** Allow discriminator to distinguish between unordered pairs of model parameter or random samples generated by generator.

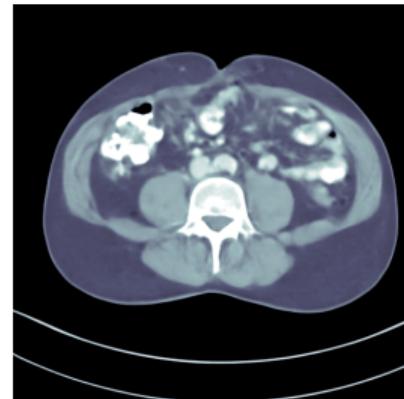
One of the images is the ground truth (phantom), can you figure out which one?



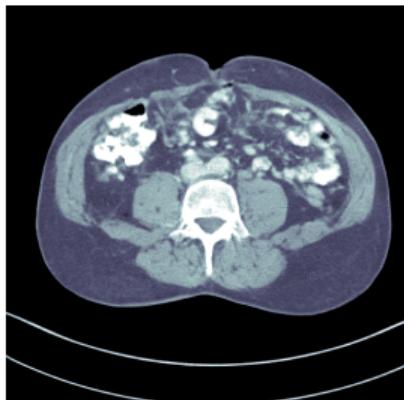
1



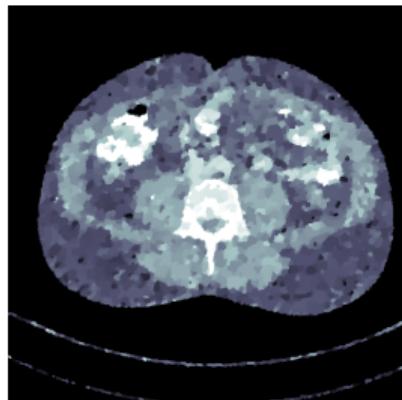
2



3

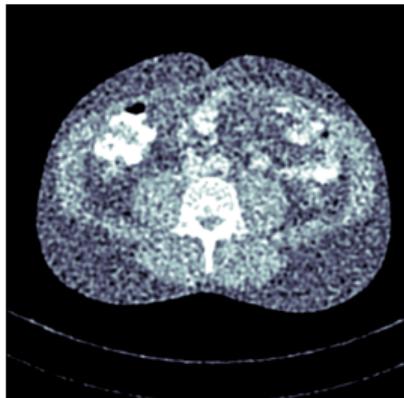


4

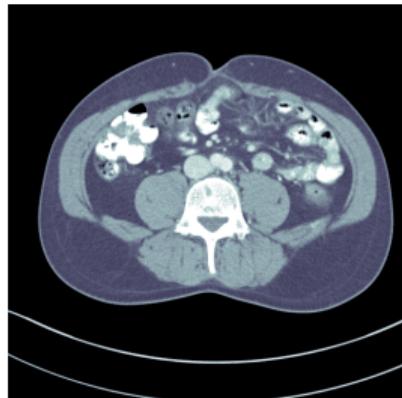


5

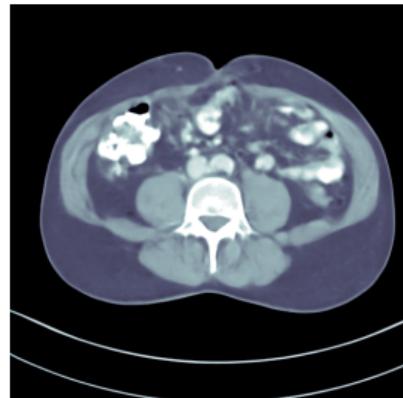
One of the images is the ground truth (phantom), can you figure out which one?



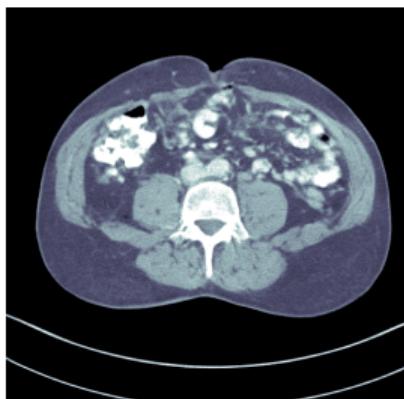
FBP



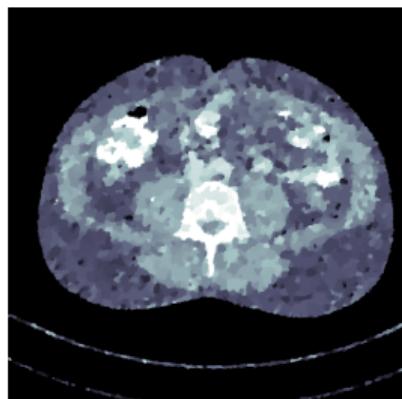
2



Conditional mean

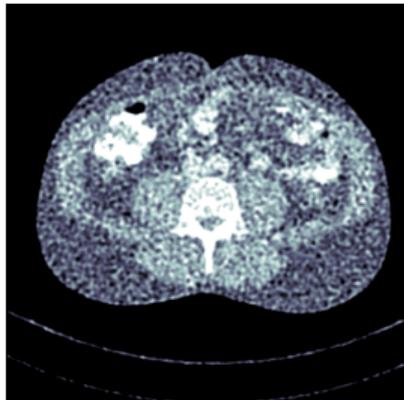


4

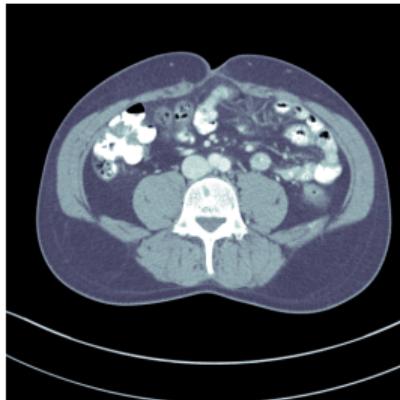


Total variation

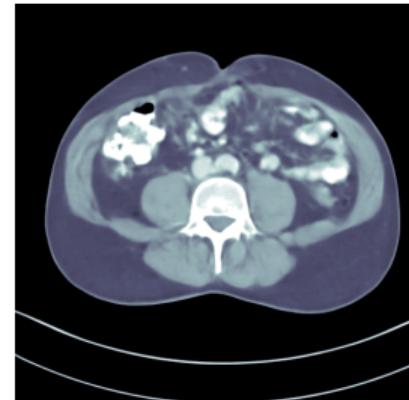
One of the images is the ground truth (phantom), can you figure out which one?



