

# Learned imaging with constraints and uncertainty quantification

Felix J. Herrmann, Ali Siahkoohi, and Gabrio Rizzuti. Learned imaging with constraints and uncertainty quantification. In: NeurIPS 2019 Deep Inverse Workshop. 2019

Presented by: Ali Siahkoohi



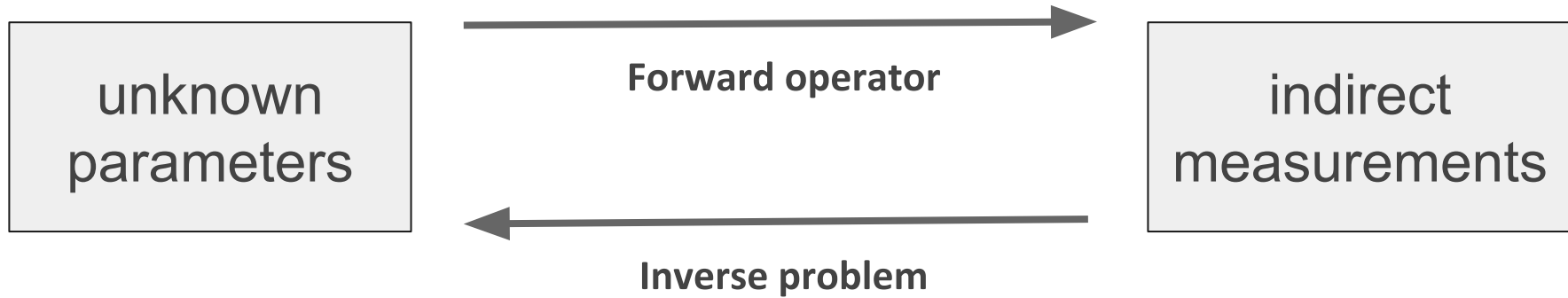
Georgia Institute of Technology

# Inverse problems

Estimate unknown parameters of a system via indirect measurements

- seismology: estimate the speed of sound in subsurface of the Earth
- medical imaging: infer visual representations of the interior of a body (X-ray radiography, MRI)

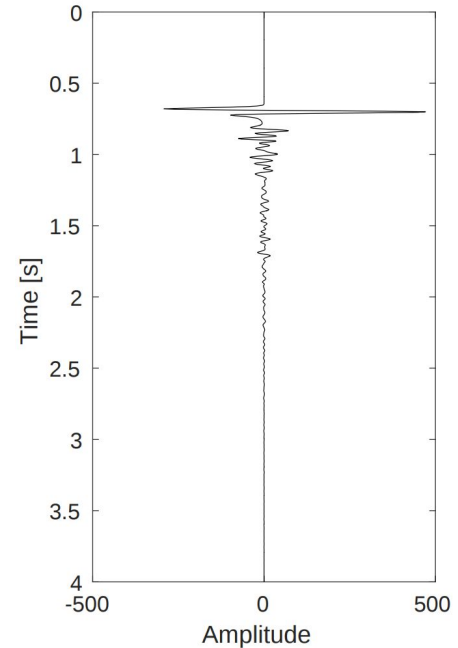
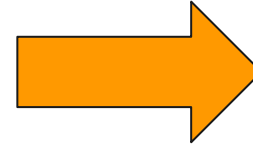
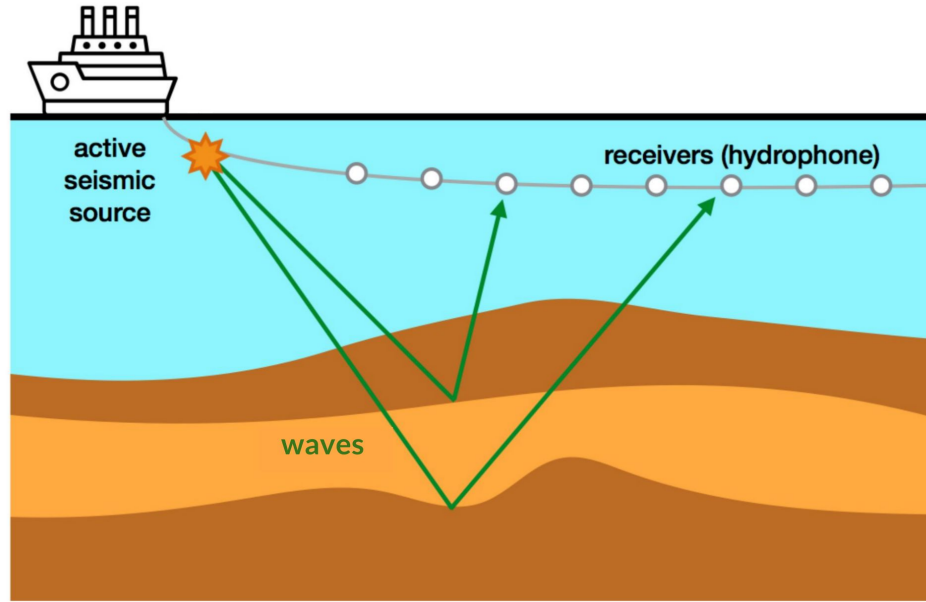
# Inverse problems



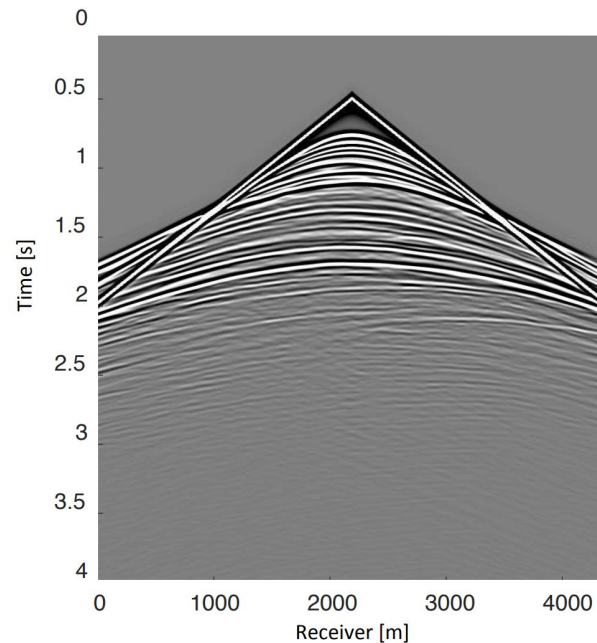
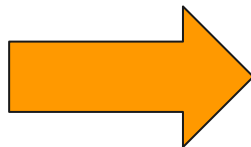
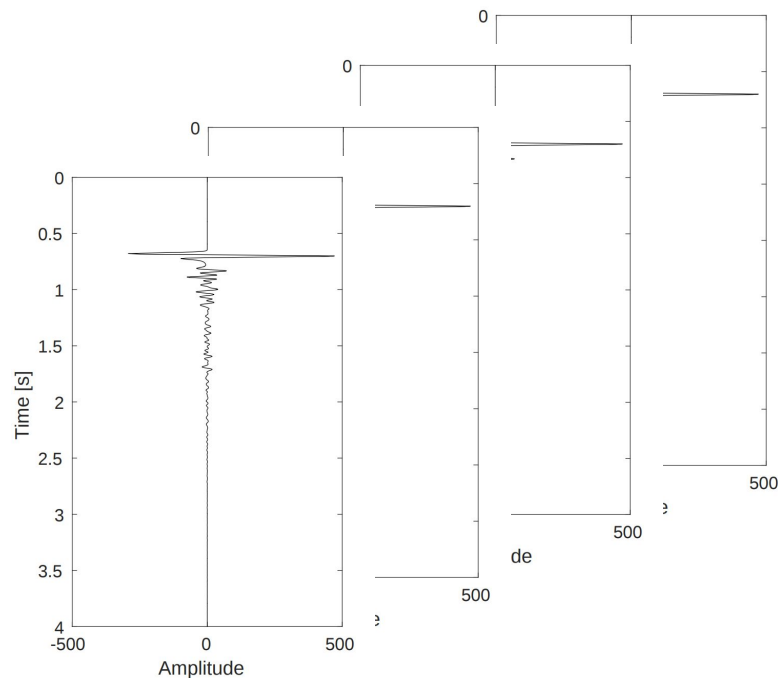
# Seismic imaging



