

Dictionary learning in geoscience

Michael Bianco

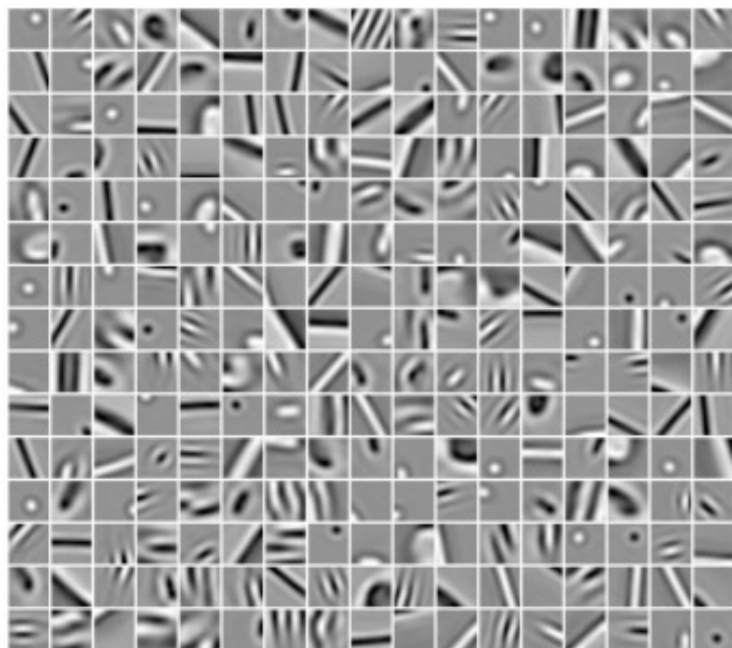
UCSD Noise Lab, Scripps Institution of Oceanography

noiselab.ucsd.edu

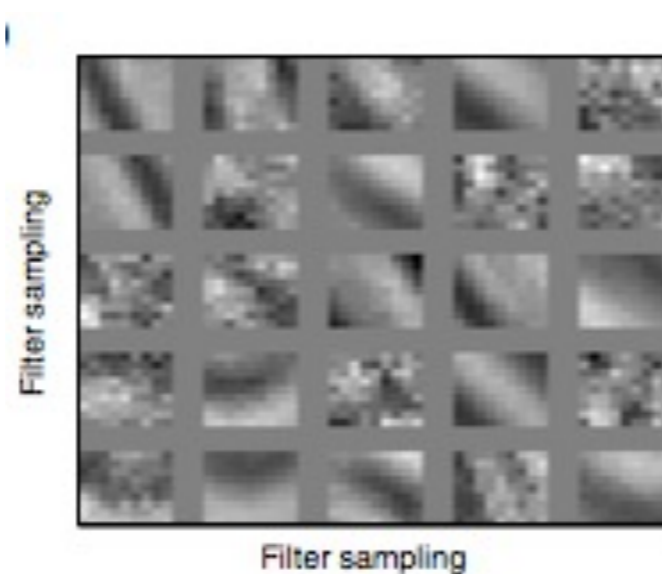
5/9/18

Dictionary learning

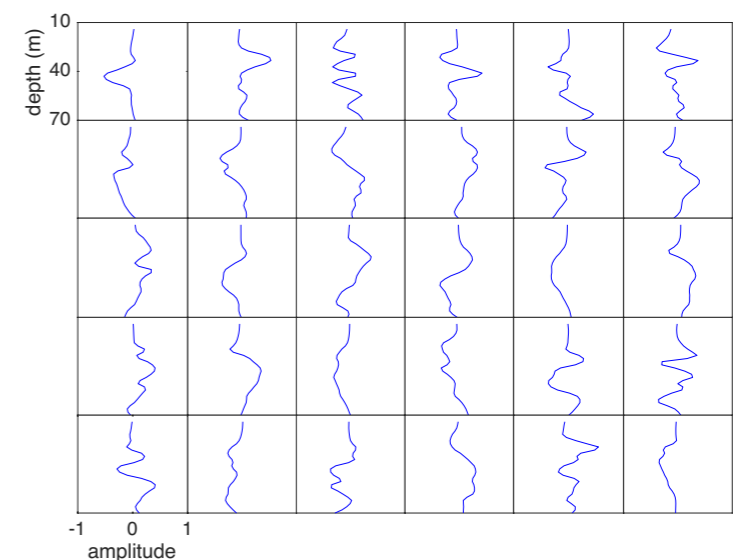
- Means of estimating sparse causes for given classes of signals, e.g. natural images, audio
- Originated in neuroscience to estimate structure of V1 visual cortex cells from natural images
- Useful for regularization of general image denoising inverse problem, but only recent applications in the geosciences
 - Seismic survey image denoising
 - Dictionary learning of ocean sound speed profiles (SSPs)