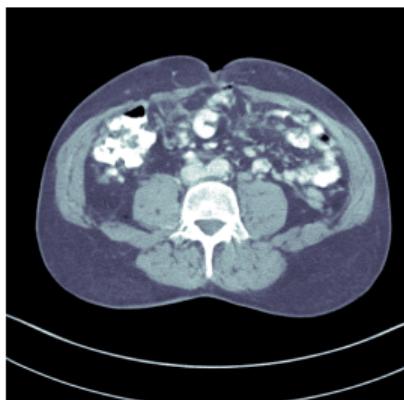
FBP



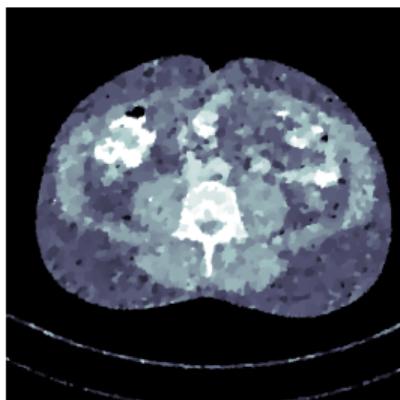
Phantom



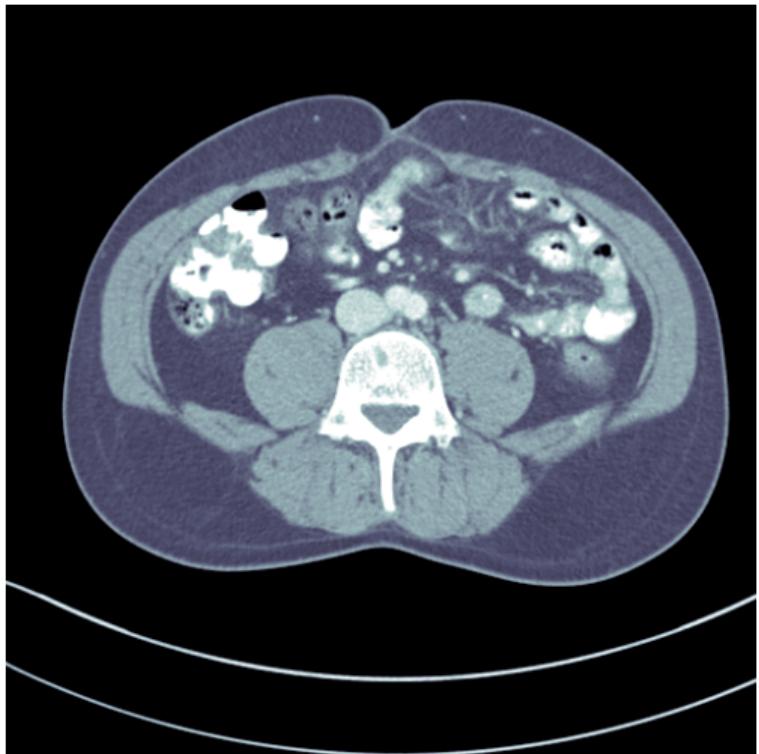
Conditional mean



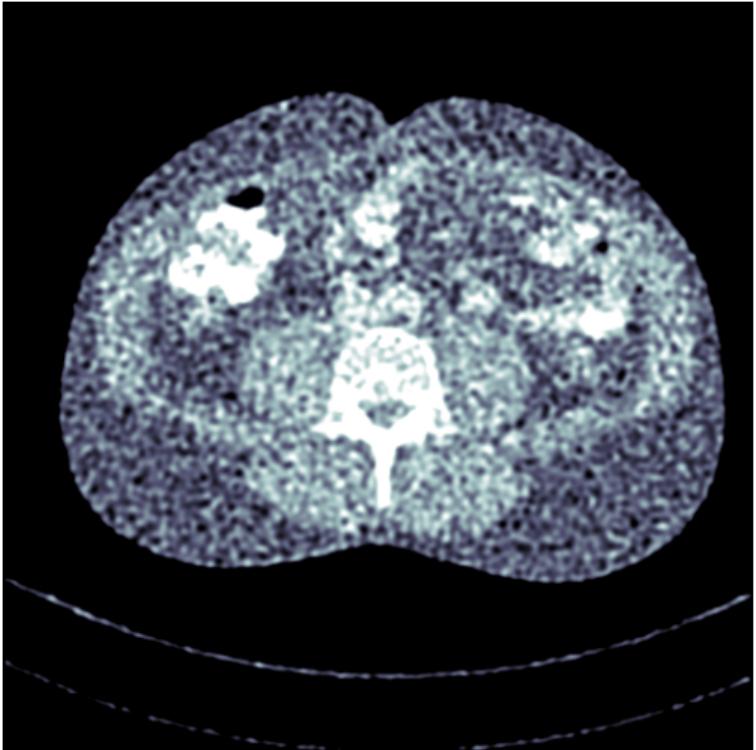
Deep posterior sample



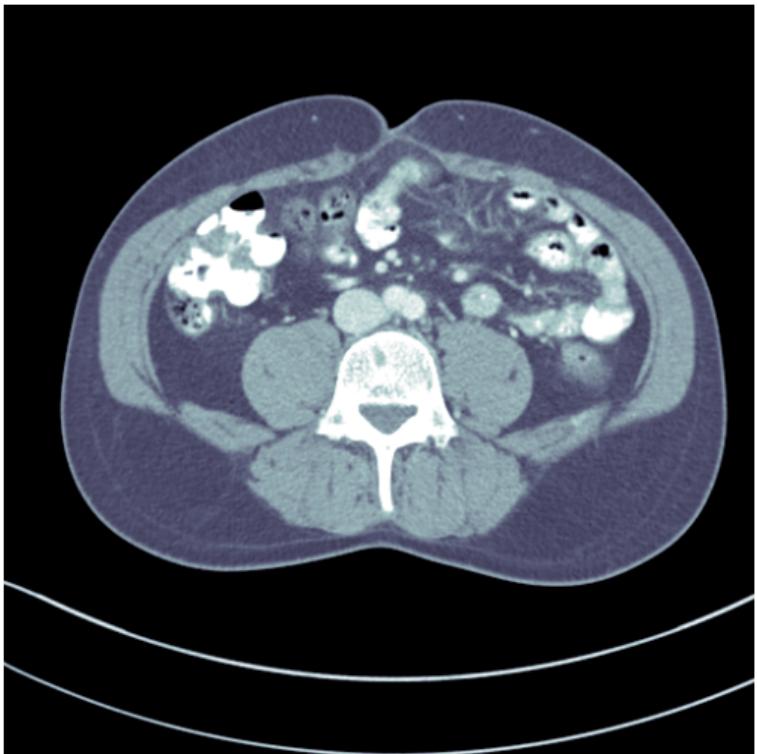
Total variation



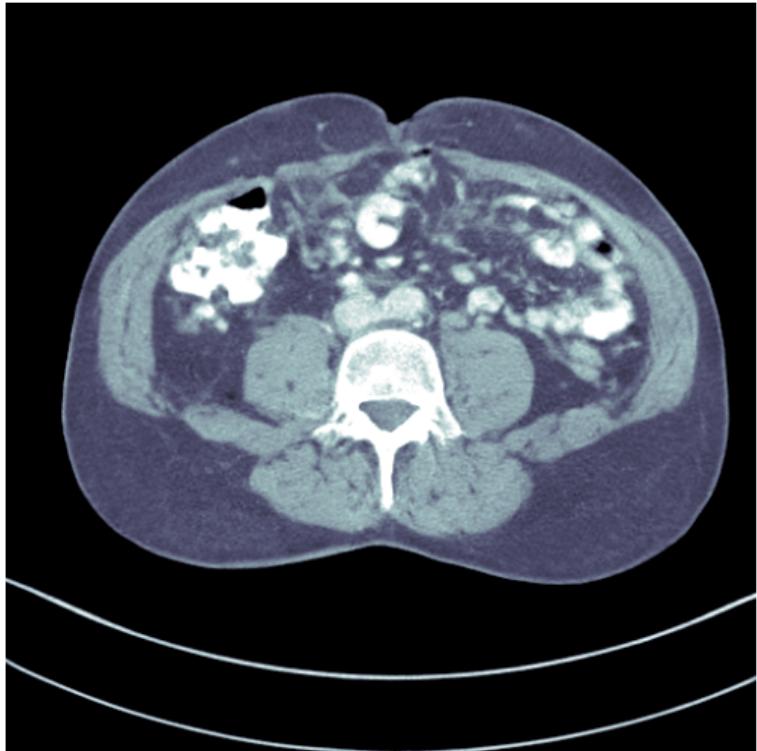