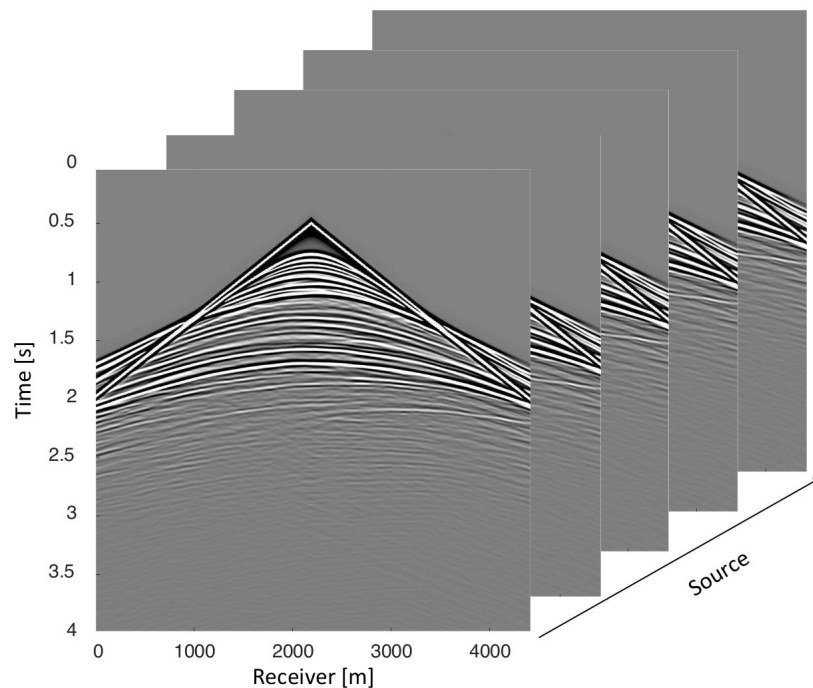
# Seismic data acquisition



# Shot record - one source location



# Seismic data volume

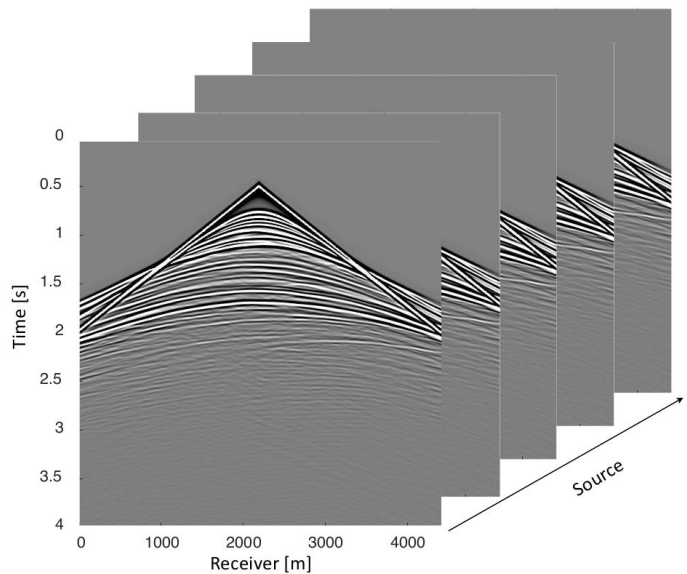


# Seismic imaging

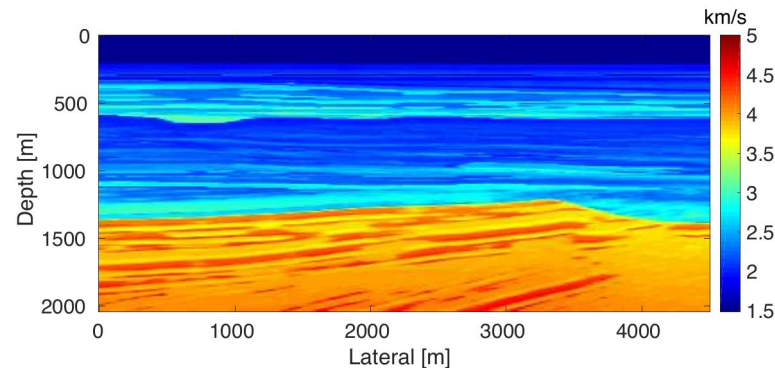
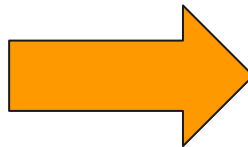




# Objective of exploration seismology



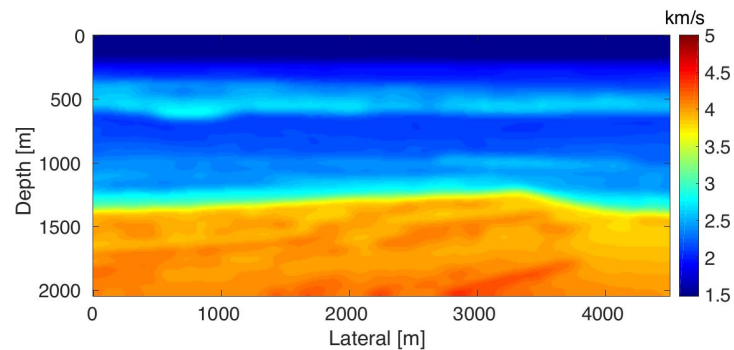
**Observed data**



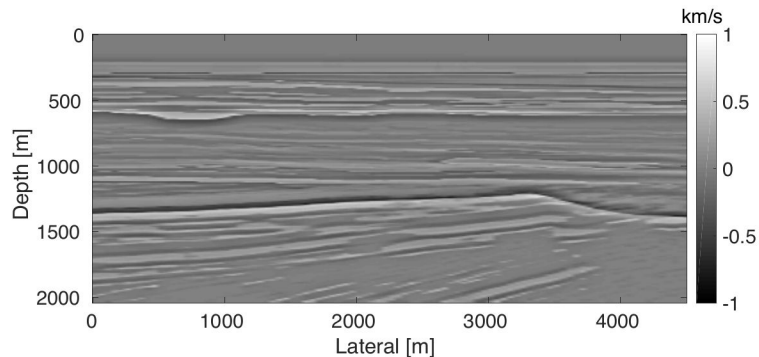
**Subsurface velocity structure**

Fang, Z., 2018. Source estimation and uncertainty quantification for wave-equation based seismic imaging and inversion (Doctoral dissertation, University of British Columbia).