Olshausen 2009

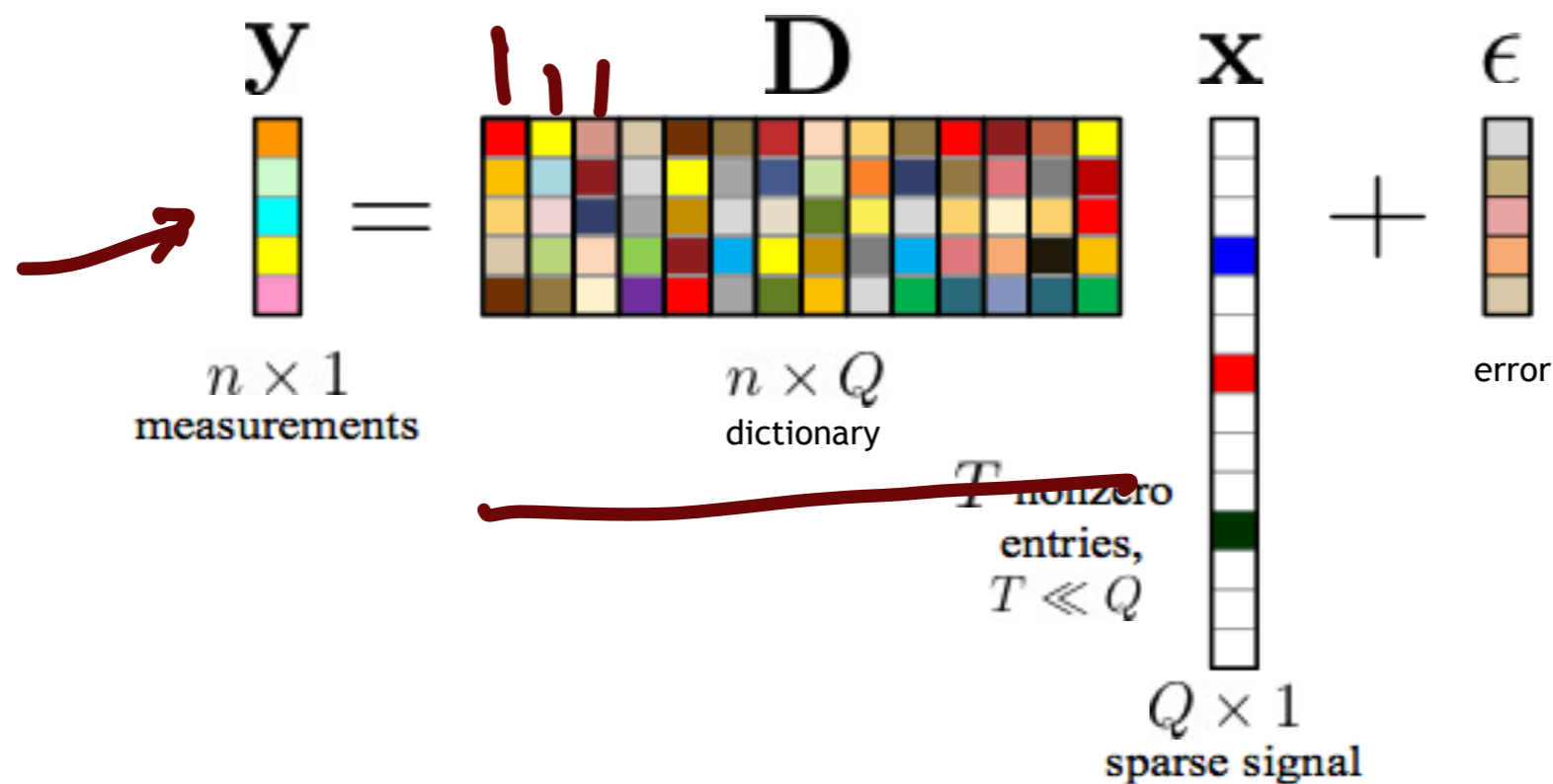


Beckouche 2014



Bianco and Gerstoft 2017

Background: sparse modeling of arbitrary signal \mathbf{y}

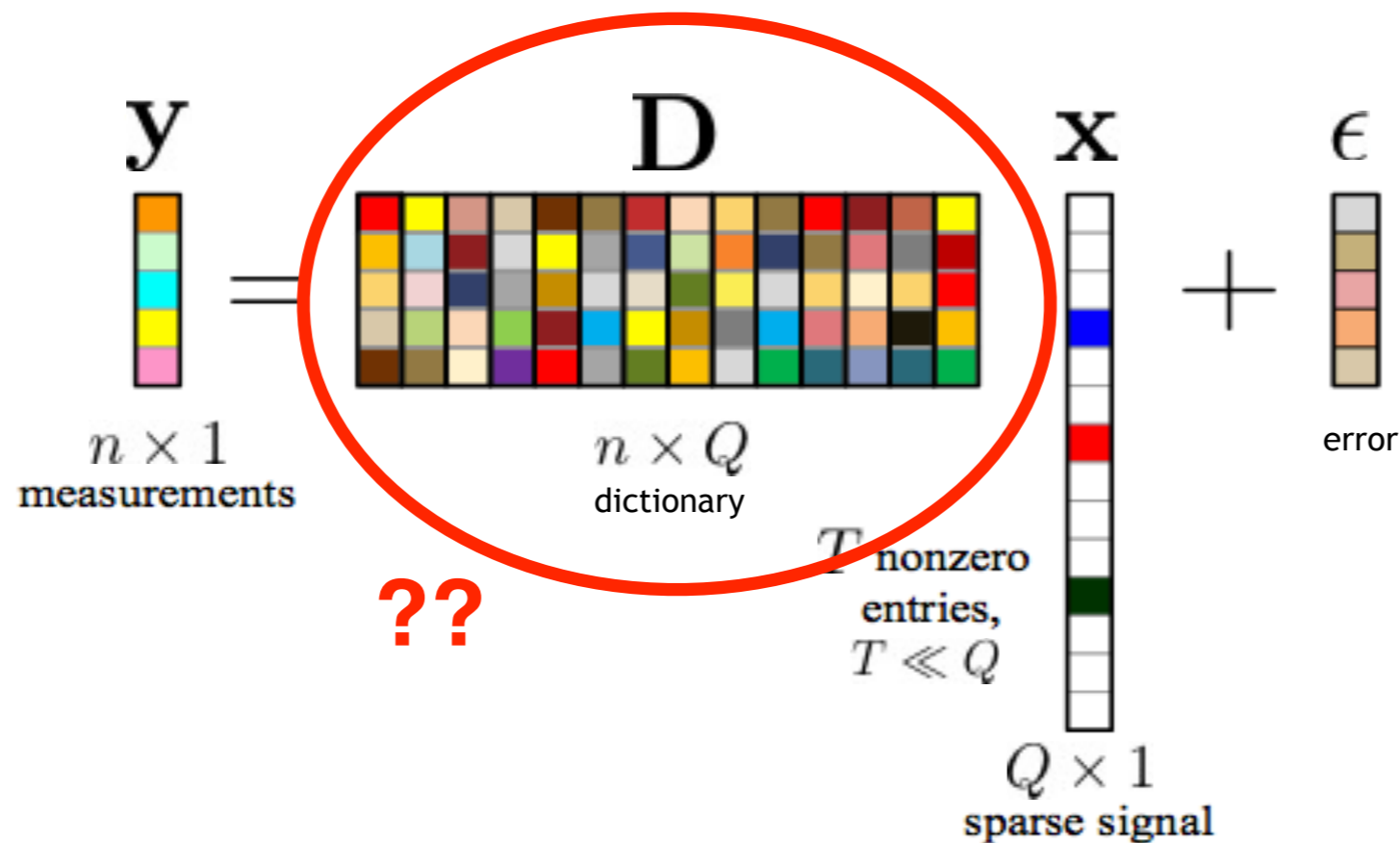


- Measurement vector \mathbf{y} is expressed as sparse linear combination of columns or "atoms" from dictionary \mathbf{D}
- \mathbf{y} could be (for example) segments of speech or vectorized 2D image patches
- Dictionary atoms represent elemental patterns that generate \mathbf{y} , e.g. wavelets or learned from the data using dictionary learning
- \mathbf{x} is estimated using sparsity inducing constraint, example " ℓ_0 -norm" regularization:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 \quad \text{subject to} \quad \|\mathbf{x}\|_0 \leq T$$

ℓ_0 -norm "counts" # non-zero coefficients

Background: sparse modeling of arbitrary signal \mathbf{y}



- Measurement vector \mathbf{y} is expressed as sparse linear combination of columns or "atoms" from dictionary \mathbf{D}
- \mathbf{y} could be (for example) segments of speech or vectorized 2D image patches
- Dictionary atoms represent elemental patterns that generate \mathbf{y} , e.g. wavelets or learned from the data using dictionary learning
- \mathbf{x} is estimated using sparsity inducing constraint, example " ℓ_0 -norm" regularization:

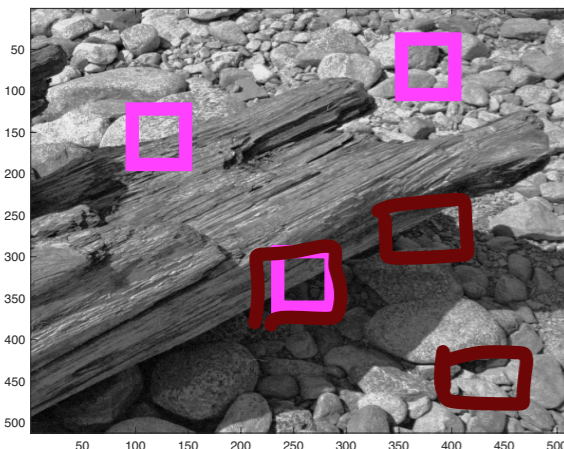
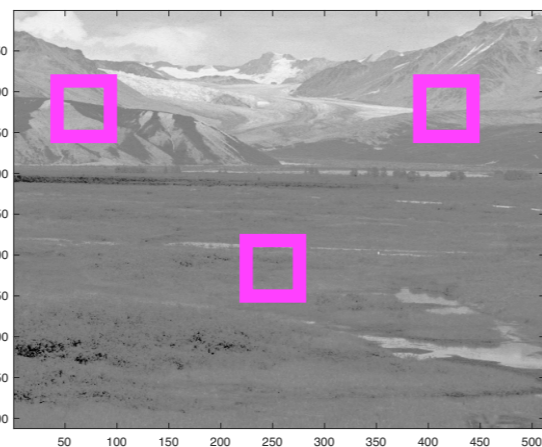
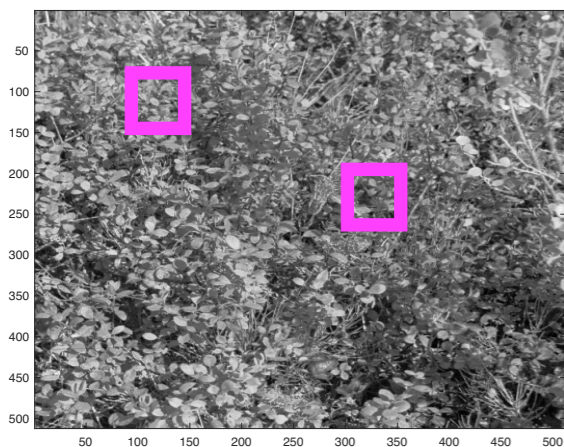
$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 \quad \text{subject to} \quad \|\mathbf{x}\|_0 \leq T$$

ℓ_0 - norm "counts" # non-zero coefficients

Background: sparsity and dictionary learning

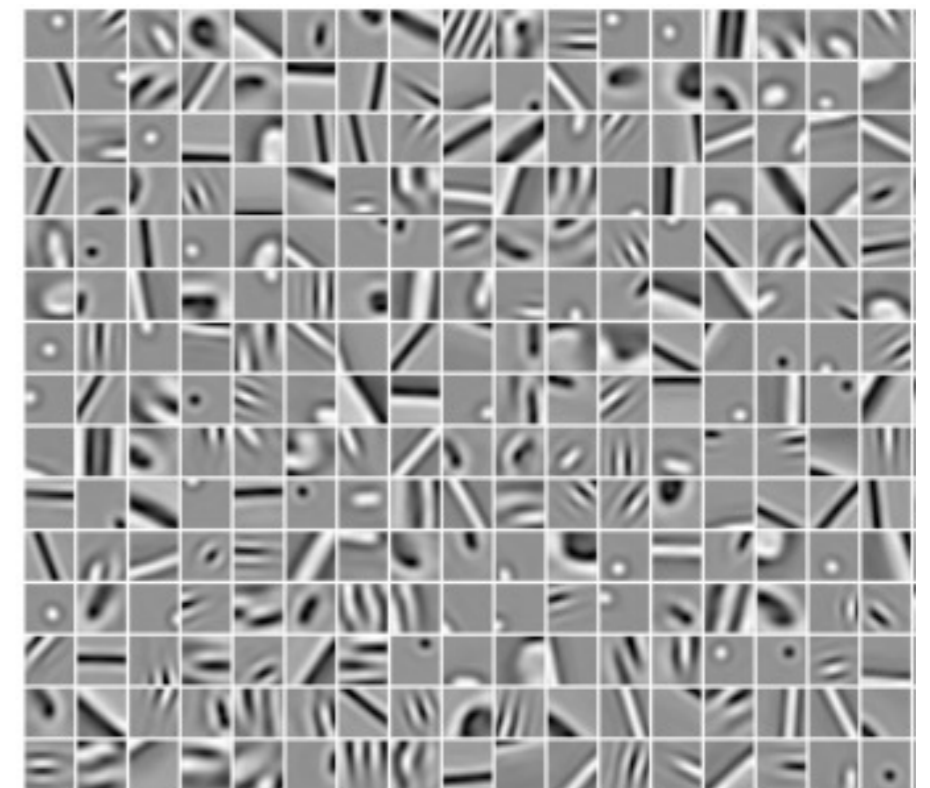
- Dictionary learning obtains "optimal" sparse modeling dictionaries directly from data
- Dictionary learning was developed in neuroscience (a.k.a. sparse coding) to help understand mammalian visual cortex structure
- Assumes (1) Redundancy in data: image patches are repetitions of a smaller set of elemental shapes; and (2) Sparsity: each patch is represented with few atoms from dictionary

"Natural" images, patches shown in magenta



- Each patch is signal \mathbf{y}_i
- Set of all patches $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_I]$