Phantom



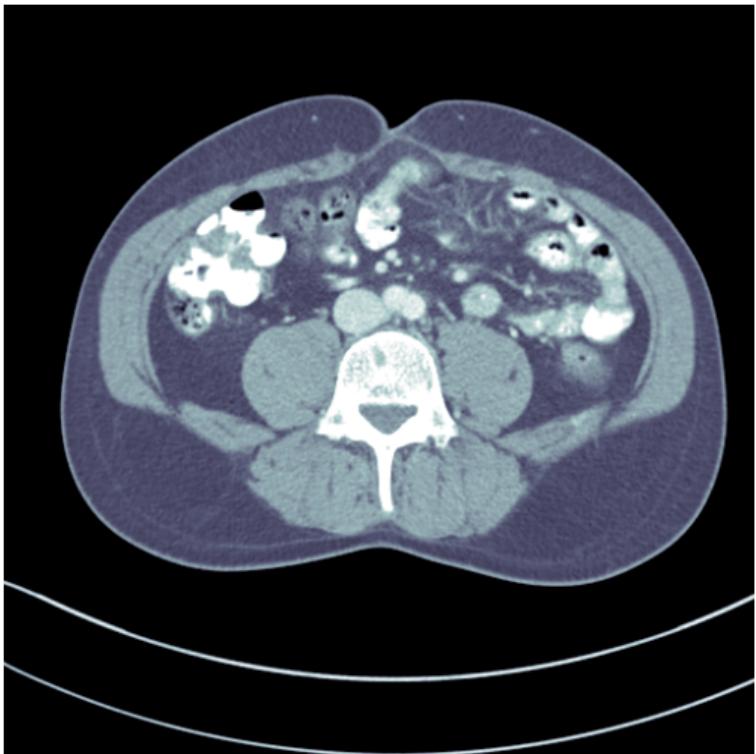
FBP reconstruction



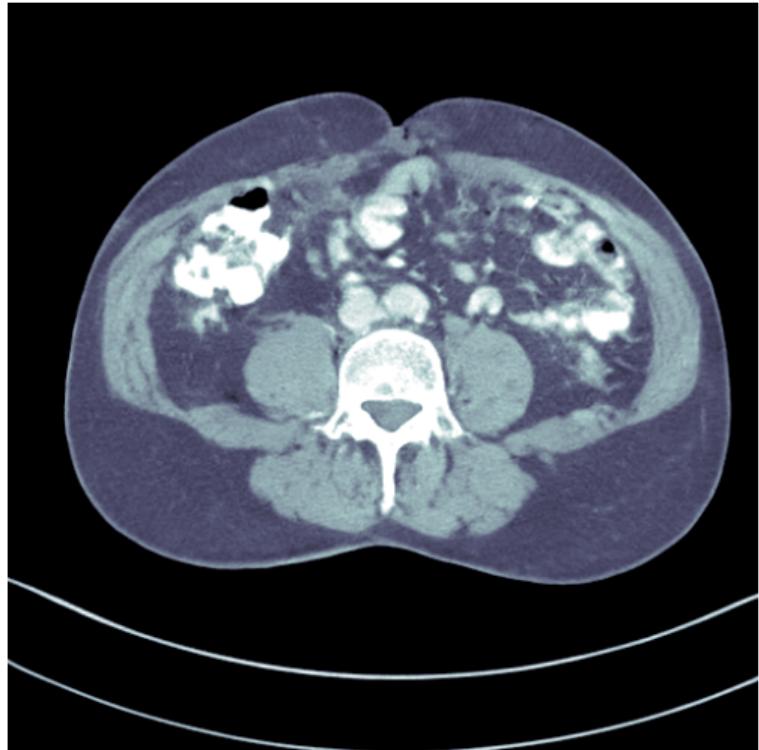
Phantom



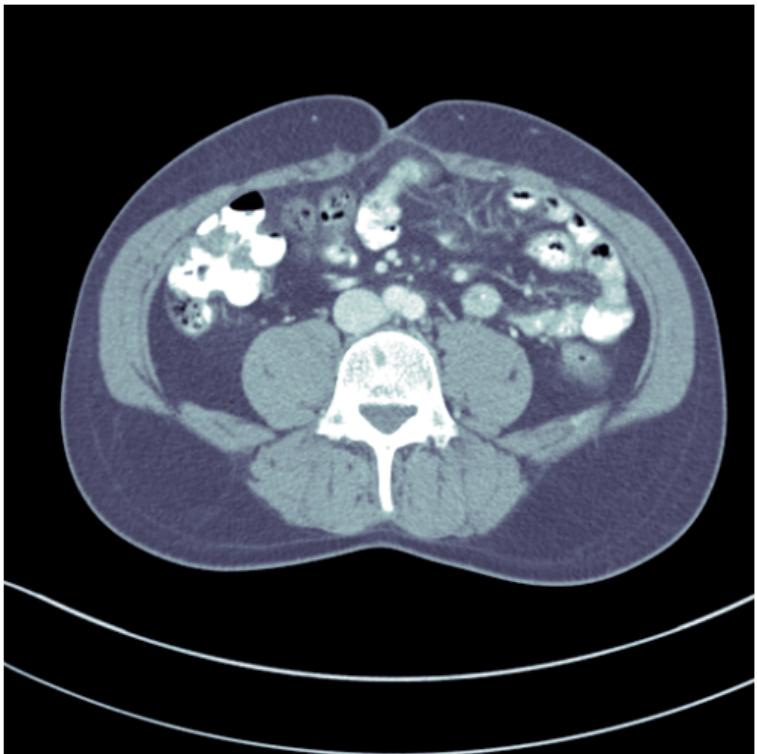
Posterior sample 1



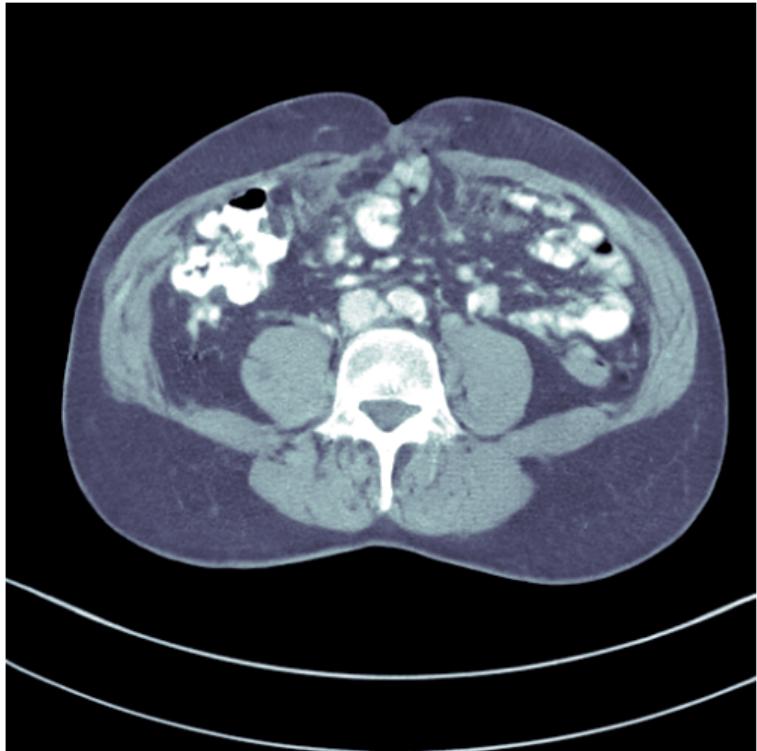
Phantom



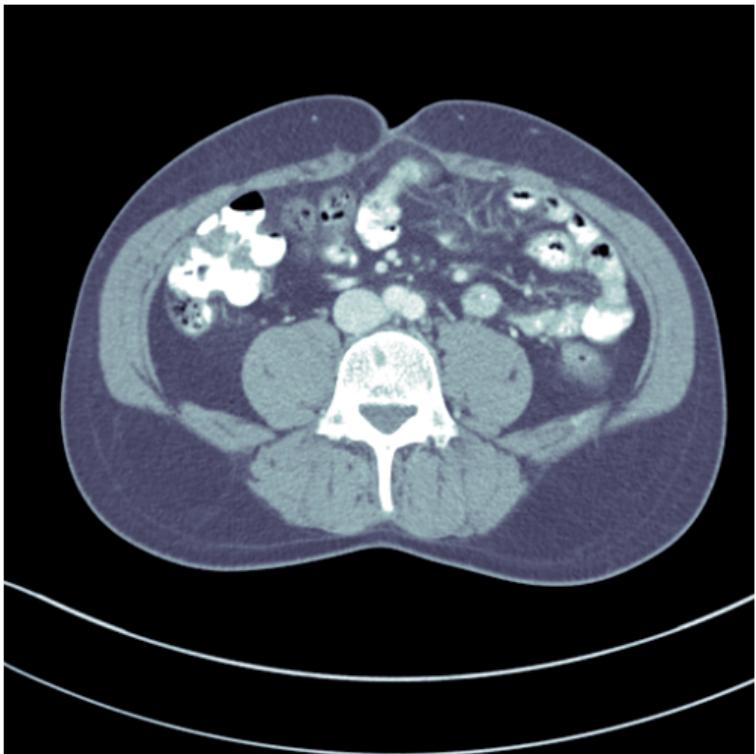
Posterior sample 2