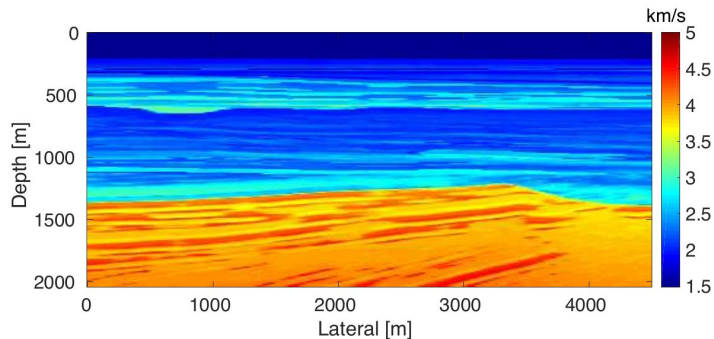
# Velocity model



**Long-wavelength structure  
(background velocity)**

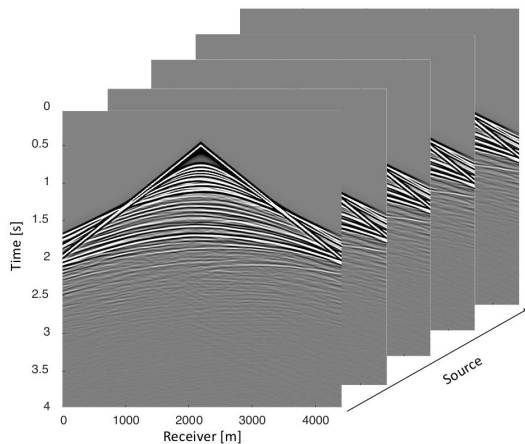


**Short-wavelength structure  
(velocity perturbation)**

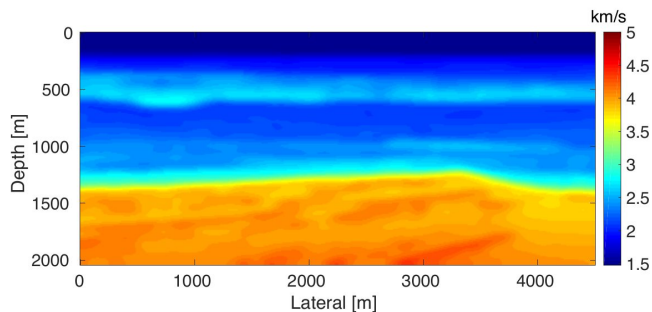
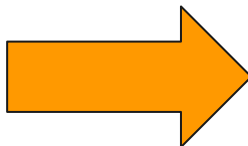


**Subsurface velocity structure**

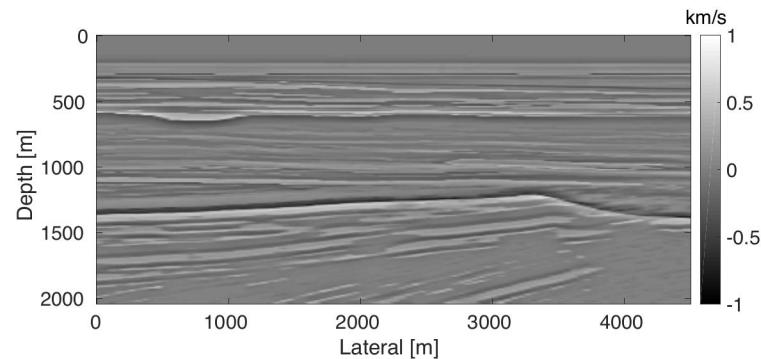
# Seismic imaging



**Observed data**



**Long-wavelength structure  
(background velocity)**

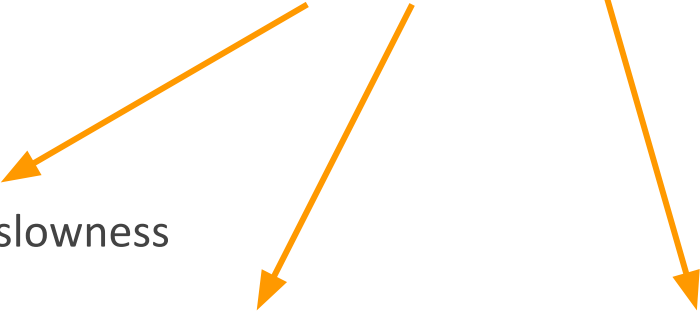


**Short-wavelength structure  
(velocity perturbation)**

# Seismic imaging



# Nonlinear forward operator

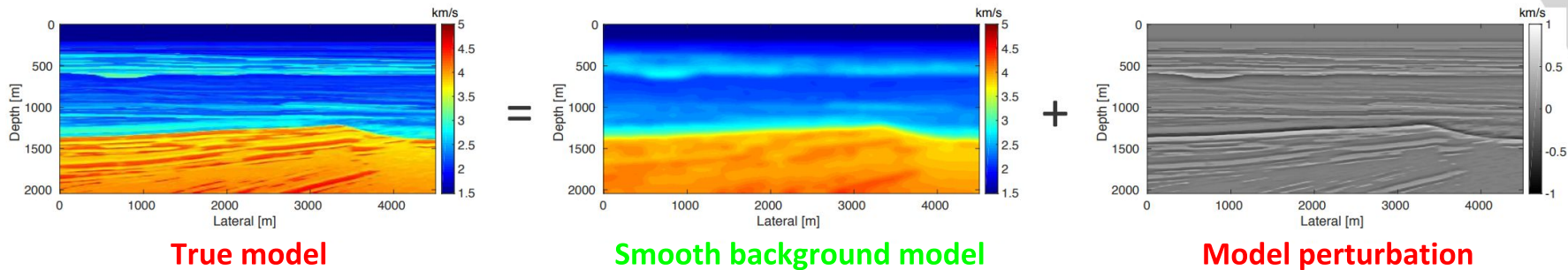
$$F(\mathbf{m}, \mathbf{q}) = \mathbf{P} \overbrace{\mathbf{A}(\mathbf{m})}^{\text{discretized wave-equation}}^{-1} \mathbf{q}$$


squared slowness

source

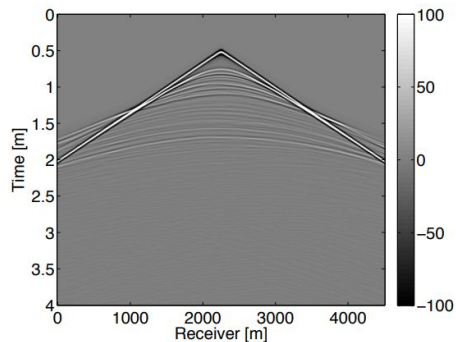
Restriction operator to the receivers location