Learn dictionary \mathbf{D} describing $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_I]$



Olshausen 2009

Sparse model for patch \mathbf{y}_i composed of few atoms from \mathbf{D}

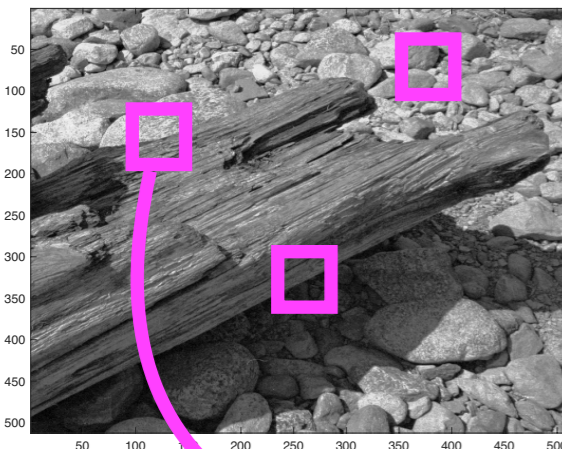
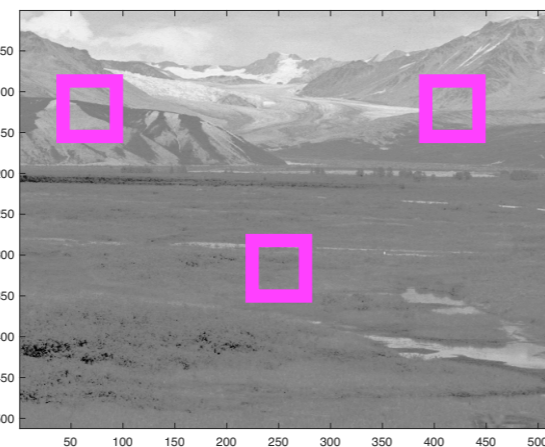
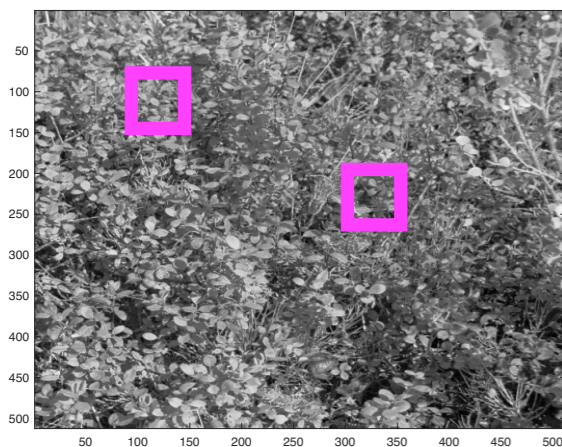
$$\hat{\mathbf{x}}_i = \arg \min_{\mathbf{x}_i} \|\mathbf{y}_i - \mathbf{D}\mathbf{x}_i\|_2 \quad \text{subject to} \quad \|\mathbf{x}_i\|_0 \leq T$$

"counts nonzero"

Background: sparsity and dictionary learning

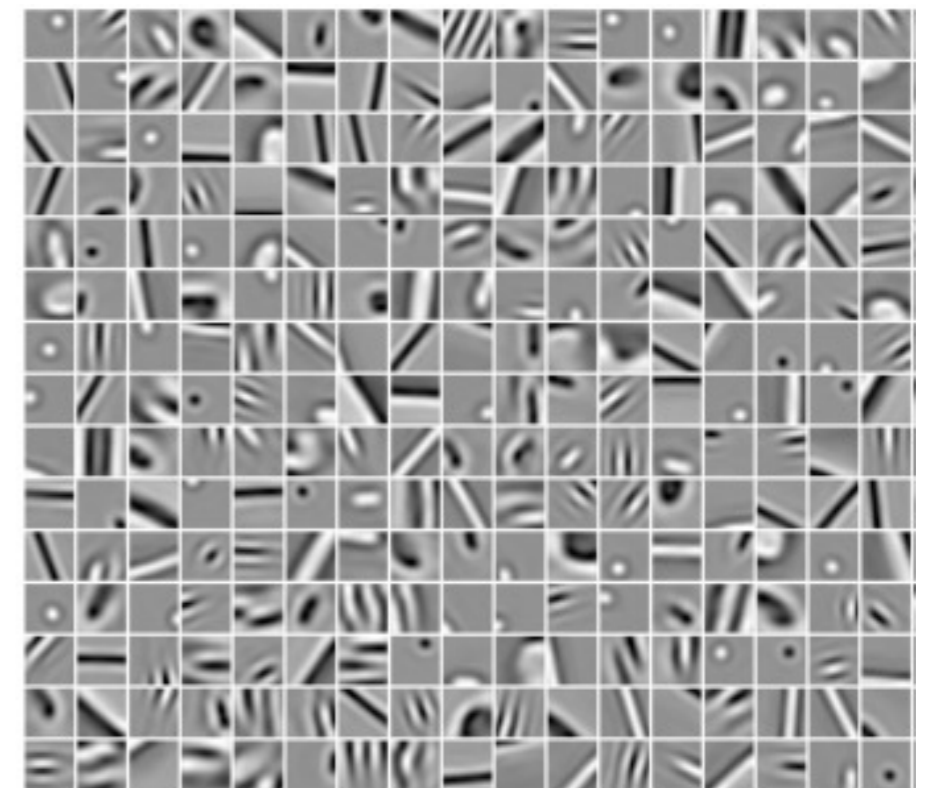
- Dictionary learning obtains "optimal" sparse modeling dictionaries directly from data
- Dictionary learning was developed in neuroscience (a.k.a. sparse coding) to help understand mammalian visual cortex structure
- Assumes (1) Redundancy in data: image patches are repetitions of a smaller set of elemental shapes; and (2) Sparsity: each patch is represented with few atoms from dictionary

"Natural" images, patches shown in magenta



- Each patch is signal \mathbf{y}_i
- Set of all patches $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_I]$

Learn dictionary \mathbf{D} describing $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_I]$



Olshausen 2009

Sparse model for patch \mathbf{y}_i composed of few atoms from \mathbf{D}

$$\mathbf{y} = \text{[patch]} = \text{[atom 1]} x_1 + \text{[atom 2]} x_2 + \dots$$

The second atom in the equation is circled in red.

Olshausen and Field 1997: image model with sparse prior

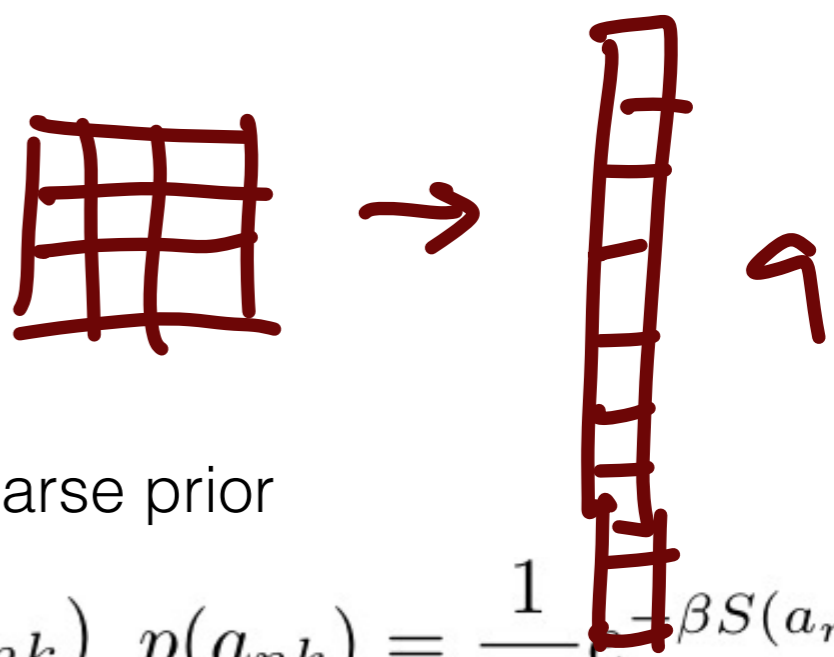
Assume that each image patch described by linear system

$$\mathbf{y}_k = \sum_n a_{nk} \phi_n = \Phi \mathbf{a}_k \quad \mathbf{y}_k = \Phi \mathbf{a}_k + \mathbf{n}$$

$y \in \mathbb{R}^k, k=64$

Goal: estimate bases Φ from observations \mathbf{y}_k
 Probability of image patch arising from bases Φ is

$$p(\mathbf{y}_k | \Phi) = \int p(\mathbf{y}_k | \mathbf{a}_k, \Phi) p(\mathbf{a}_k) d\mathbf{a}_k, \text{ with}$$