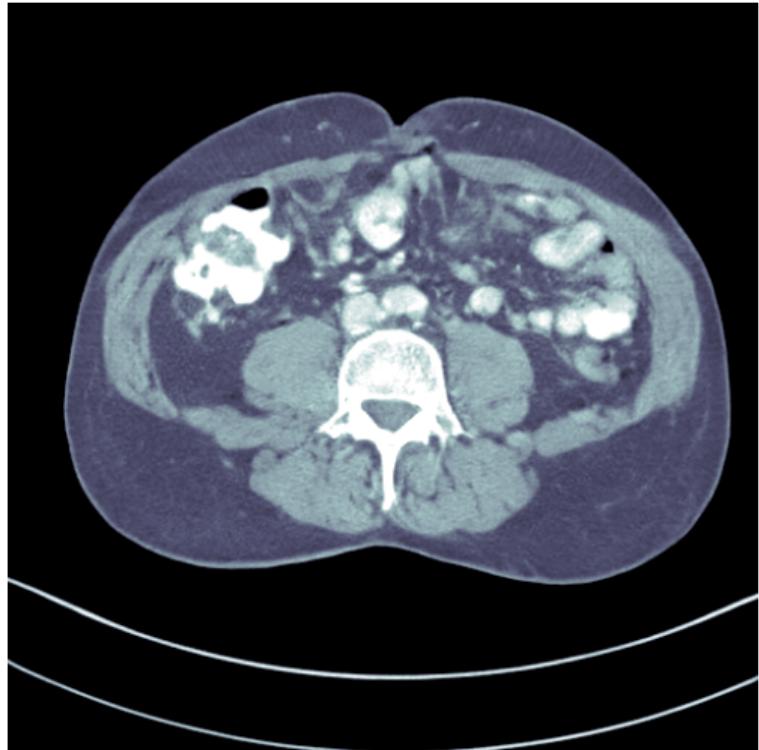
Phantom



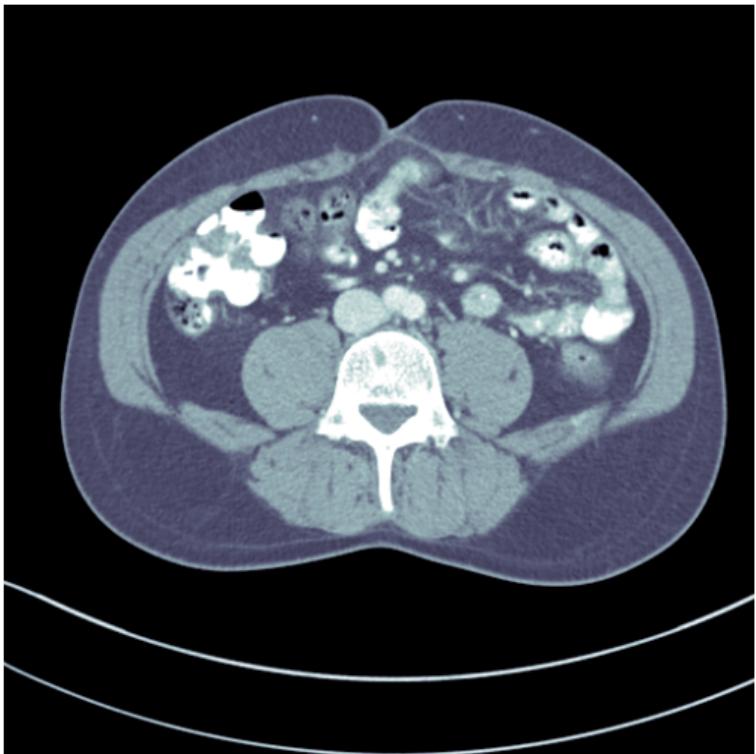
Posterior sample 3



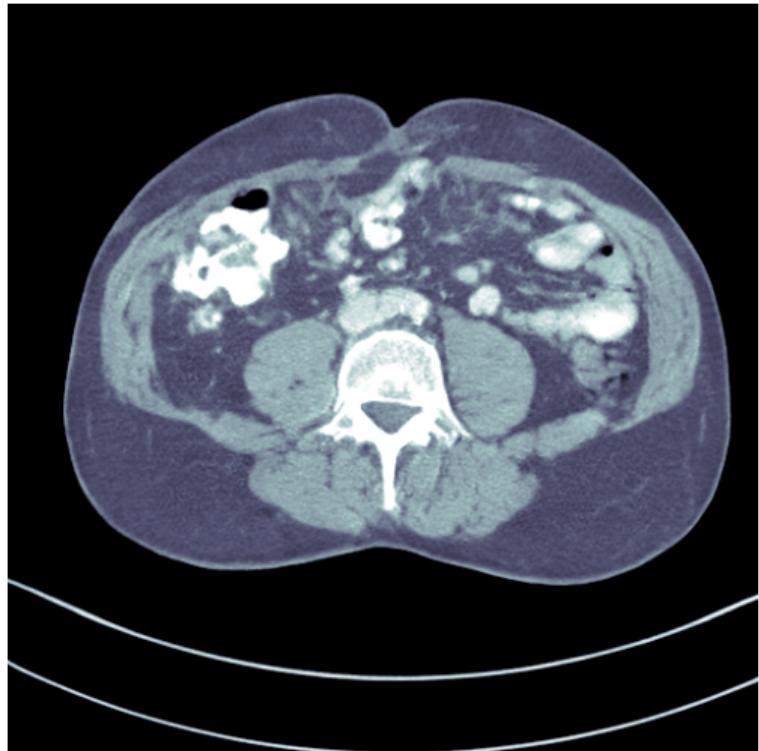
Phantom



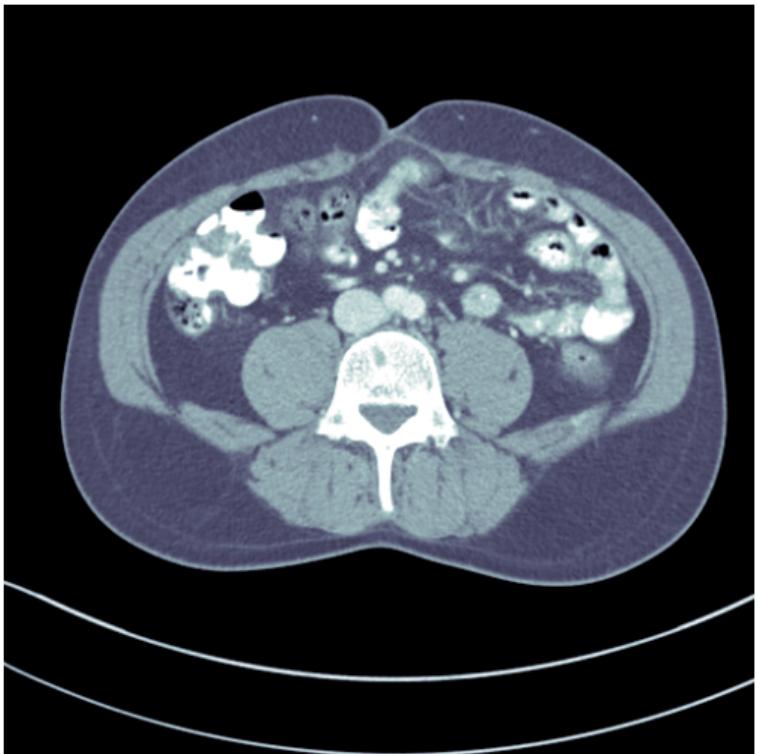
Posterior sample 4



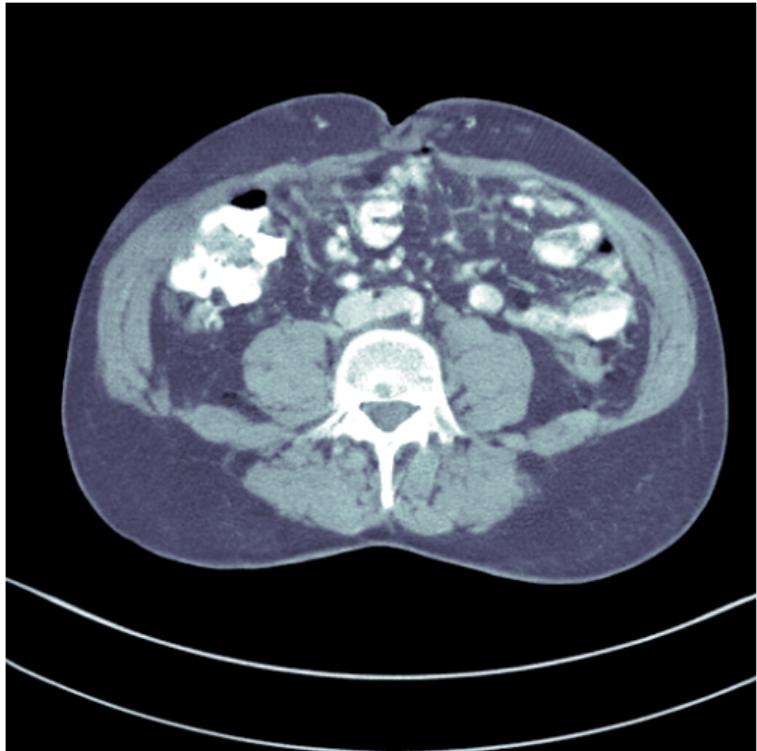
Phantom



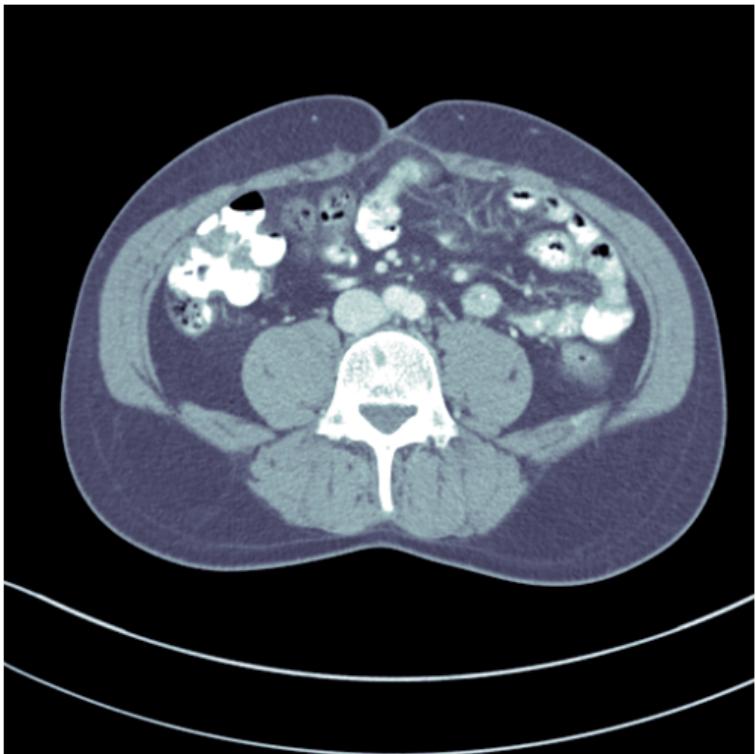