# Taylor series expansion

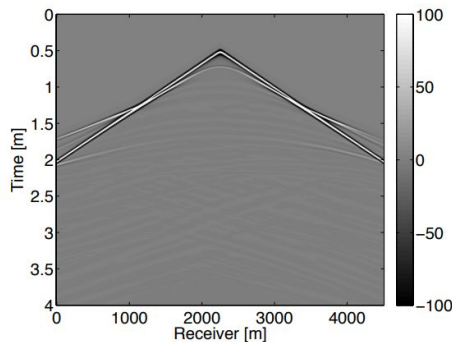


$$F(\underbrace{\mathbf{m}_0 + \delta\mathbf{m}}_{\text{True model}}, \mathbf{q}) = F(\underbrace{\mathbf{m}_0}_{\text{Smooth background model}}, \mathbf{q}) + \nabla F(\mathbf{m}_0, \mathbf{q})^T \underbrace{\delta\mathbf{m}}_{\text{Model perturbation}} + \mathcal{O}(\|\delta\mathbf{m}\|_2^2)$$

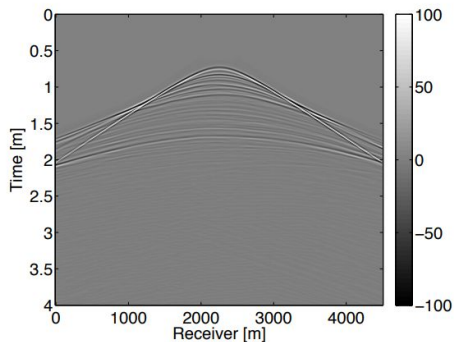
# Taylor series expansion



Observed data



Predicted data



Data residual

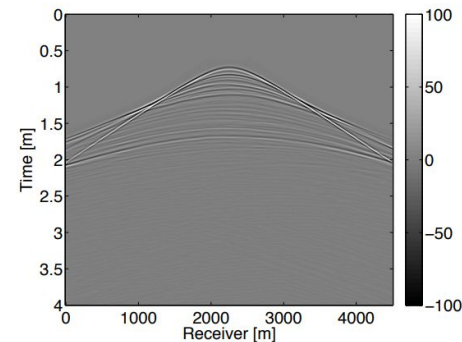
$$F(\mathbf{m}_0 + \delta\mathbf{m}, \mathbf{q}) = F(\mathbf{m}_0, \mathbf{q}) + \nabla F(\mathbf{m}_0, \mathbf{q})^T \delta\mathbf{m} + \mathcal{O}(\|\delta\mathbf{m}\|_2^2)$$

Diagram illustrating the Taylor series expansion of the forward model  $F(\mathbf{m}, \mathbf{q})$  around the initial model  $\mathbf{m}_0$ . The expansion is shown as the sum of three terms, each corresponding to one of the plots above:

- Observed data**:  $F(\mathbf{m}_0 + \delta\mathbf{m}, \mathbf{q})$
- Predicted data**:  $F(\mathbf{m}_0, \mathbf{q})$
- Data residual**:  $\nabla F(\mathbf{m}_0, \mathbf{q})^T \delta\mathbf{m} + \mathcal{O}(\|\delta\mathbf{m}\|_2^2)$

# Taylor series expansion

$$\nabla F(\mathbf{m}_0, \mathbf{q})^T \delta \mathbf{m} + \mathcal{O}(\|\delta \mathbf{m}\|_2^2) =$$

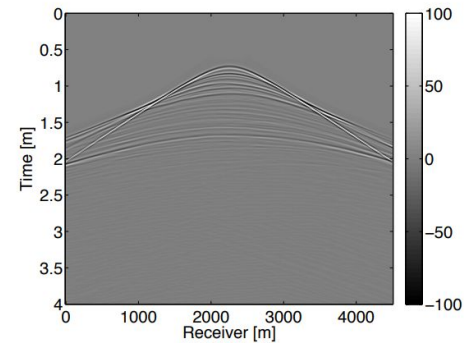


Data residual



# Linearization error

$$\nabla F(\mathbf{m}_0, \mathbf{q})^T \delta \mathbf{m} \approx$$



Data residual

# Seismic imaging

$$\nabla F(\mathbf{m}_0, \mathbf{q})^T \delta \mathbf{m} \simeq \delta \mathbf{d}$$



Linearized born forward modeling operator

$$\mathbf{J}(\mathbf{m}_0, \mathbf{q}) \equiv \nabla F(\mathbf{m}_0, \mathbf{q})^T$$

# Seismic imaging



# Least squares seismic imaging

$$\arg \min_{\delta \mathbf{m}} \quad \frac{1}{2} \sum_{i=1}^N \|\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i) \delta \mathbf{m} - \delta \mathbf{d}_i\|_2^2 \quad \text{subject to} \quad \delta \mathbf{m} \in \mathcal{C}$$

where,  $\delta \mathbf{d}_i = \mathbf{d}_{\text{obs},i} - F(\mathbf{m}_0, \mathbf{q}_i)$

- $N$  Number of source experiments
- $\mathbf{q}_i$  Source signature in  $i^{\text{th}}$  source experiments
- $\mathbf{d}_{\text{obs},i}$  Observed data in  $i^{\text{th}}$  source experiments
- $\delta \mathbf{d}_i$  Data residual in  $i^{\text{th}}$  source experiments
- $\mathcal{C}$  A constraint set encoding our prior knowledge

# Computational challenges

Applying  $\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i)$  and  $\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i)^T$  is expensive,

- Involves two PDE solves:


$$\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i) = \nabla F(\mathbf{m}_0, \mathbf{q})^T = -\mathbf{P}\mathbf{A}(\mathbf{m}_0)^{-1} [\nabla \mathbf{A}(\mathbf{m}_0) (\mathbf{A}(\mathbf{m}_0)^{-1} \mathbf{q})]$$

Many source experiments,

- $N$  is large
- Evaluating the objective function requires  $2N$  PDE solves.

# Need for regularization/prior

Solving inconsistent system of equations,

- linearization error, i.e.,  $\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i)\delta\mathbf{m} + \mathcal{O}(\|\delta\mathbf{m}\|_2^2) = \delta\mathbf{d}_i$
- noise in observed data, i.e.,  $\mathbf{d}_{\text{obs},i} = F(\mathbf{m}, \mathbf{q}_i) + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \sim p_{\text{noise}}(\boldsymbol{\epsilon})$

Good choices for regularization/prior is challenging,

- Hand crafted priors bias solution towards handcrafted priors

# Stochastic optimization edit

## Stochastic optimization (over source experiments)

- Approximate objective/gradient with a minibatch of source experiments

$$\sum_{i=1}^N \|\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i) \delta \mathbf{m} - \delta \mathbf{d}_i\|_2^2 \simeq \frac{N}{n} \sum_{i=1}^n \|\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i) \delta \mathbf{m} - \delta \mathbf{d}_i\|_2^2$$

- Typically we need 2-3 passes over the entire set of source experiments

# Regularization/prior

## Handcrafted quadratic penalty terms

- Tikhonov regularization, a.k.a (weighted) ridge regression

$$\arg \min_{\delta \mathbf{m}} \quad \frac{1}{2} \sum_{i=1}^N \|\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i) \delta \mathbf{m} - \delta \mathbf{d}_i\|_2^2 + \frac{\lambda^2}{2} \|\mathbf{R} \delta \mathbf{m}\|_2^2$$