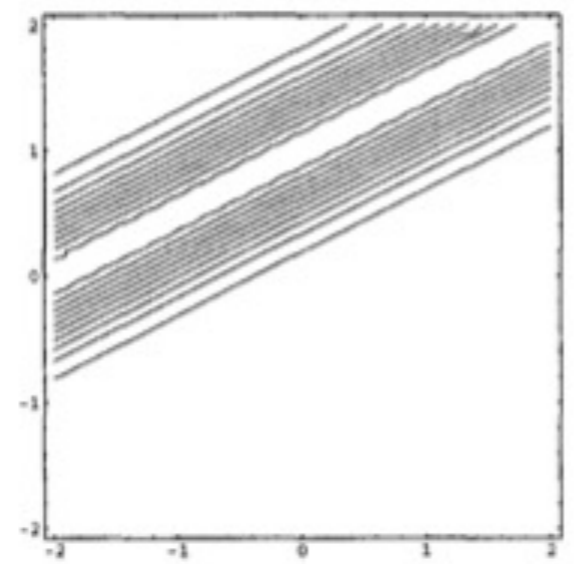
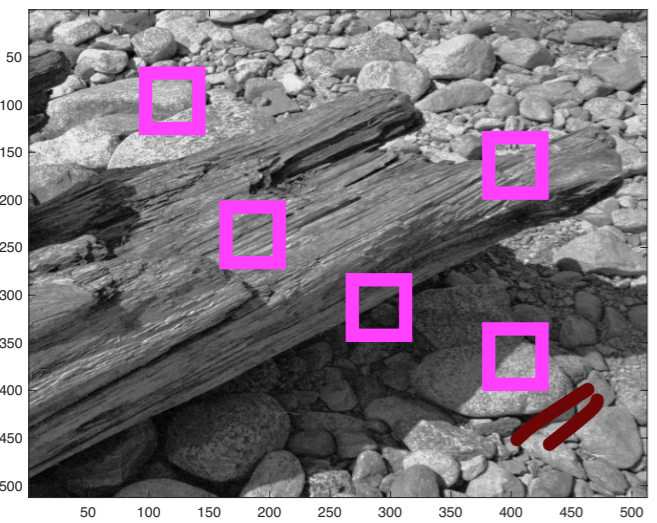


Likelihood

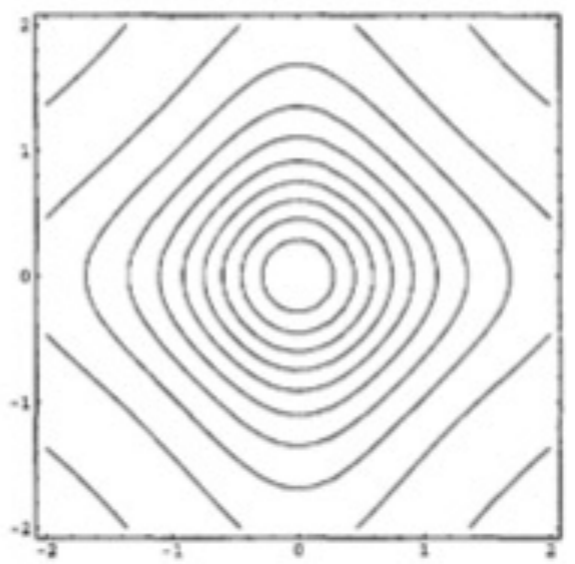
Independent, sparse prior

$$p(\mathbf{y}_k | \mathbf{a}_k, \Phi) = \frac{1}{Z_\sigma} e^{-\frac{\|\mathbf{y}_k - \Phi \mathbf{a}_k\|_2^2}{2\sigma^2}} \quad p(\mathbf{a}_k) = \prod p(a_{nk}) \quad p(a_{nk}) = \frac{1}{Z_\beta} e^{-\beta S(a_{nk})}$$

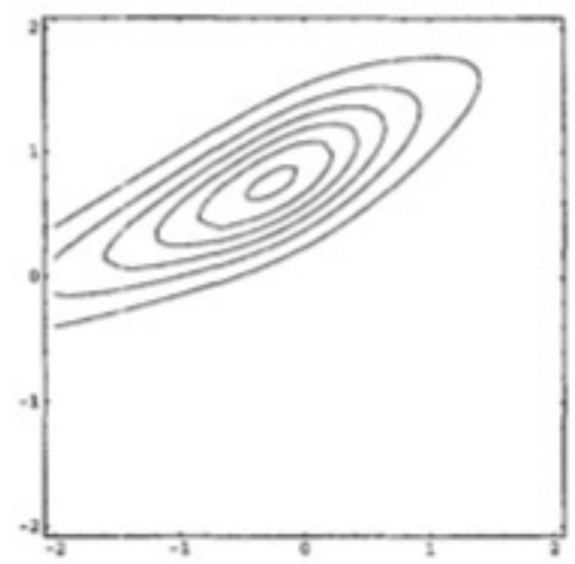
Image patches \mathbf{y}_k



Likelihood



Prior



Posterior

Olshausen and Field 1997- sparse prior induces sparse coefficients

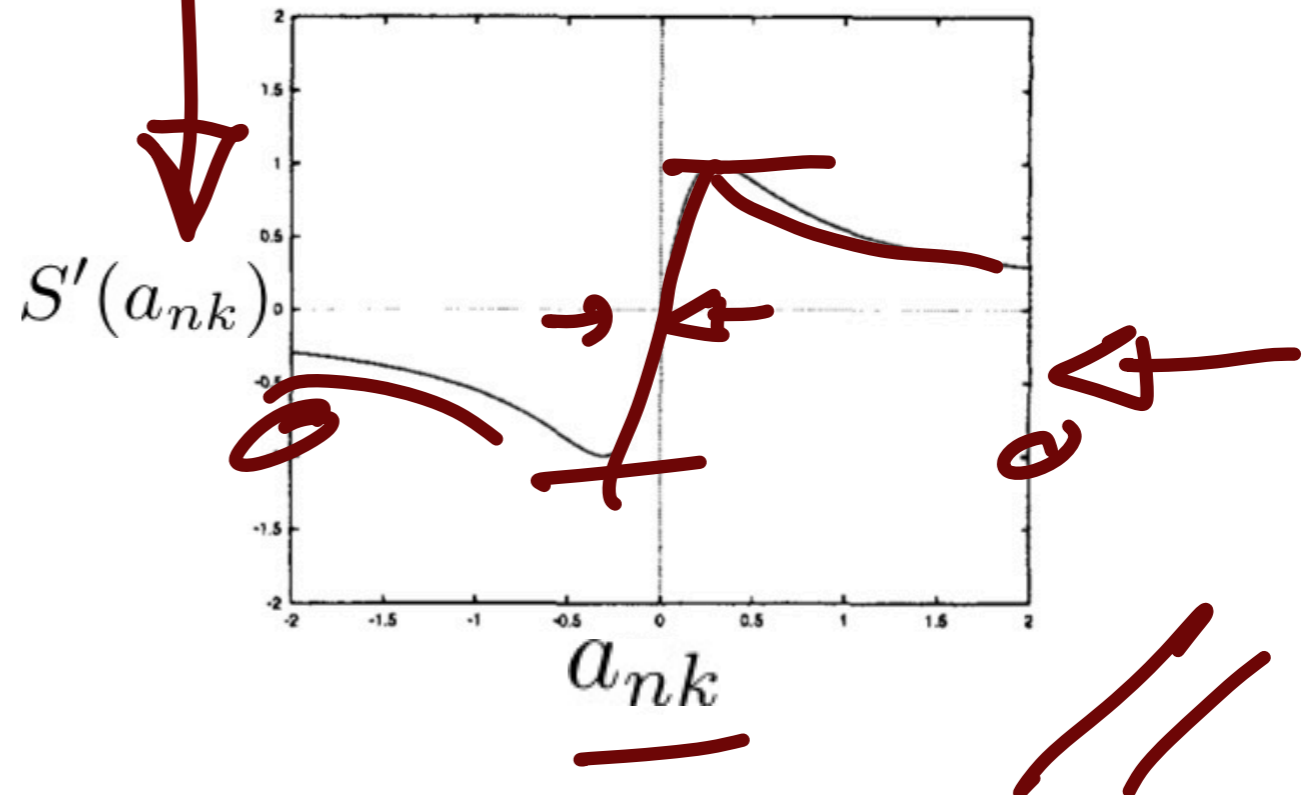
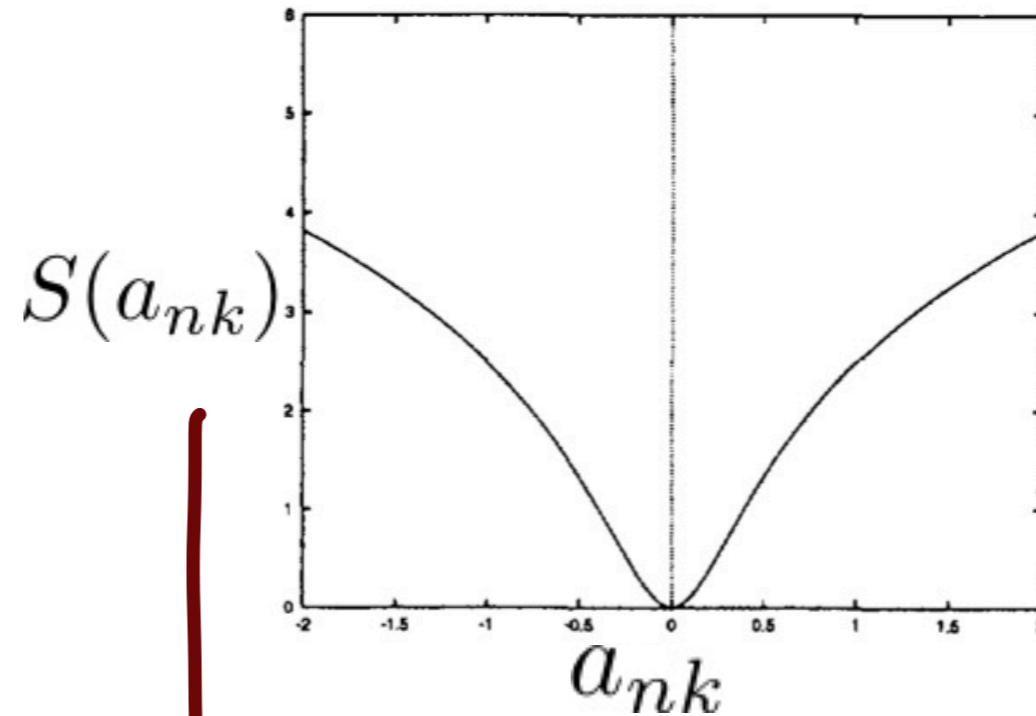
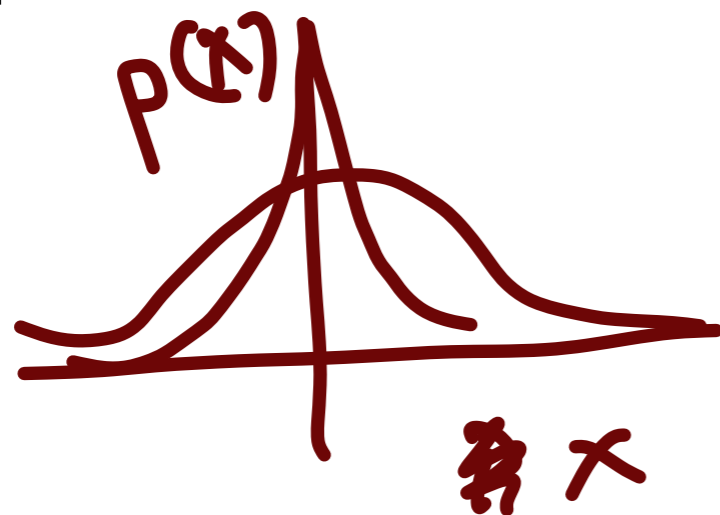
Sparsity inducing prior