Posterior sample 5



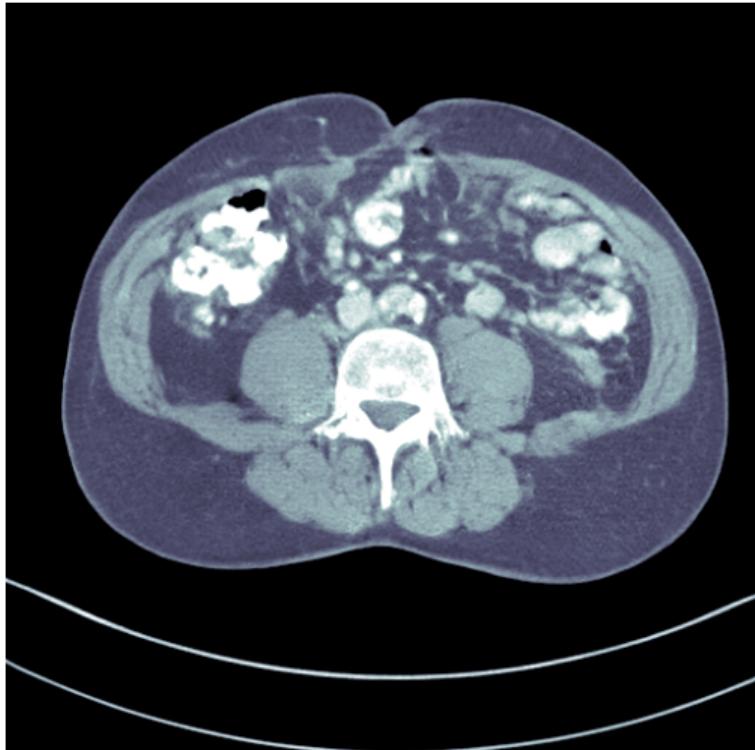
Phantom



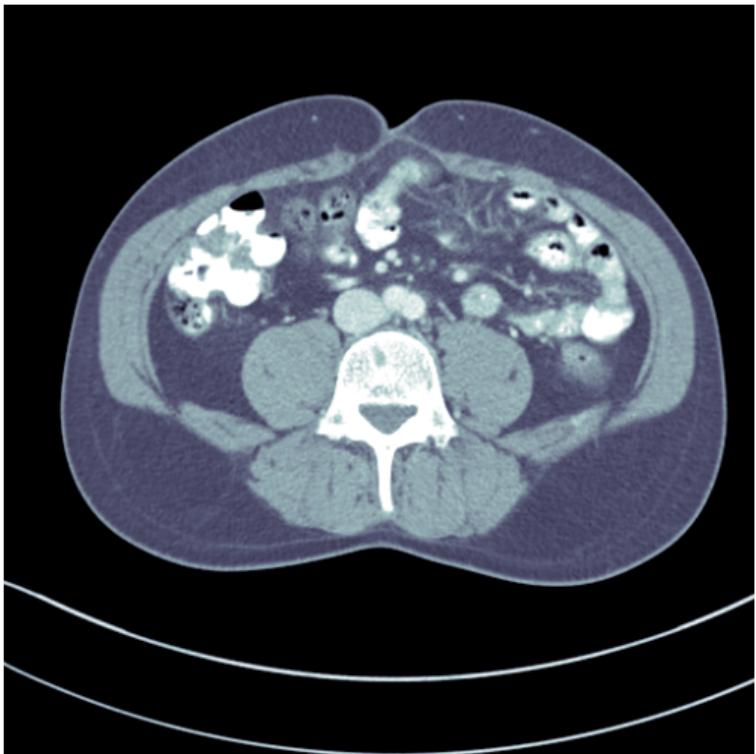
Posterior sample 6



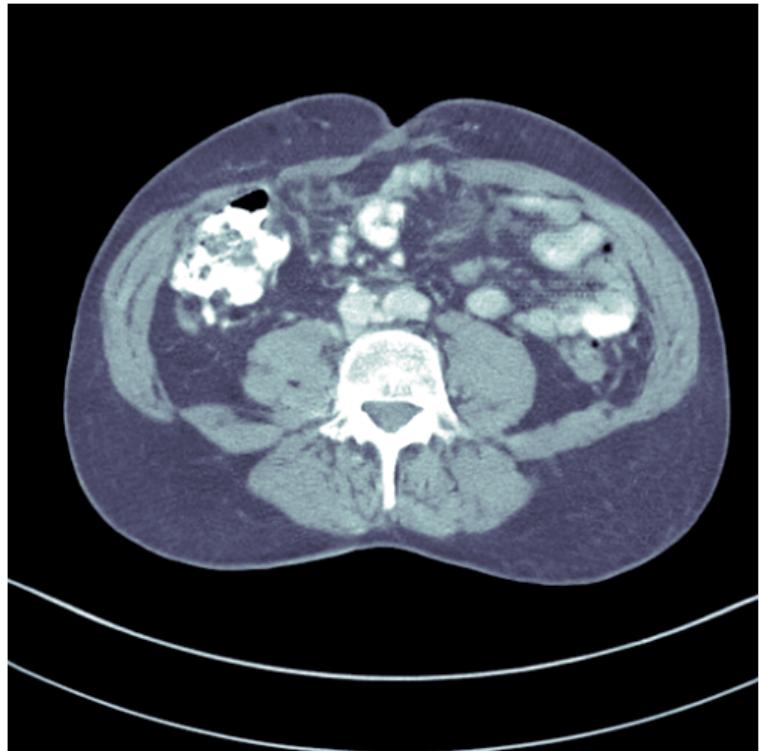
Phantom



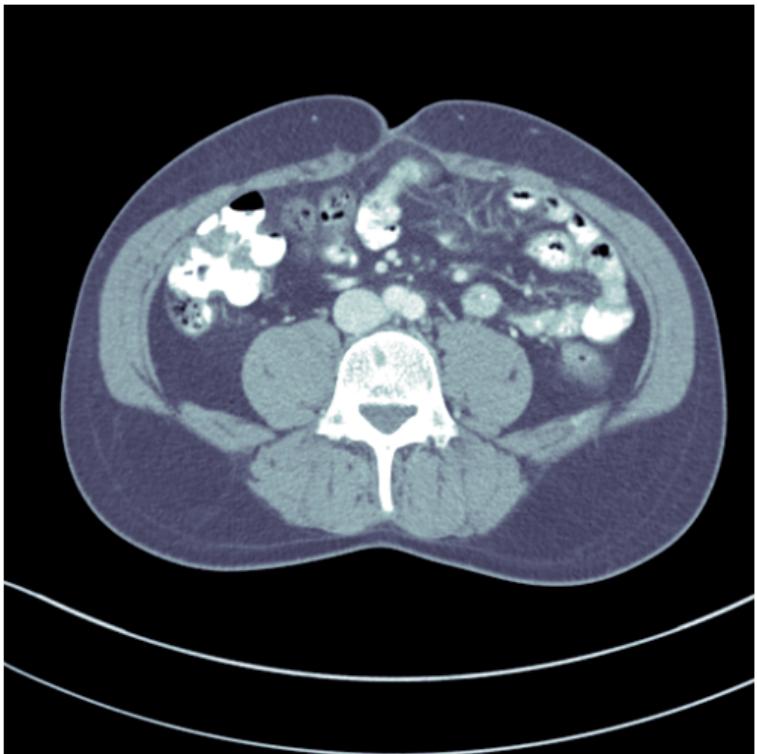
Posterior sample 7



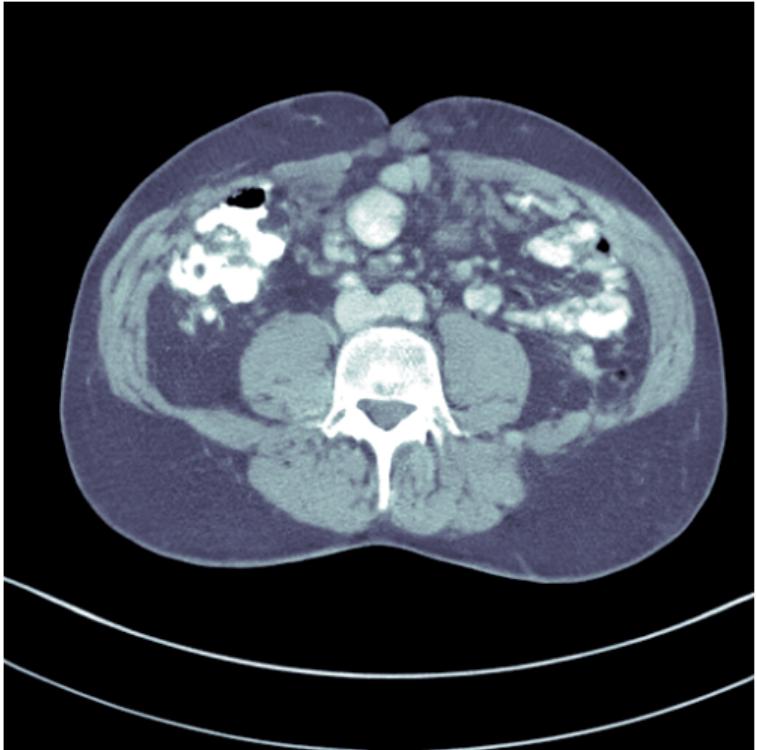