- may adversely affect gradient & Hessian
- no guarantees that all model iterates are (physically) feasible
- biases solution towards handcrafted prior



# Regularization/prior

## Handcrafted constraints

- Total variation constraint

$$\arg \min_{\delta \mathbf{m}} \quad \frac{1}{2} \sum_{i=1}^N \|\mathbf{J}(\mathbf{m}_0, \mathbf{q}_i) \delta \mathbf{m} - \delta \mathbf{d}_i\|_2^2 \quad \text{subject to} \quad \delta \mathbf{m} \in \mathcal{C}$$

- (physical) constraints are satisfied during every iteration
- biases solution towards handcrafted prior

# Generative models as priors

If we have access to sufficient samples from the ground truth distribution

- train a generative model, e.g. GAN, VAE, to sample the distribution

$$\mathbf{x} \sim p_X^{\text{gen}}(\mathbf{x}) \Rightarrow \mathbf{x} = \mathbf{g}(\mathbf{z}), \mathbf{z} \sim p_Z(\mathbf{z})$$

$$p_X^{\text{gen}}(\mathbf{x}) = p_X^{\text{true}}(\mathbf{x})$$

- use the *pre-trained* generative model as prior

# Bayesian Inference

Statistical approach for formulating & solving inverse problems

- quantifies uncertainty in inference and motivates regularization

$$\begin{aligned} p_X^{\text{post}}(\mathbf{x}|\mathbf{y}) &= \frac{1}{\mathbb{Z}} p^{\text{l}}(\mathbf{y}|\mathbf{x}) p_X^{\text{prior}}(\mathbf{x}) \\ &= \frac{1}{\mathbb{Z}} p_{\eta}(\hat{\mathbf{y}} - \mathbf{f}(\mathbf{x})) p_X^{\text{prior}}(\mathbf{x}) \end{aligned}$$

$p_X^{\text{prior}}$ : prior distribution.

$p_{\eta}$ : distribution for the noise in measurement

# Inference with generative model

Use the *pre-trained* generative model as prior in Bayesian inference

$$p_X^{\text{post}}(\mathbf{x}|\mathbf{y}) = \frac{1}{\mathbb{Z}} p_\eta(\hat{\mathbf{y}} - \mathbf{f}(\mathbf{x})) p_X^{\text{gen}}(\mathbf{x})$$

Performing inference for an arbitrary function:

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim p_X^{\text{post}}} [l(\mathbf{x})] &= \frac{1}{\mathbb{Z}} \mathbb{E}_{\mathbf{x} \sim p_X^{\text{gen}}} [l(\mathbf{x}) p_\eta(\hat{\mathbf{y}} - \mathbf{f}(\mathbf{x}))], \\ &= \frac{1}{\mathbb{Z}} \mathbb{E}_{\mathbf{z} \sim p_Z} [l(\mathbf{g}(\mathbf{z})) p_\eta(\hat{\mathbf{y}} - \mathbf{f}(\mathbf{g}(\mathbf{z})))], \\ &= \mathbb{E}_{\mathbf{z} \sim p_Z^{\text{post}}} [l(\mathbf{g}(\mathbf{z}))], \end{aligned}$$

# Approximate inference

Drawing samples from the posterior can be done by performing MCMC on ([more on this later](#))

$$p_Z^{\text{post}}(\mathbf{z}|\mathbf{y}) \equiv \frac{1}{\mathbb{Z}} p_{\eta}(\hat{\mathbf{y}} - \mathbf{f}(\mathbf{g}(\mathbf{z}))) p_Z(\mathbf{z})$$

i.e.,  $p_Z^{\text{mcmc}}(\mathbf{z}|\mathbf{y}) \approx p_Z^{\text{post}}(\mathbf{z}|\mathbf{y})$ . Next, for any function:

$$\overline{l(\mathbf{x})} \equiv \mathbb{E}_{\mathbf{x} \sim p_X^{\text{post}}} [l(\mathbf{x})] \approx \frac{1}{N_{\text{samp}}} \sum_{n=1}^{N_{\text{samp}}} l(\mathbf{g}(\mathbf{z})), \quad \mathbf{z} \sim p_Z^{\text{mcmc}}(\mathbf{z}|\mathbf{y}).$$

# Why generative model as prior?

The prior is no longer hand-crafted

- solution of the inverse problem is not biased towards the prior
  - **Only if**  $p_X^{\text{gen}}(\mathbf{x}) = p_X^{\text{true}}(\mathbf{x})$
- performing MCMC on the latent variable is easier
  - **much smaller dimension**

**Big assumption:** we have access to samples from the ground truth distribution

- **limits the application in seismic inversion/imaging**

## Limitations in seismic inverse problems

Due to Earth's heterogeneity, we do not have access to samples from true prior

$$\{\boldsymbol{x}_i\}_{i=1}^N \sim p_{prior}(\boldsymbol{x})$$

i.e., limited chance of pre-training a generative model

We have many source experiments  $\mathbf{d}_{\text{obs},i}$  explaining one and only one unknown,

i.e.,  $(\mathbf{y}_1, \mathbf{x}), (\mathbf{y}_2, \mathbf{x}), \dots, (\mathbf{y}_N, \mathbf{x})$

# Deep Image Prior

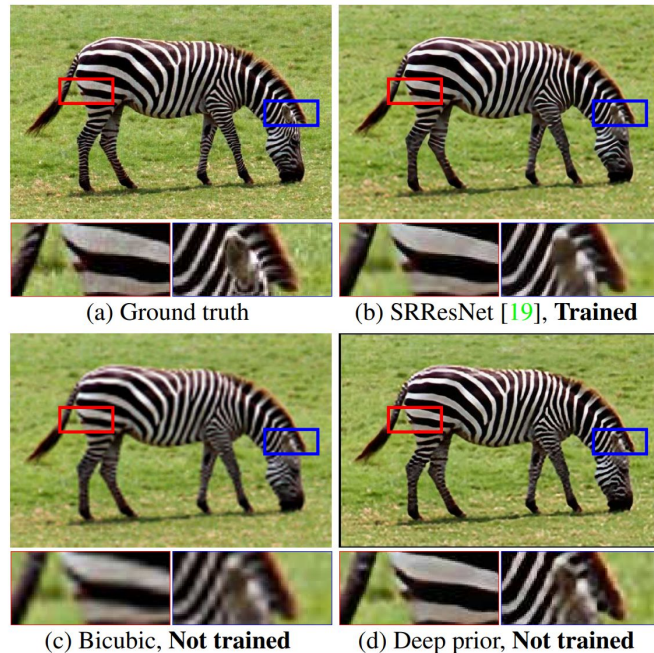
An *untrained* neural net as a prior

$$w^* = \arg \min_w \|y - AG(z; w)\|^2$$

where  $y = Ax^* + \eta$  is some noisy observation

$z \in \mathbb{R}^k$  fixed random vector

$G(z; w): \mathbb{R}^k \rightarrow \mathbb{R}^n$  *Untrained network*





## Deep Image Prior

$$\underset{w}{\text{minimize}} \frac{1}{2} \|y - Ag(z, w)\|_2^2$$

is equivalent to:

$$\underset{x, w}{\text{minimize}} \frac{1}{2} \|y - Ax\|_2^2 \quad \text{subject to} \quad x = g(z, w)$$