$$p(a_{nk}) = \frac{1}{Z_\beta} e^{-\beta S(a_{nk})}$$

$$S(a_{nk}) = \ln(1 + a_{nk}^2)$$

"Cauchy distribution"

Derivative of prior induces sparsity in solution, as we'll see...



Olshausen and Field 1997 - derivation of Error function

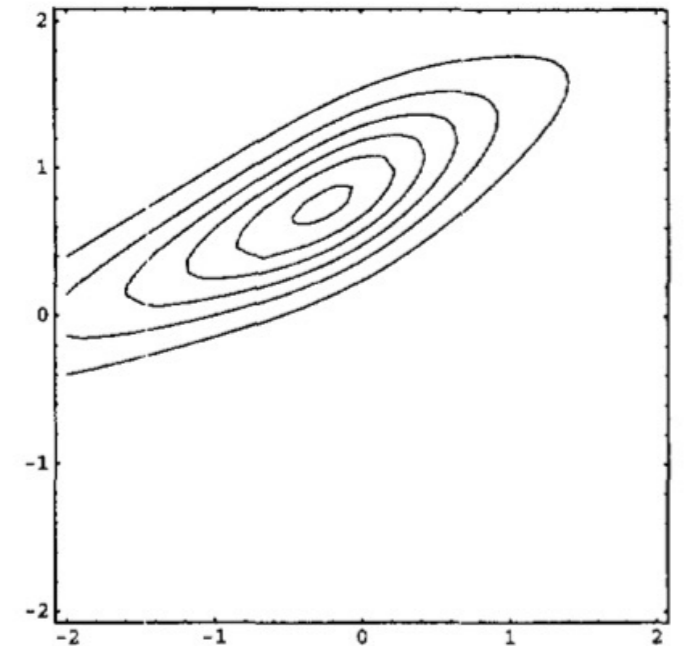
Learn basis functions Φ by minimizing Kullback-Leibler (KL) divergence between true images and those reproduced by model

$$KL = \int \underline{p^*(\mathbf{y}_k)} \ln \frac{\underline{p^*(\mathbf{y}_k)}}{\underline{p(\mathbf{y}_k|\Phi)}} d\mathbf{y}_k$$

Since $p^*(\mathbf{y}_k)$ is fixed, KL is minimized by maximizing log-likelihood (or minimizing negative log-likelihood) of image patches generated from model, hence

$$\{\hat{\Phi}, \hat{\mathbf{a}}_k\} = \arg \min_{\Phi} \left[\min_{\mathbf{a}_k} E(\mathbf{y}_k, \mathbf{a}_k | \Phi) \right]$$

$$p(\mathbf{y}_k, \mathbf{a}_k | \Phi) = p(\mathbf{y}_k | \mathbf{a}_k, \Phi) p(\mathbf{a}_k)$$



$$E(\mathbf{y}_k, \mathbf{a}_k | \Phi) = -\ln \underline{p(\mathbf{y}_k | \mathbf{a}_k, \Phi) p(\mathbf{a}_k)}$$

Given: $p(\mathbf{y}_k | \mathbf{a}_k, \Phi) = \frac{1}{Z_\sigma} e^{-\frac{\|\mathbf{y}_k - \Phi \mathbf{a}_k\|_2^2}{2\sigma^2}}$ $p(\mathbf{a}_k) = \prod_n p(a_{nk})$ $p(a_{nk}) = \frac{1}{Z_\beta} e^{-\beta S(a_{nk})}$

$$E(\mathbf{y}_k, \mathbf{a}_k | \Phi) = \|\mathbf{y}_k - \Phi \mathbf{a}_k\|_2^2 + \lambda \sum_n S(a_{nk})$$

Olshausen and Field 1997 - derivation of Error function cont'd

Learn basis functions Φ by minimizing Kullback-Leibler (KL) divergence between true images and those reproduced by model

$$KL = \int p^*(\mathbf{y}_k) \ln \frac{p^*(\mathbf{y}_k)}{p(\mathbf{y}_k|\Phi)} d\mathbf{y}_k \quad \text{Min KL}$$

$$= \int p^*(y) [\ln p^*(y) - \ln p(y|\Phi)] dy$$

$$= \int p^*(y) \ln p^*(y) dy - \int p^*(y) \ln p(y|\Phi) dy$$

$$\text{Min KL} = \text{Min} \quad \quad = \text{max} \int p^*(y) \ln p(y|\Phi) dy$$

$\text{max} \langle \ln p(y|\Phi) \rangle$

Olshausen and Field 1997 - derivation of Error function cont'd

Learn basis functions Φ by minimizing Kullback-Leibler (KL) divergence between true images and those reproduced by model

$$\{\hat{\Phi}, \hat{\mathbf{a}}_k\} = \arg \min_{\Phi} [\min_{\mathbf{a}_k} E(\mathbf{y}_k, \mathbf{a}_k | \Phi)]$$

$$E(\mathbf{y}_k, \mathbf{a}_k | \Phi) = -\ln p(\mathbf{y}_k | \mathbf{a}_k, \Phi) p(\mathbf{a}_k)$$

Given: $p(\mathbf{y}_k | \mathbf{a}_k, \Phi) = \frac{1}{Z_\sigma} e^{-\frac{\|\mathbf{y}_k - \Phi \mathbf{a}_k\|_2^2}{2\sigma^2}}$ $p(\mathbf{a}_k) = \prod p(a_{nk})$ $p(a_{nk}) = \frac{1}{Z_\beta} e^{-\beta S(a_{nk})}$

$$p(\mathbf{y} | \Phi) = \int_{-\infty}^{\infty} p(\mathbf{y} | \mathbf{a}, \Phi) p(\mathbf{a}) d\mathbf{a}$$

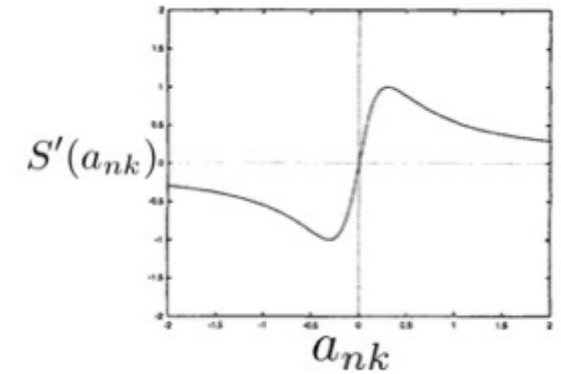
$$\hat{\Phi} = \arg \max_{\Phi} \int_{\mathbf{a}} \ln p(\mathbf{y} | \mathbf{a}, \Phi) p(\mathbf{a})$$

Obtain: $E(\mathbf{y}_k, \mathbf{a}_k | \Phi) = \|\mathbf{y}_k - \Phi \mathbf{a}_k\|_2^2 + \lambda \sum_n S(a_{nk})$

Olshausen and Field 1997 - gradients for network model

Rewriting Error function, take derivatives to find gradient

$$E(\mathbf{y}_k, \mathbf{a} | \Phi) = \sum_m (y_{mk} - \sum_n \phi_{mn} a_{nk})^2 + \lambda \sum_n S(a_{nk})$$



$$\{\hat{\Phi}, \hat{\mathbf{a}}_k\} = \arg \min_{\Phi} \left[\min_{\mathbf{a}_k} E(\mathbf{y}_k, \mathbf{a} | \Phi) \right]$$