Phantom



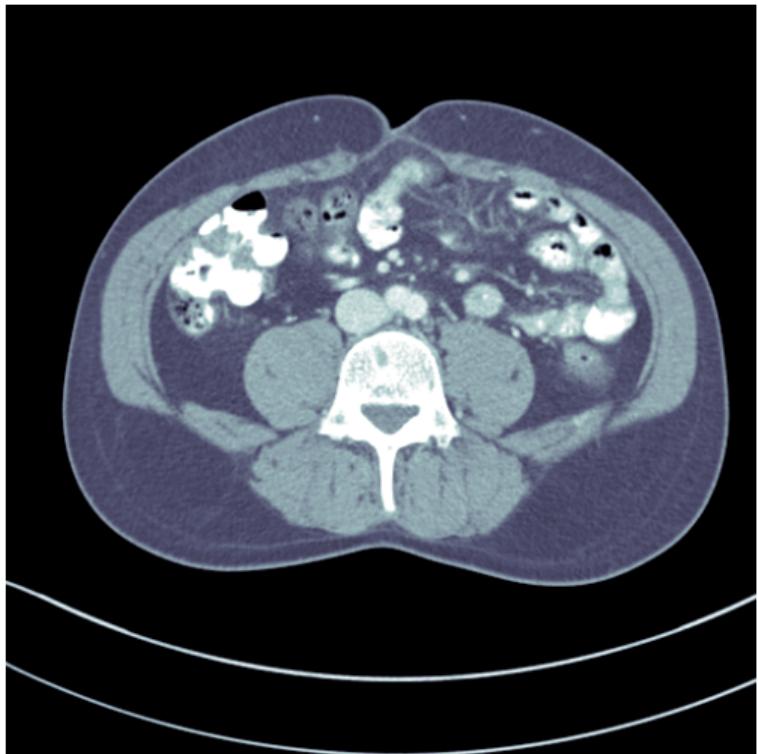
Posterior sample 8



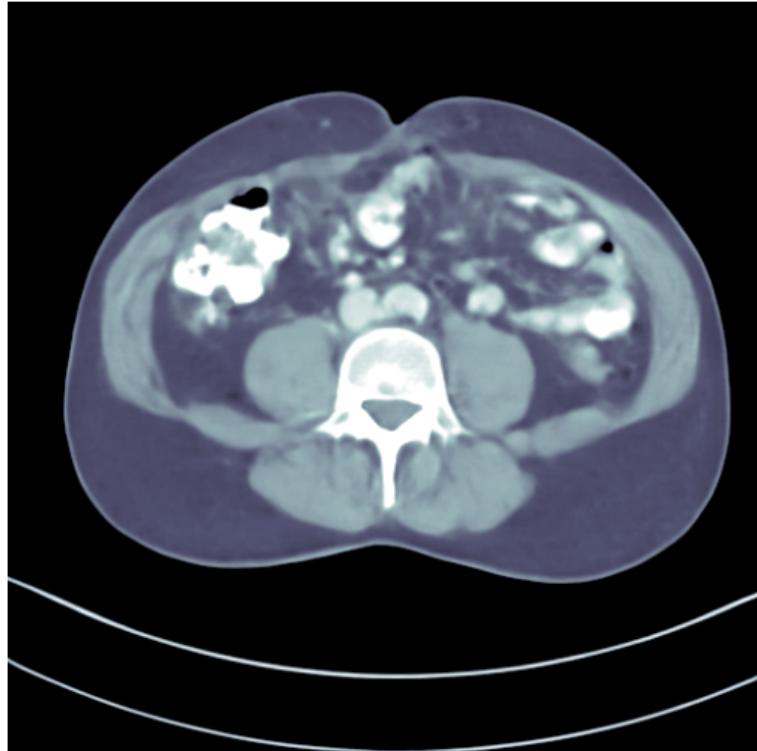
Phantom



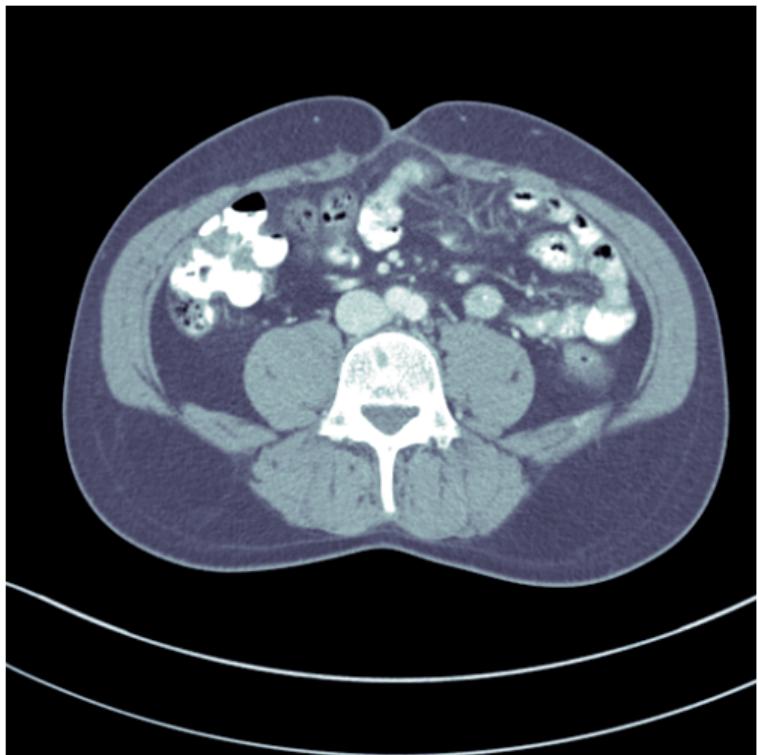
Posterior sample 9



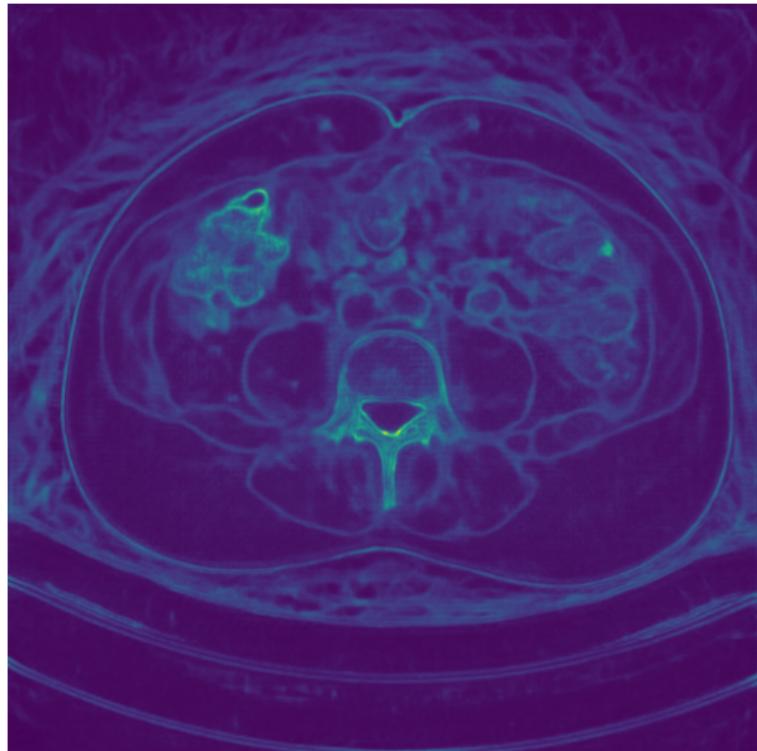
Phantom



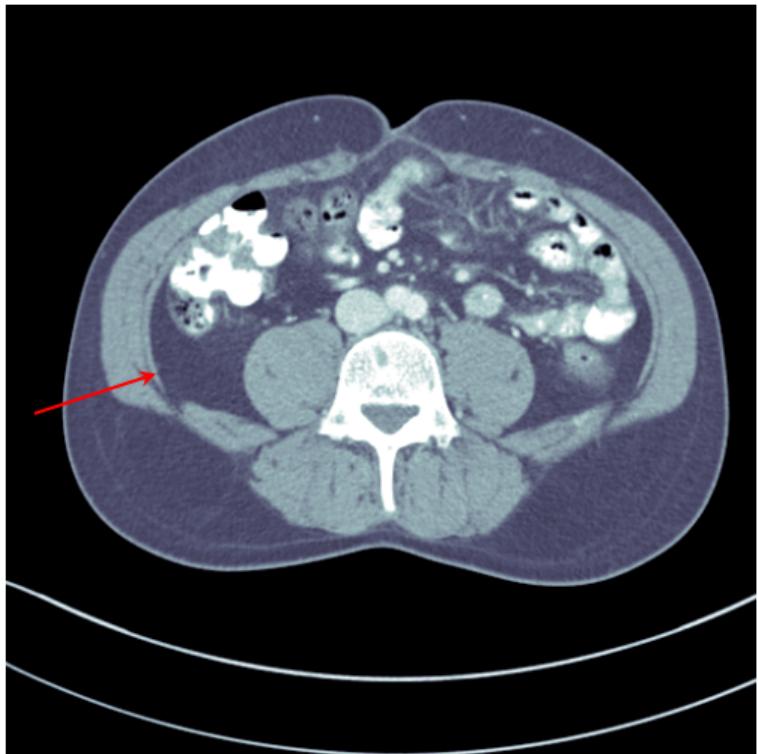
Conditional mean (1000 samples)