- network topology acts as a generic prior for “natural images”
- does not bias solution (untrained network adds little info)
- Non-convex and no guarantee that is (physically) feasible

## Constrained weak deep prior

Replace w/ *weak* constraint:

$$\underset{x, w}{\text{minimize}} \frac{1}{2} \|y - Ax\|_2^2 + \frac{\lambda^2}{2} \|x - g(z, w)\|_2^2$$

and add *strong* handcrafted constraints (TV-norm):

$$\text{→} \underset{x \in \mathcal{C}, w}{\text{minimize}} \frac{1}{2} \|y - Ax\|_2^2 + \frac{\lambda^2}{2} \|x - g(z, w)\|_2^2$$

- ▶ imposes solution to remain in range generator w/ error proportional to  $\lambda$
- ▶ strong constraints balanced w/ deep prior
- ▶ **strong handcrafted constraints guarantee (physically) feasibility**
- ▶ **does not sample posterior when training w/ *single*  $z$**

# Alternating back-propagation

A method for training generative models

- based on EM algorithm
- Does not need an extra network, e.g., recognition network in VAE, discriminator in GAN
- able to learn from incomplete or indirect data
  - This may prove difficult or less convenient for VAE and GAN
- no need for  $x_i$ 's

## Setup

- training set of data vectors  $\{Y_i, i = 1, \dots, n\}$
- each  $Y_i$  has a corresponding  $Z_i$  (not observed)
- Probabilistic model:  $[Y|Z, W] \sim p(Y|Z, W)$

$$Y = f(Z; W) + \epsilon,$$

$$Z \sim N(0, I_d), \epsilon \sim N(0, \sigma^2 I_D), d < D$$

**Goal:** Find parameters  $W$  that maximizes

$$p(Y; W) = \int p(Z)p(Y|Z, W)dZ \approx \sum_{i=1}^n \log p(Y_i; W) = \\ \sum_{i=1}^n \log \int p(Y_i, Z_i; W)dZ_i$$

## EM algorithm to train generative models

$$\begin{aligned}\frac{\partial}{\partial W} \log p(Y; W) &= \frac{1}{p(Y; W)} \frac{\partial}{\partial W} p(Y; W) \\&= \frac{1}{p(Y; W)} \int \frac{\partial}{\partial W} p(Z, Y; W) dZ = \frac{p(Z|Y; W)}{p(Z, Y; W)} \int \frac{\partial}{\partial W} p(Z, Y; W) dZ \\&= \int p(Z|Y; W) \frac{1}{p(Z, Y; W)} \frac{\partial}{\partial W} p(Z, Y; W) dZ = \int p(Z|Y; W) \frac{\partial}{\partial W} \log p(Z, Y; W) dZ \\&= \mathbb{E}_{p(Z|Y; W)} \left[ \frac{\partial}{\partial W} \log p(Z, Y; W) \right]\end{aligned}$$

E-step: Approximate the expectation by drawing samples from  $p(Z|Y, W)$  with

**Stochastic Gradient Langevin Dynamics**

## Stochastic Gradient Langevin Dynamics

- Task: Given a target distribution  $d\mu = e^{-V(\mathbf{x})}d\mathbf{x}$ , generate samples from  $\mu$ .
  - ▷ Most samples should be gathered around the minimum of  $V$
  - ▷ We do not want convergence to the minimum

- Idea: Use Gradient Descent + noise

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \beta_k \nabla V(\mathbf{X}^k) + \text{noise}$$

- Question: What amount of noise should we add so that  $\mathbf{X}^\infty \sim \pi$ ?

# Large scale sampling

- A scalable framework: First-order sampling (assuming access to  $\nabla V$ ).

## Step 1. Langevin Dynamics

$$d\mathbf{X}_t = -\nabla V(\mathbf{X}_t)dt + \sqrt{2}d\mathbf{B}_t \quad \Rightarrow \quad \mathbf{X}_\infty \sim e^{-V}.$$

## Step 2. Discretize

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \beta_k \nabla V(\mathbf{x}^k) + \sqrt{2\beta_k} \boldsymbol{\xi}^k$$

- ▷  $\beta_k$  step-size,  $\boldsymbol{\xi}^k$  standard normal
- ▷ strong analogy to gradient descent method

## Recall

### Complete data log-likelihood

$$\begin{aligned} Y &= f(Z; W) + \epsilon, \\ Z &\sim \mathcal{N}(0, I_d), \epsilon \sim \mathcal{N}(0, \sigma^2 I_D), d < D \end{aligned} \quad \Rightarrow \quad \begin{aligned} \log p(Y, Z; W) &= \log [p(Z)p(Y|Z, W)] \\ &= -\frac{1}{2\sigma^2} \|Y - f(Z; W)\|^2 - \frac{1}{2} \|Z\|^2 + \text{const.} \end{aligned}$$

The posterior of  $Z$  can be written as (Bayes' rule)

$$p(Z|Y, W) = p(Y, Z; W)/p(Y; W) \propto p(Z)p(Y|Z, W)$$



## E-step: Langevin dynamics

$$p(Z|Y, W) = p(Y, Z; W)/p(Y; W) \propto -\frac{1}{2\sigma^2}\|Y - f(Z; W)\|^2 - \frac{1}{2}\|Z\|^2$$

Samples from  $p(Z|Y, W)$  are obtained from iterates:

$$Z_{\tau+1} = Z_{\tau} + sU_{\tau} + \frac{s^2}{2} \left[ \frac{1}{\sigma^2} (Y - f(Z_{\tau}; W)) \frac{\partial}{\partial Z} f(Z_{\tau}; W) - Z_{\tau} \right],$$

$\tau$  time step for the Langevin sampling

$s$  step size

$U_{\tau} \sim N(0, I_d)$

## M-step: backpropagation

Now that we have samples from  $p(Z|Y, W)$  we perform optimization of  $W$

$$\begin{aligned}\mathbb{E}_{p(Z|Y;W)} \left[ \frac{\partial}{\partial W} \log p(Z, Y; W) \right] &\approx \sum_{i=1}^n \frac{\partial}{\partial W} \log p(Y_i, Z_i; W) \\ &= - \sum_{i=1}^n \frac{\partial}{\partial W} \frac{1}{2\sigma^2} \|Y_i - f(Z_i; W)\|^2 \\ &= \sum_{i=1}^n \frac{1}{\sigma^2} (Y_i - f(Z_i; W)) \frac{\partial}{\partial W} f(Z_i; W)\end{aligned}$$