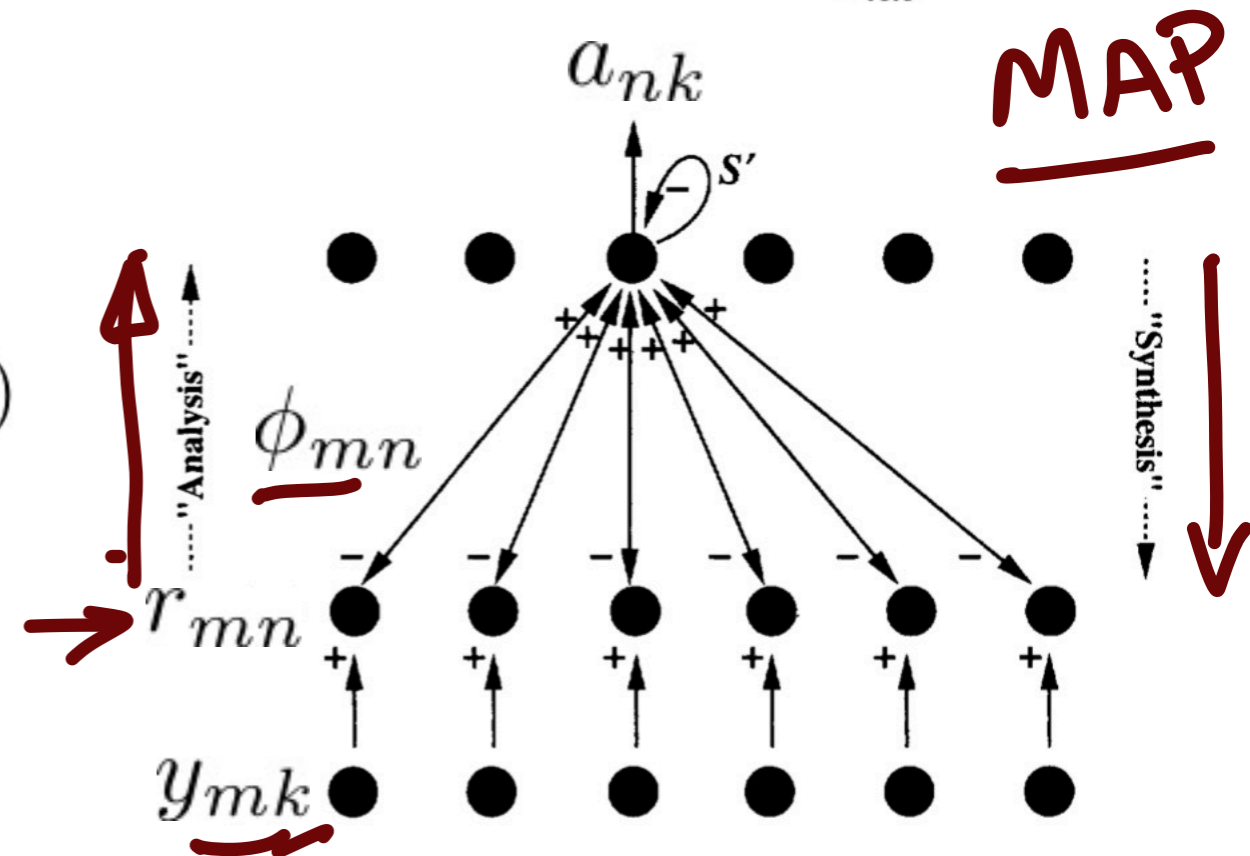
Update to a_{nk} with network (inner loop)

$$\dot{a}_{nk} = - \frac{dE}{da_{nk}} = \sum_m \phi_{mn} r_{mn} - \lambda S'(a_{nk})$$

with $r_{mn} = y_{mk} - \sum_n \phi_{mn} a_{nk}$

Update to ϕ_{mn} with gradient descent (outer loop)

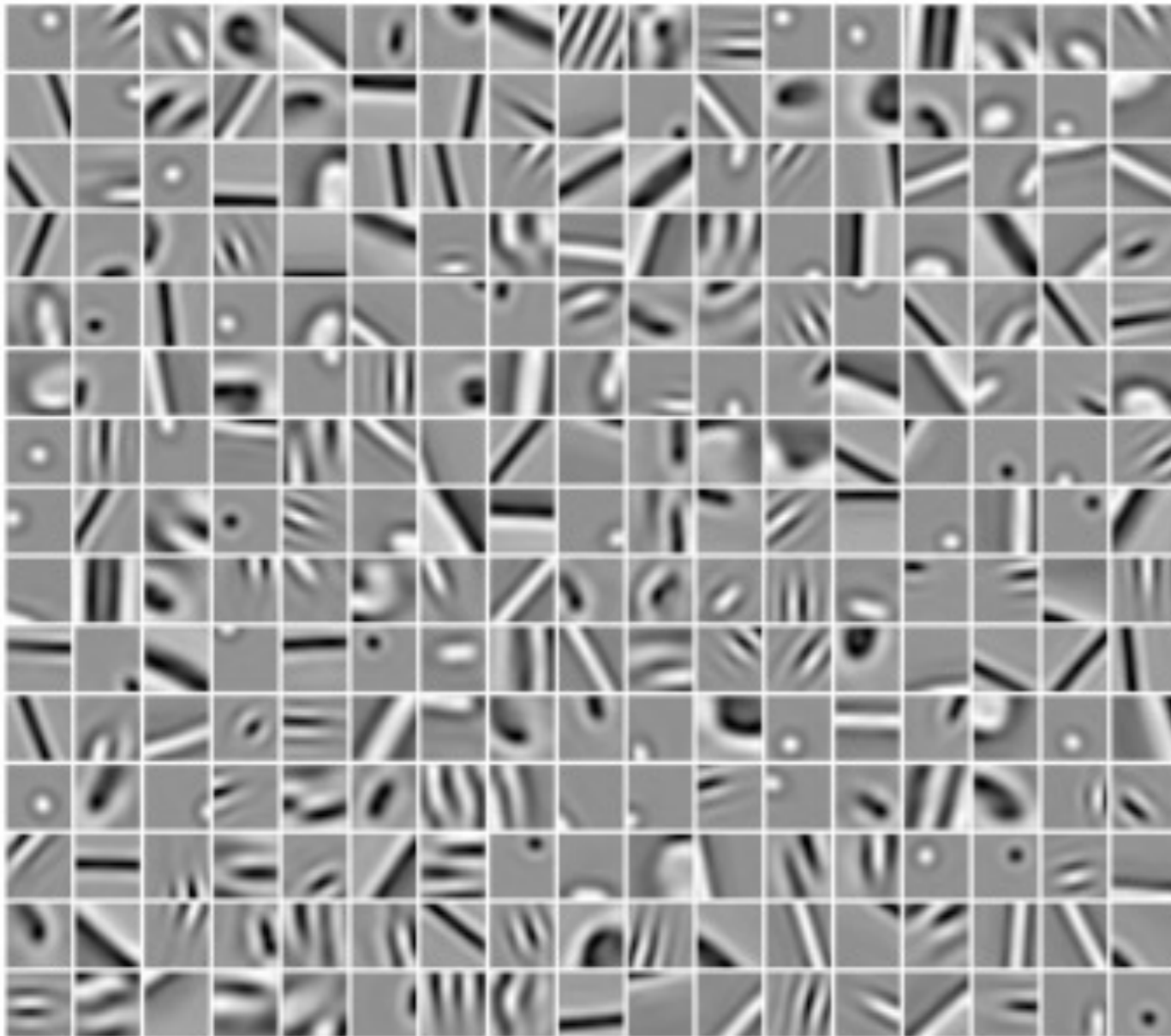
$$\Delta \phi_{mn} = \eta \langle a_{nk} r_{mn} \rangle$$



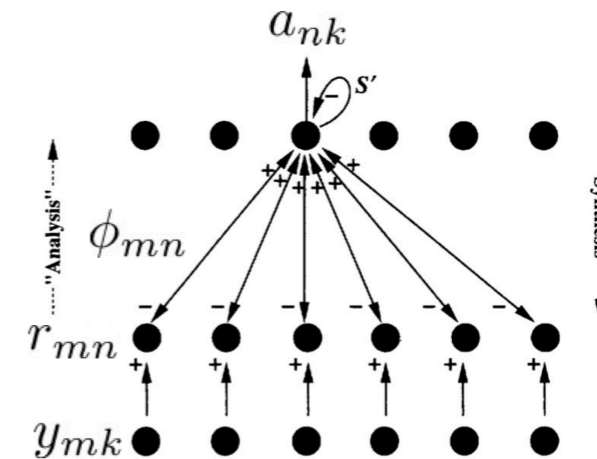
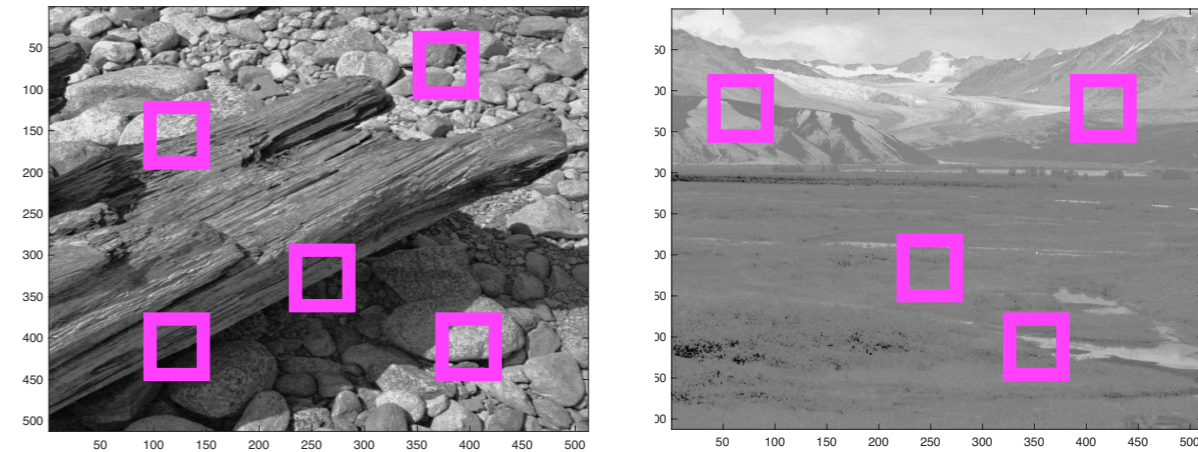
"Hebbian" update

From Olshausen '97 method, obtain dictionary atoms that resemble cells from mammalian visual cortex

Dictionary elements ϕ_k



Natural image patches



"Hebbian" update



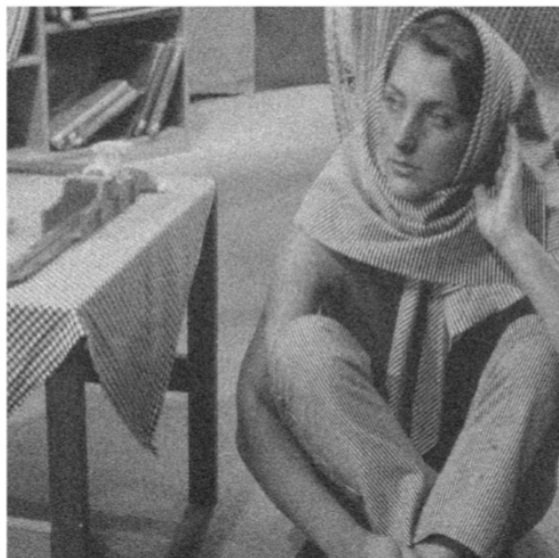
Nice to have atoms like cells, but what else is dictionary learning useful for?