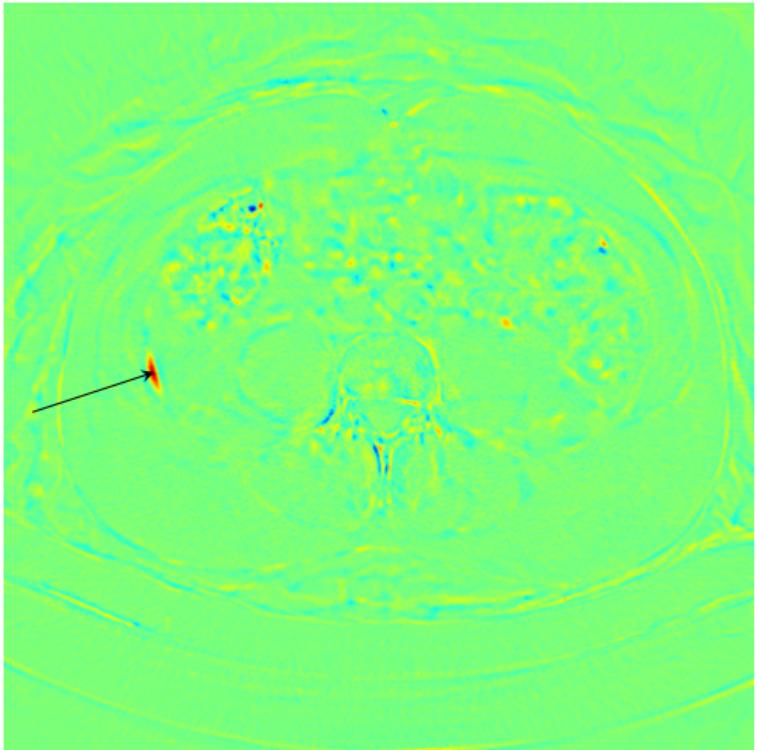
Phantom



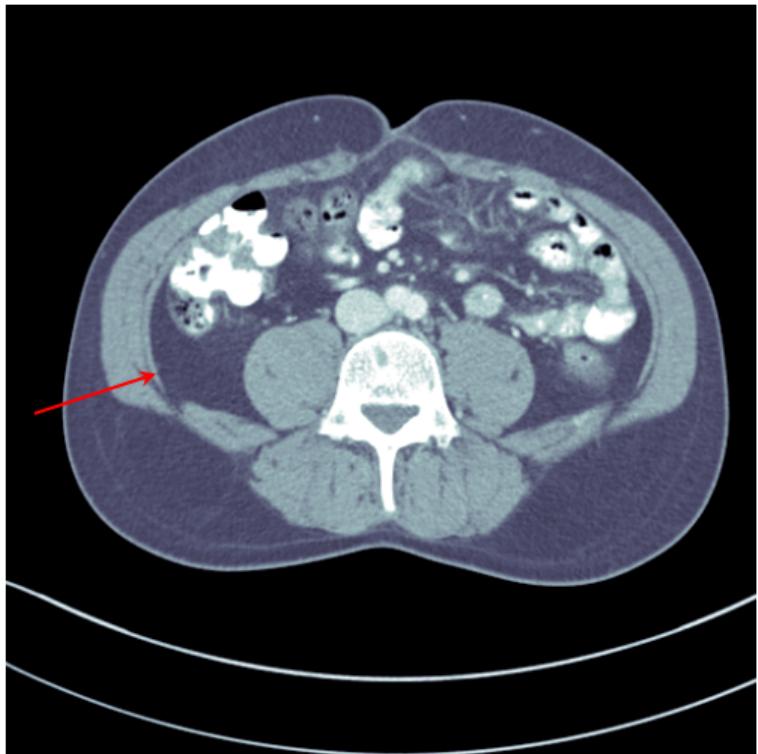
Standard deviation



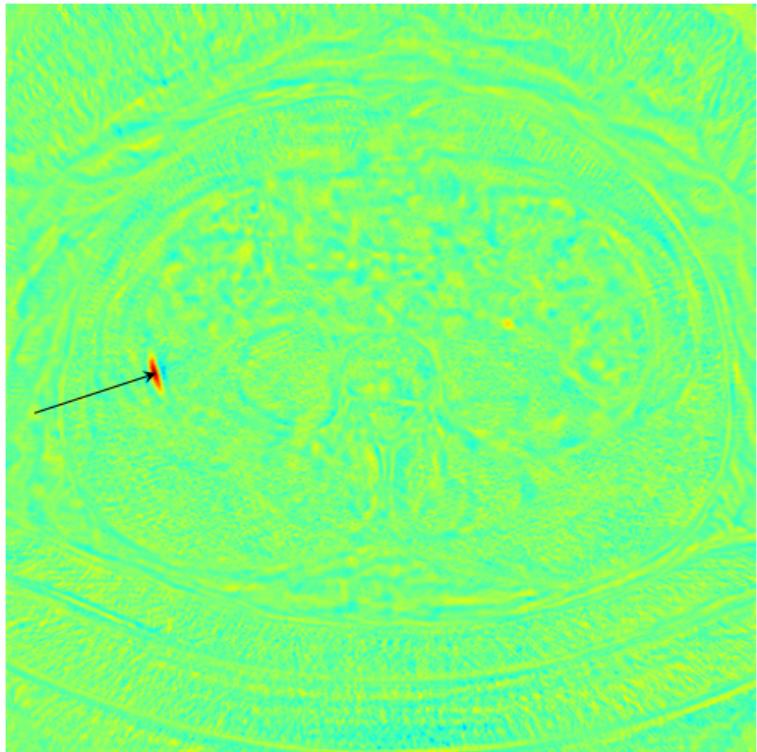
Phantom



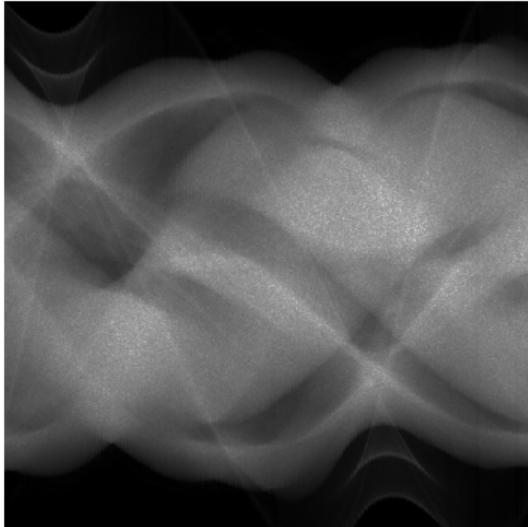
Covariance w.r.t. single point



Phantom



Correlation w.r.t. single point

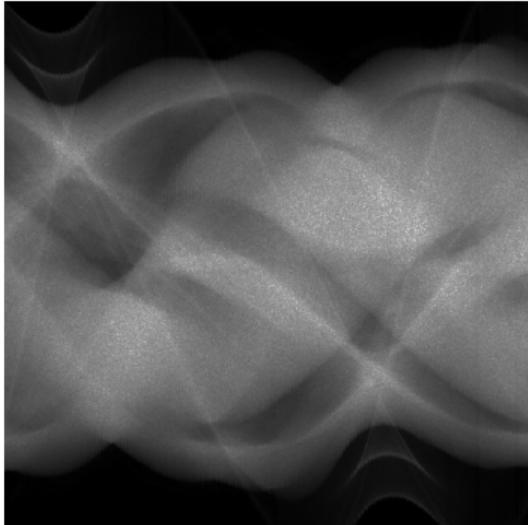


Data



FBP

- **Case:** Patient with suspected metastasis to the liver.
- **Data:** Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- **Liver lesion:**  $\triangle$  = difference in average contrast between ROI and liver.



Data

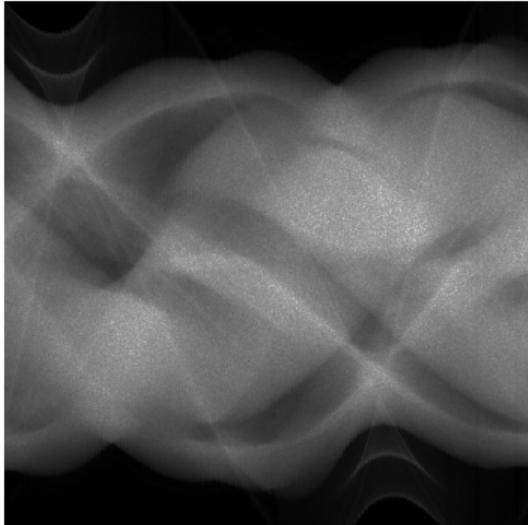


FBP

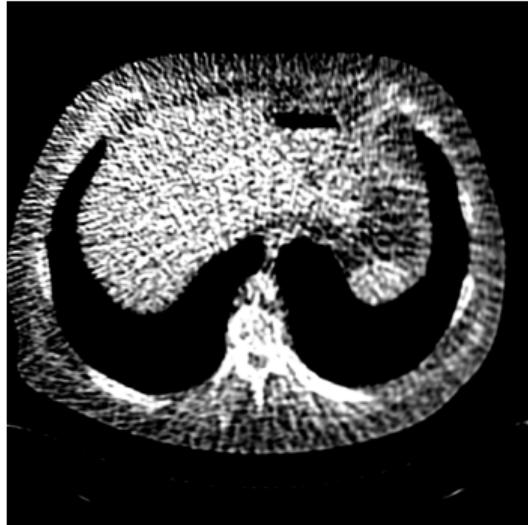


Posterior mean

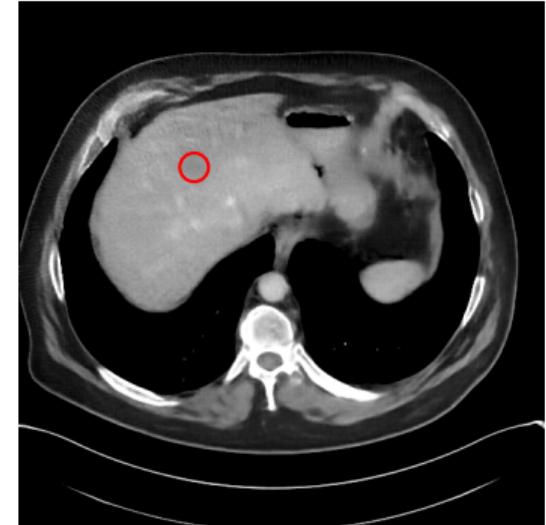
- **Case:** Patient with suspected metastasis to the liver.
- **Data:** Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- **Liver lesion:**  $\triangle$  = difference in average contrast between ROI and liver.



Data

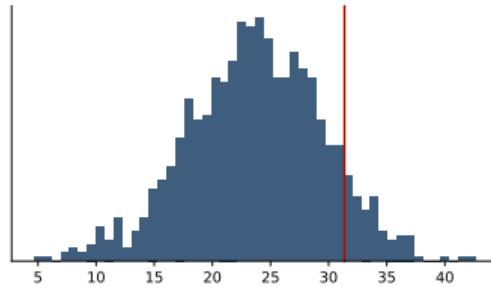


FBP

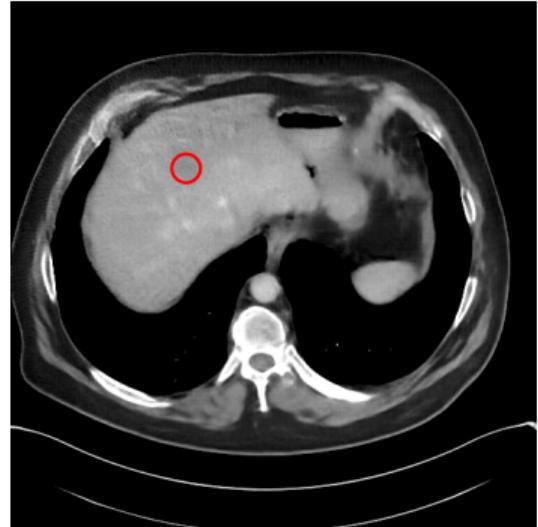


Posterior mean with ROI

- **Case:** Patient with suspected metastasis to the liver.
- **Data:** Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- **Liver lesion:**  $\triangle$  = difference in average contrast between ROI and liver.