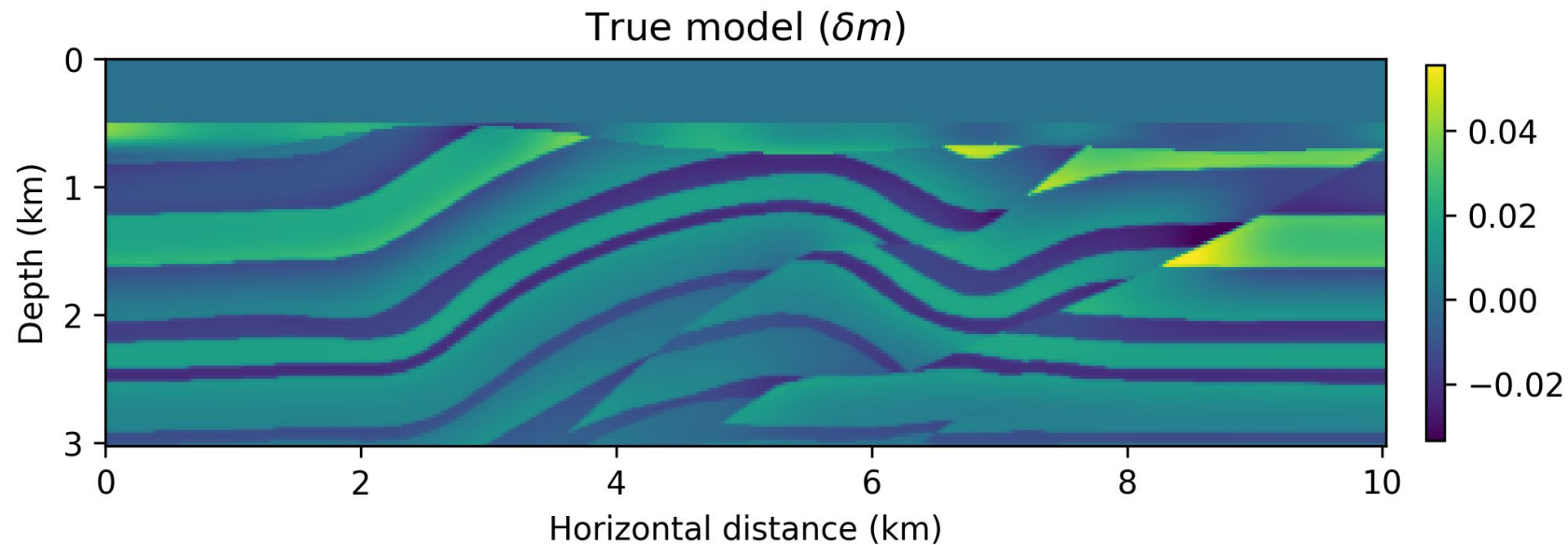
## After training/inversion: inference

After training keep  $w^* = w$  fixed

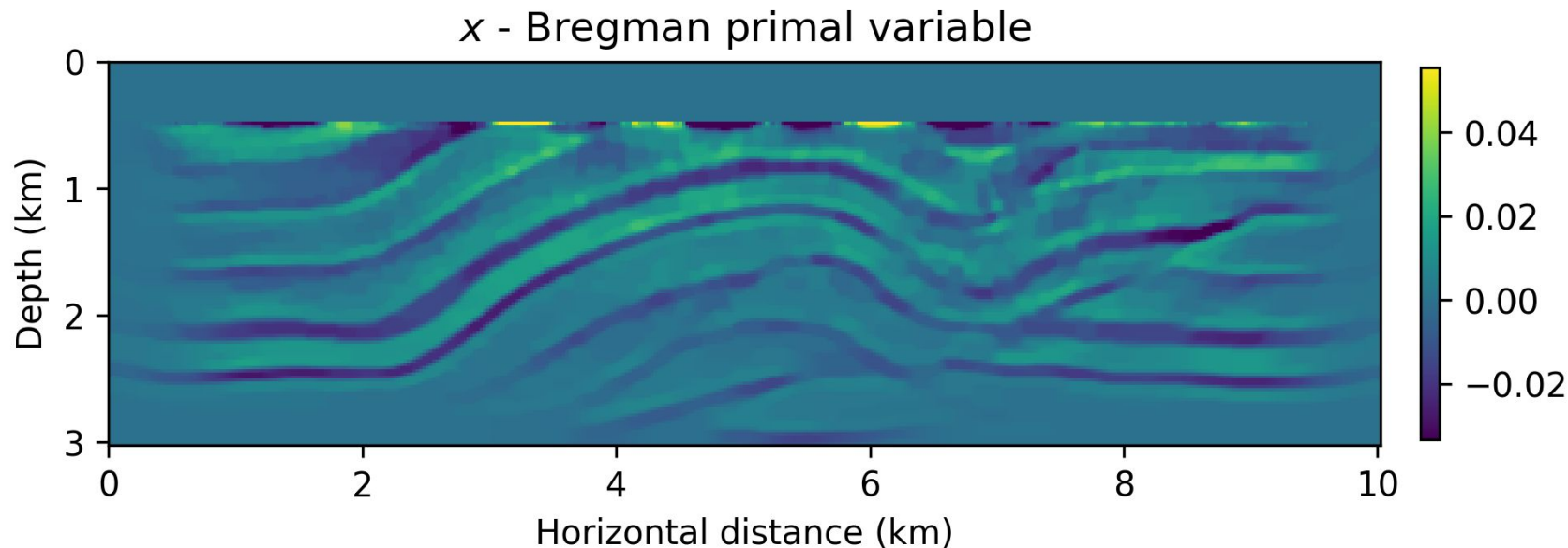
- ▶ draw samples from posterior via  $x \sim g(z, w^*)$  with  $z \sim N(0,1)$
- ▶ derive statistical properties posterior simply via

$$\mathbb{E}_z \left\{ f(g(z, w^*)) \right\} \approx \frac{1}{K} \sum_{k=1}^K f(g(z_k, w^*)), \quad z_k \sim N(0,1)$$