Nice to have atoms like cells, but what else is dictionary learning useful for?

Image restoration tasks

Denoising

Noisy Image (22.1307 dB, $\sigma=20$)



Denoised Image Using Adaptive Dictionary (30.8295 dB)



Elad 2006

Inpainting (a.k.a. matrix completion)



Mairal 2009

Olshausen and Field 1997 - gradients for network model

Can be rephrased with Laplacian prior

$$\hat{\Phi} = \arg \min_{\Phi} \sum_k \min_{\mathbf{a}_k} \{ \|\Phi \mathbf{a}_k - \mathbf{y}_k\|_2^2 + \lambda \|\mathbf{a}_k\|_1 \}$$

"Cauchy"

"Laplacian"

Coefficients calculated using gradient descent, then dictionary updated by

$$\Phi^{(i+1)} = \Phi^{(i)} - \eta \sum_k (\Phi^{(i)} \mathbf{a}_k - \mathbf{y}_k) \mathbf{a}_k^T$$

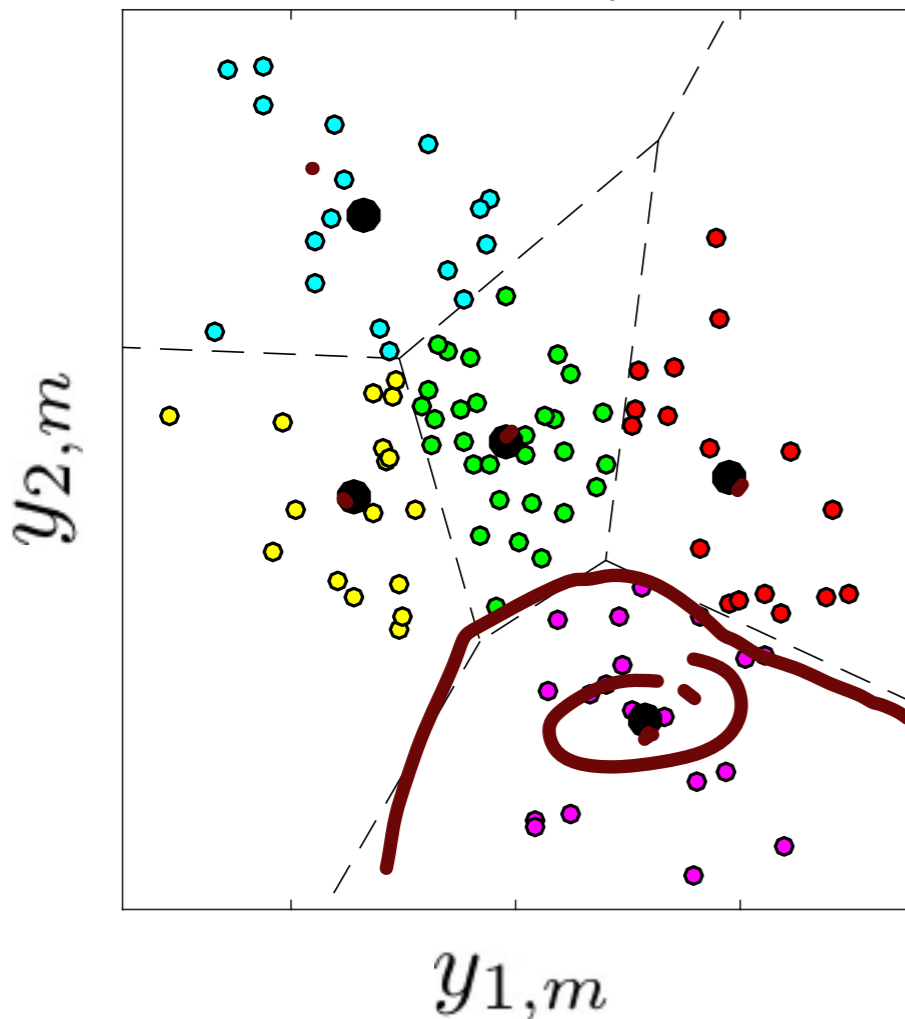
.....

This idea of iterative refinement is familiar: solving for coefficients, then updating basis functions

K-means.

Vector Quantization and K-means

2D example



Vector quantization (VQ): means of compressing a set of data observations $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_M]$ using a nearest neighbor metric with codebook $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_N]$

$$R_n = \{i \mid \forall l \neq n, \|\mathbf{y}_i - \mathbf{c}_n\|_2 < \|\mathbf{y}_i - \mathbf{c}_l\|_2\}$$

$$S_n(\mathbf{y}) = \begin{cases} 1 & \text{if } \mathbf{y} \in R_n \\ 0 & \text{otherwise,} \end{cases} \quad \hat{\mathbf{y}}_m = \sum_{i=1}^N S_i(\mathbf{y}_m) \mathbf{c}_i$$

K-means: finds optimal codebook for VQ 

Given: training vectors $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_M] \in \mathbb{R}^{K \times M}$

Initialize: index $i = 0$, codebook $\mathbf{C}^0 = [\mathbf{c}_1^0, \dots, \mathbf{c}_N^0] \in \mathbb{R}^{K \times N}$,
 MSE^0

I: Update codebook

1. Partition \mathbf{Y} into N regions (R_1, \dots, R_N) by

$$R_n = \{i \mid \forall l \neq n, \|\mathbf{y}_i - \mathbf{c}_n^i\|_2 < \|\mathbf{y}_i - \mathbf{c}_l^i\|_2\}$$

2. Make code vectors centroids of \mathbf{y}_j in partitions R_n

$$\mathbf{c}_n^{i+1} = \frac{1}{|R_n^i|} \sum_{j \in R_n^i} \mathbf{y}_j$$

II. Check error

1. Calculate MSE^{i+1} from updated codebook \mathbf{C}^{i+1}

2. If $|\text{MSE}^{i+1} - \text{MSE}^i| < \eta$

$i = i + 1$, return to I

else

end

Relationship to sparse coding

$T=1$

Sparse processor

$$\hat{\mathbf{x}}_m = \arg \min_{\mathbf{x}_m} \underbrace{\|\mathbf{y}_m - \mathbf{Q}\mathbf{x}_m\|_2}_{\text{L2 norm}} \text{ subject to } \underbrace{\|\mathbf{x}_m\|_0}_{\text{L0 norm}} \leq T \leftarrow$$

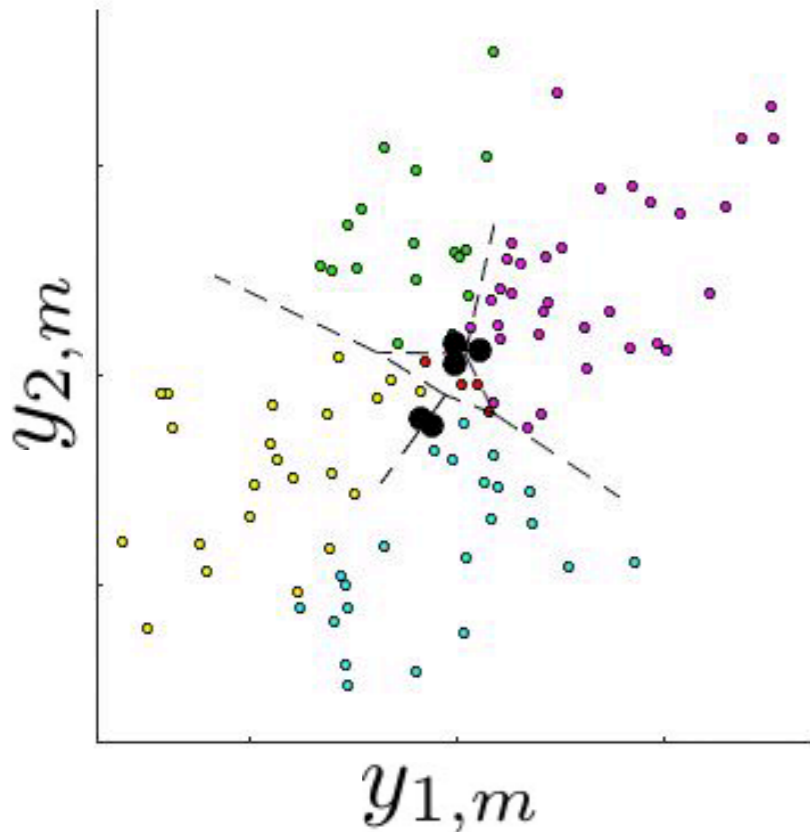
VQ operators

$$R_n = \{i \mid \forall l \neq n, \|\mathbf{y}_i - \mathbf{c}_n\|_2 < \|\mathbf{y}_i - \mathbf{c}_l\|_2\} \quad \hat{\mathbf{y}}_m = \sum_{i=1}^N S_i(\mathbf{y}_m) \mathbf{c}_i \quad S_n(\mathbf{y}) = \begin{cases} 1 & \text{if } \mathbf{y} \in R_n \\ 0 & \text{otherwise,} \end{cases}$$