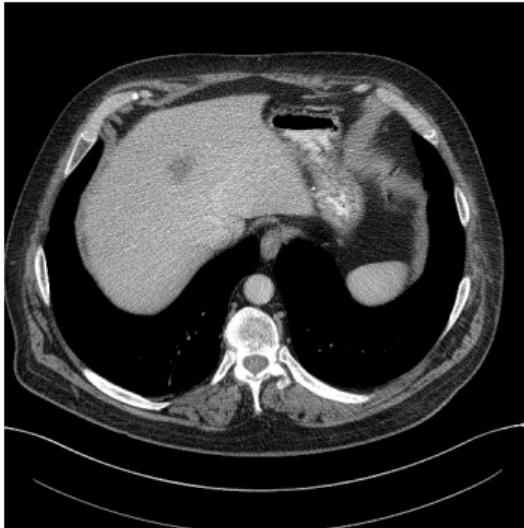


Histogram of  $\Delta$

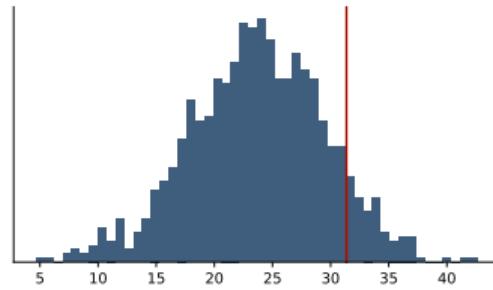


Posterior mean with ROI

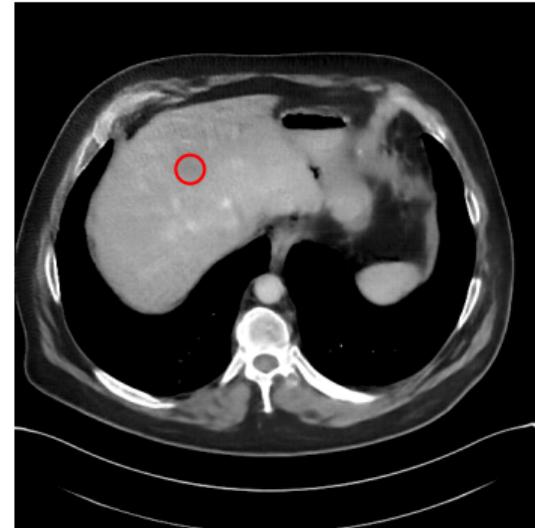
- **Case:** Patient with suspected metastasis to the liver.
- **Data:** Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- **Liver lesion:**  $\Delta =$  difference in average contrast between ROI and liver.
- **Hypothesis test:** Based on 1 000 samples, the ROI contains a lesion at 95%.



Normal dose image



Histogram of  $\Delta$



Posterior mean with ROI

- **Case:** Patient with suspected metastasis to the liver.
- **Data:** Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- **Liver lesion:**  $\Delta$  = difference in average contrast between ROI and liver.
- **Hypothesis test:** Based on 1 000 samples, the ROI contains a lesion at 95%.

# Deep posterior sampling

## Numerical performance

- Images:  $512 \times 512$  pixel.
- Fan beam data: 1000 angles, 1000 pixels/angle, additive noise (20% white noise).
- Generating a sample from the posterior takes about 20 ms on a ‘gaming’ PC.
- Training the generator and discriminator used 2000 pairs  $(x_i, y_i)$  from 9 patients.
- Generator and discriminator over-parametrised, about  $10^8$  parameters (generator has 2× parameters compared to discriminator)



Thank you for your attention!

# References

## References I

- Adler, J., Lunz, S., Verdier, O., Schönlieb, C.-B., & Öktem, O. (2018a). Task adapted reconstruction for inverse problems. *ArXiv e-print, cs.CV(1809.00948)*.
- Adler, J., Lunz, S., Verdier, O., Schönlieb, C.-B., & Öktem, O. (2018b). Task adapted reconstruction for inverse problems. In *NIPS 2018 meets medical imaging: Workshop within the 32nd annual conference on neural information processing systems (NIPS 2018)*. (Also available as ArXiv e-print:  
<https://arxiv.org/abs/1809.00948>)
- Adler, J., & Öktem, O. (2017). Solving ill-posed inverse problems using iterative deep neural networks. *Inverse Problems*, 33(12), 124007 (24pp). (Special issue: 100 Years of the Radon transform)
- Adler, J., & Öktem, O. (2018a). Deep Bayesian inversion: Computational uncertainty quantification for large scale inverse problems. *ArXiv e-print, stat.ML(1811.05910)*.

## References II

- Adler, J., & Öktem, O. (2018b). Learned primal-dual reconstruction. *IEEE Transactions on Medical Imaging*, 37(6), 1322–1332.
- Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein generative adversarial networks. *Proceedings of Machine Learning Research*, 70, 214–223.  
(Proceedings of the 34th International Conference on Machine Learning)
- Banert, S., Ringh, A., Adler, J., Karlsson, J., & Öktem, O. (2018). Data-driven nonsmooth optimization. *ArXiv*, 1808.00946.
- Dalca, A. V., Balakrishnan, G., Guttag, J., & Sabuncu, M. R. (2018). Unsupervised learning for fast probabilistic diffeomorphic registration. In F. A. F., J. A. Schnabel, C. Davatzikos, L. C. Alberola, & G. Fichtinger (Eds.), *21st international conference on medical image computing and computer assisted intervention MICCAI 2018* (Vol. 11070, pp. 729–738). Springer Verlag.

## References III

- Farabet, C., Couprie, C., Najman, L., & LeCun, Y. (2013). Learning hierarchical features for scene labeling. *IEEE transactions on pattern analysis and machine intelligence*, 35(8), 1915–1929.
- Ghosal, S., & Ray, N. (2017). Deep deformable registration: Enhancing accuracy by fully convolutional neural net. *Pattern Recognition Letters*, 94, 81–86.
- Giryes, R., Eldar, Y. C., Bronstein, A. M., & Sapiro, G. (2017). Tradeoffs between convergence speed and reconstruction accuracy in inverse problems. *ArXiv e-print, cs.NA(1605.09232)*.
- Gregor, K., & LeCun, Y. (2010). Learning fast approximations of sparse coding. In *Proceedings of the 27th International Conference on International Conference on Machine Learning (ICML '10)* (pp. 399–406).
- Guo, Y., Liu, Y., Georgiou, T., & Lew, M. S. (2018). A review of semantic segmentation using deep neural networks. *International Journal of Multimedia Information Retrieval*, 7(2), 87–93.

## References IV

- Hammernik, K., Klatzer, E., T. ans Kobler, Recht, M. P., Sodickson, D. K., Pock, T., & Knoll, F. (2018). Learning a variational network for reconstruction of accelerated MRI data. *Magnetic Resonance in Medicine*, 79(6), 3055–3071.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 770–778).
- Hsieh, J.-T., Zhao, S., Eismann, S., Mirabella, L., & Ermon, S. (2019). Learning neural PDE solvers with convergence guarantees. In *Seventh international conference on learning representations (ICLR 2019)*.
- Karpathy, A., & Fei-Fei, L. (2017). Deep visual-semantic alignments for generating image descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4), 664–676.

## References V

- Li, C. Y., Liang, X., Hu, Z., & Xing, E. P. (2018). Hybrid retrieval-generation reinforced agent for medical image report generation. *ArXiv e-print*, cs.CV(1805.08298).
- Li, H., Schwab, J., Antholzer, S., & Haltmeier, M. (2018). NETT: Solving inverse problems with deep neural networks. *ArXiv*, 1803.00092.
- Lunz, S., Öktem, O., & Schönlieb, C.-B. (2018). Adversarial regularizers in inverse problems. *ArXiv*, 1805.11572. (NIPS 2018)
- Mao, X., Li, Q., Xie, H., Lau, R. Y., Wang, Z., & Smolley, S. P. (2016). Least squares generative adversarial networks. *ArXiv*, 1611.04076.
- Mardani, M., Gong, E., Cheng, J. Y., Vasanawala, S. S., Zaharchuk, G., Xing, L., & Pauly, J. M. (2019). Deep generative adversarial neural networks for compressive sensing MRI. *IEEE Transactions on Medical Imaging*, 38(1), 167–179.
- Oymak, S., & Soltanolkotabi, M. (2017). Fast and reliable parameter estimation from nonlinear observations. *SIAM Journal on Optimization*, 27(4), 2276–2300.

## References VI

- Putzky, P., & Welling, M. (2017). *Recurrent inference machines for solving inverse problems*. Retrieved from <https://openreview.net/pdf?id=HkS01P91g>
- Rizzuti, G., Siahkoohi, A., & Herrmann, F. J. (2019). *Learned iterative solvers for the Helmholtz equation*. Submitted to 81st EAGE Conference & Exhibition 2019. (Available from <https://www.slim.eos.ubc.ca/content/learned-iterative-solvers-helmholtz-equation>)
- Romano, Y., Elad, M., & Milanfar, P. (2017). The little engine that could: Regularization by denoising (RED). *SIAM Journal on Imaging Sciences*, 10(4), 1804–1844.
- Romano, Y., Isidoro, J., & Milanfar, P. (2017). RAISR: Rapid and accurate image super resolution. *IEEE Transactions on Computational Imaging*, 3(1), 110–125.
- Schwab, J., Antholzer, S., & Haltmeier, M. (2018). Deep null space learning for inverse problems: Convergence analysis and rates. *ArXiv*, 1806.06137.

## References VII

- Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., & LeCun, Y. (2013). OverFeat: Integrated recognition, localization and detection using convolutional networks. *ArXiv e-print, cs.CV(1312.6229)*.
- Syu, N.-S., Chen, Y.-S., & Chuang, Y.-Y. (2018). Learning deep convolutional networks for demosaicing. *ArXiv e-print, cs.CV(1802.03769)*.
- Thoma, M. (2016). A survey of semantic segmentation. *ArXiv e-print, cs.CV(1602.06541)*.
- Wolterink, J. M., Dinkla, A. M., Savenije, M. H. F., Seevinck, P. R., van den Berg, C. A. T., & Işgum, I. (2017). Deep MR to CT synthesis using unpaired data. In S. Tsafaris, A. Gooya, A. Frangi, & J. Prince (Eds.), *Simulation and synthesis in medical imaging. SASHIMI 2017* (Vol. 10557, pp. 14–23).
- Xie, J., Xu, L., & Chen, E. (2012). Image denoising and inpainting with deep neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems (NIPS 2012)* (pp. 341–349).

## References VIII

Zhu, B., Liu, J. Z., Cauley, S. F., Rosen, B. R., & Rosen, M. S. (2018). Image reconstruction by domain-transform manifold learning. *Nature*, 555, 487–492.