# True perturbation

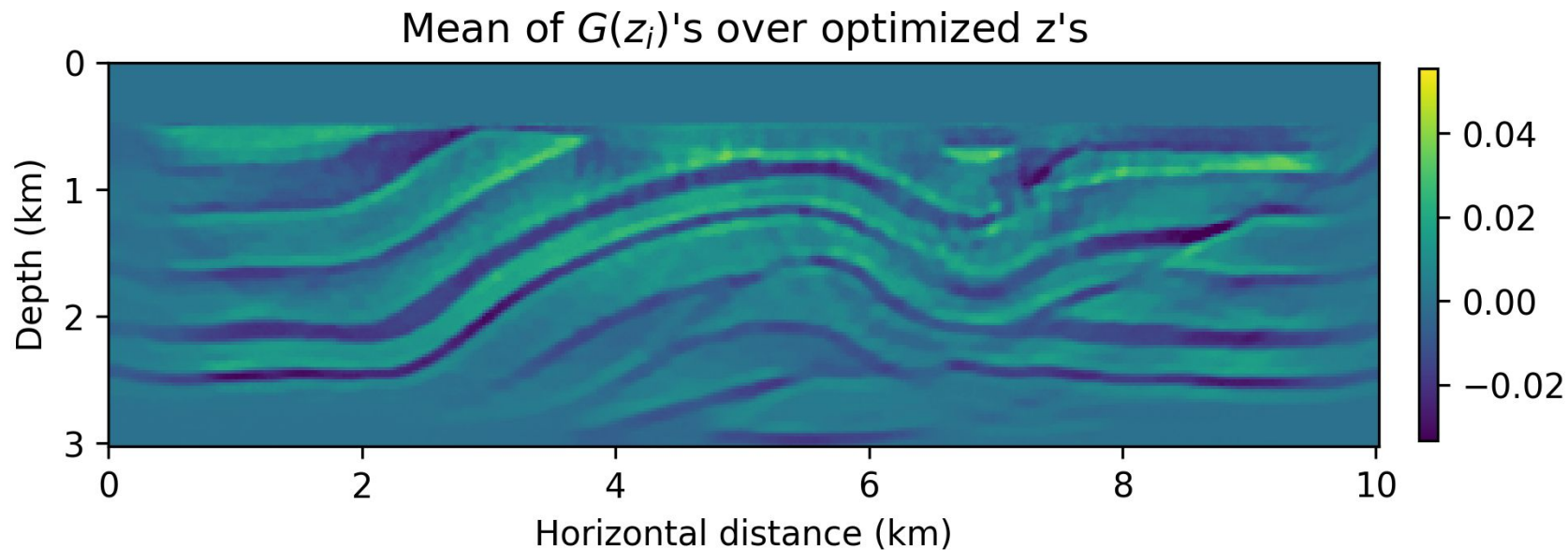


# Slack variable in weak formulation

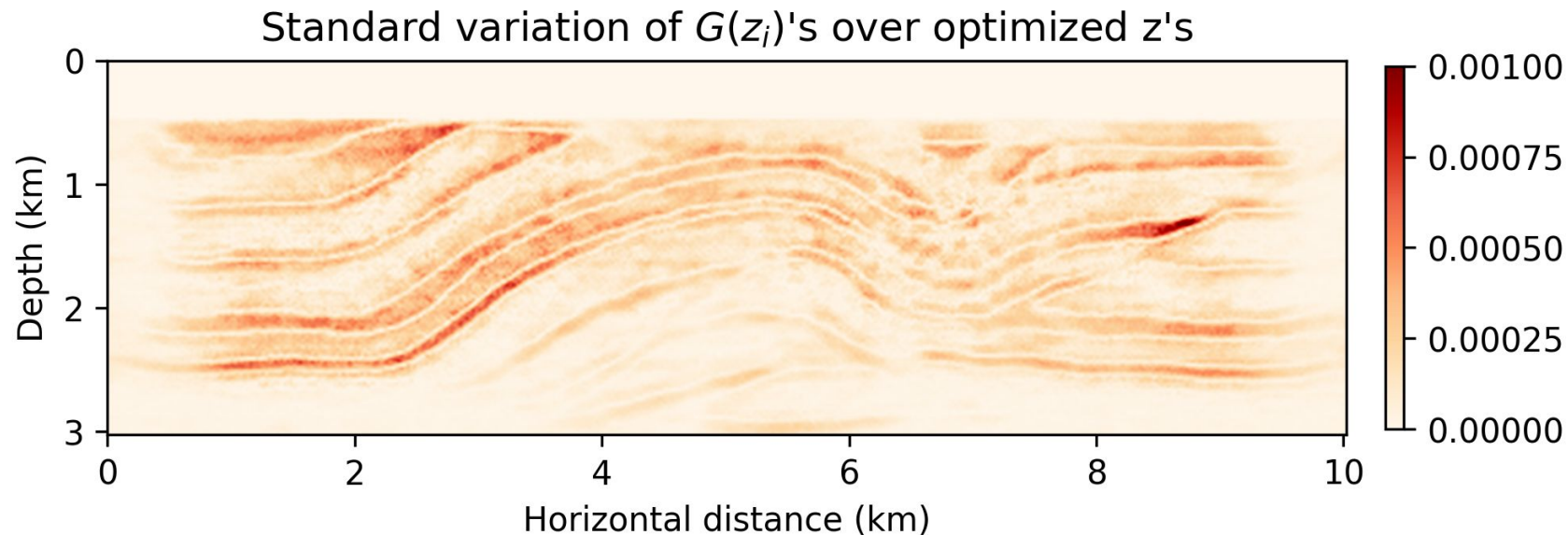


→ minimize  $\frac{1}{2} \|y - Ax\|_2^2 + \frac{\lambda^2}{2} \|x - g(z, w)\|_2^2$

# Mean

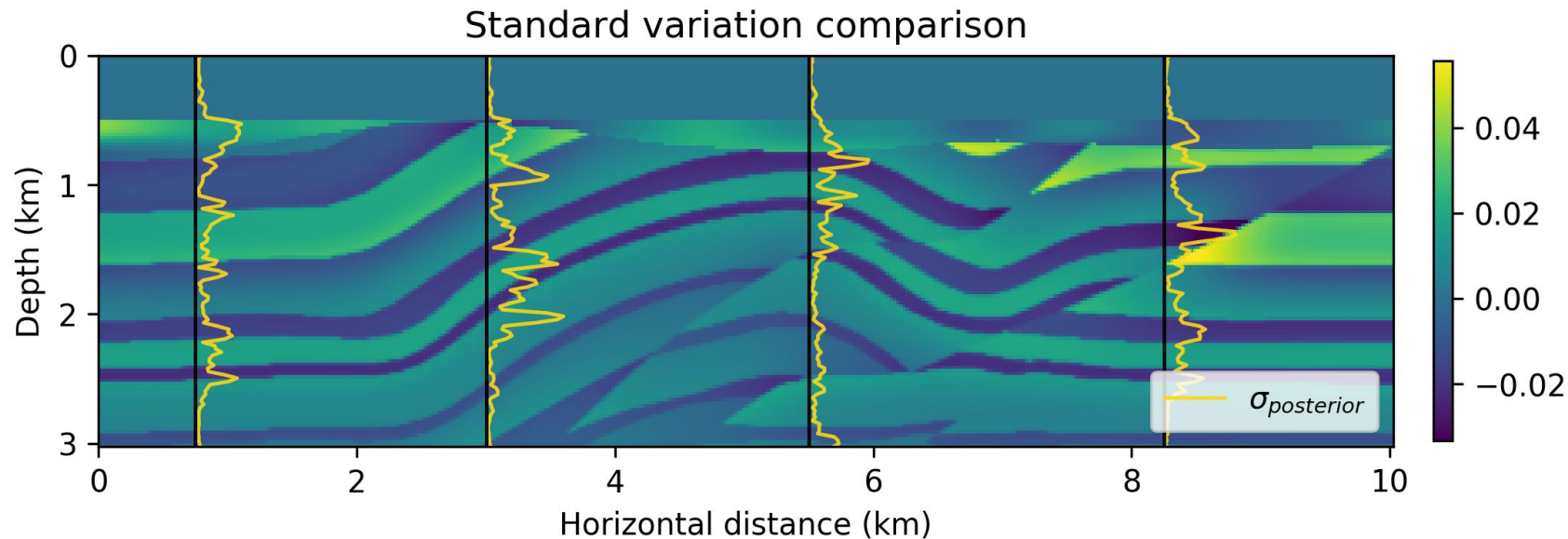


## Pointwise variance





# Pointwise variance: traces





# Observations

Deep priors:

- ▶ regularize via their network topology
- ▶ made feasible in combination w/ strong handcrafted constraints
- ▶ reap information on the posterior during inversion

Computationally feasible fully-data driven framework that provides UQ

More theory & (nonlinear) examples needed...