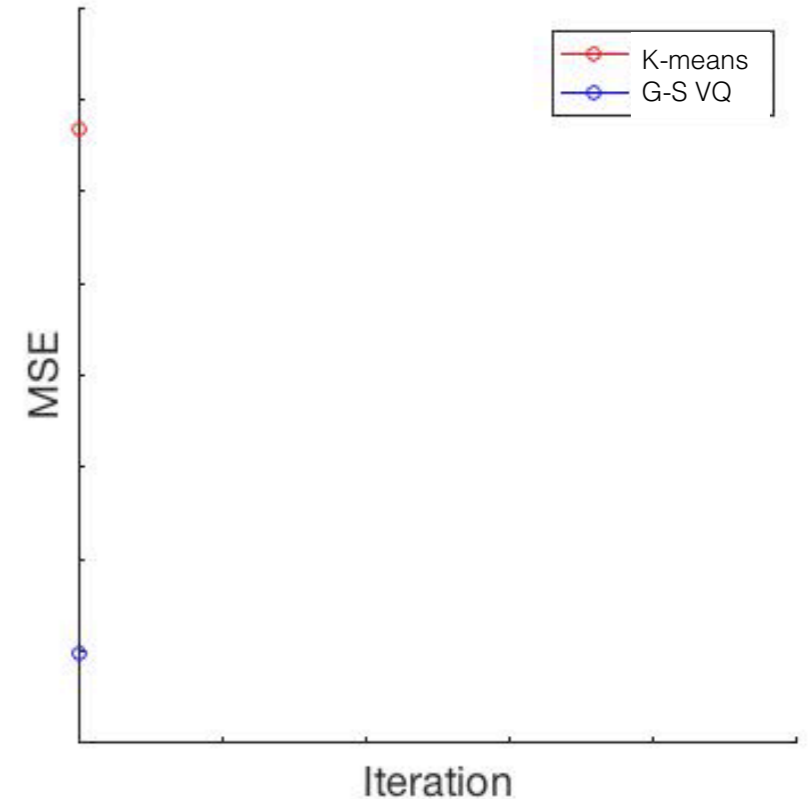
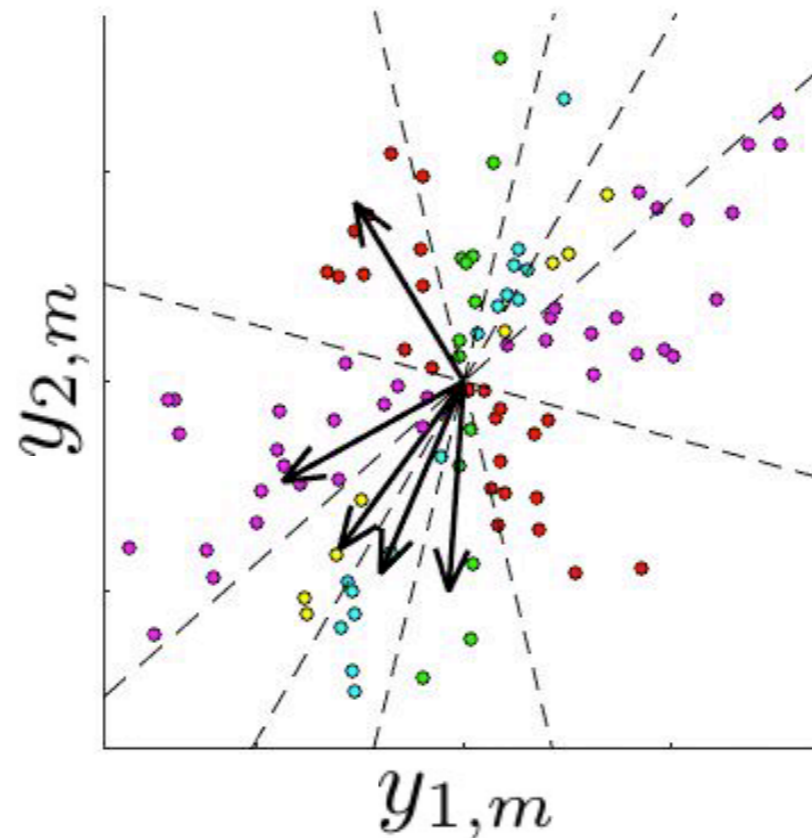
Dictionary learning objective

$$\min_{\mathbf{Q}} \{ \min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 \text{ subject to } \forall m, \|\mathbf{x}_m\|_0 \leq T \}$$

K-means



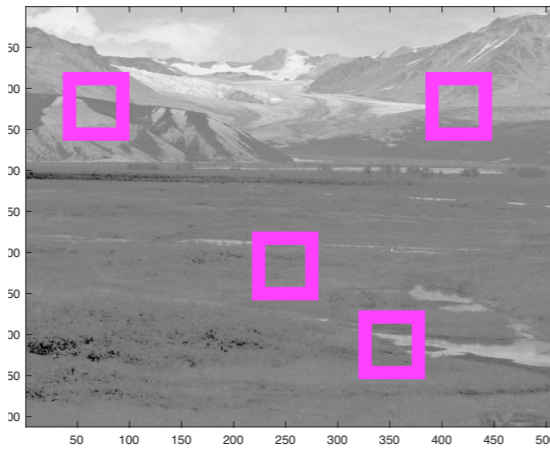
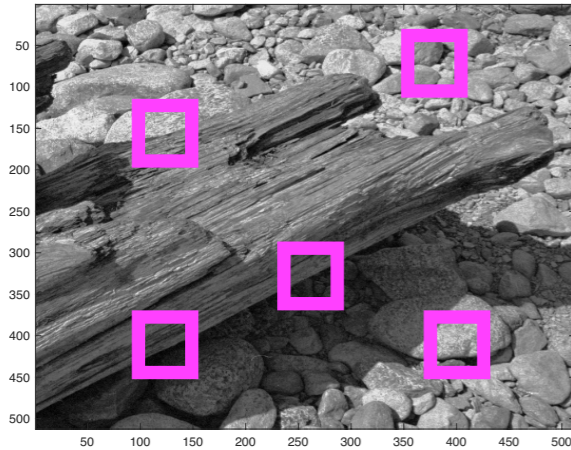
Gain-shape VQ



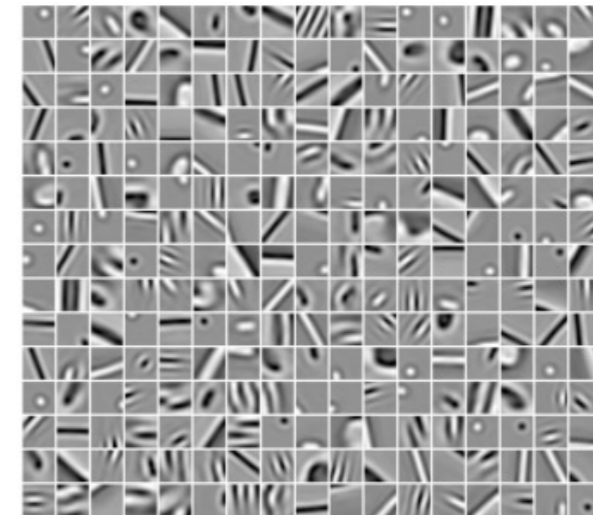
Legend:
 - K-means (red circle)
 - G-S VQ (blue circle)

Background: a basic dictionary learning framework

Given set of patches $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_I]$, learn dictionary \mathbf{D} describing them



Patches shown in **magenta**



Dictionary \mathbf{D}

Dictionary learning objective

$$\min_{\mathbf{D}} \{ \min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \text{ subject to } \|\mathbf{x}_i\|_0 \leq T \forall i \}$$

Objective solved as simple optimization problem *"alt. min."*

- 1. Solve for sparse coefficients $\mathbf{X} = [\hat{\mathbf{x}}_{\ell_0,1}, \dots, \hat{\mathbf{x}}_{\ell_0,I}]$ using sparse solver
- 2. Solve for dictionary \mathbf{D} using sparse coefficients from step (1)..... repeat until convergence

MOD algorithm: Extending K-means to dictionary learning problem

Method of Optimal Directions (MOD) [Engan 2000]

$$\min_{\mathbf{Q}} \{ \min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 \text{ subject to } \forall m, \|\mathbf{x}_m\|_0 \leq T \}$$

MOD algorithm:

1. COEFFICIENTS: Solve for coefficients $\mathbf{X}=[\mathbf{x}_1 \dots \mathbf{x}_i]$ for fixed \mathbf{Q} using orthogonal matching pursuit (OMP)
2. DICTIONARY UPDATE: Solve for dictionary $\mathbf{Q}=[\mathbf{q}_1 \dots \mathbf{q}_i]$, by inverting the coefficient matrix \mathbf{X} , and normalizing dictionary entries to have unit norm.

$$\hat{\mathbf{Q}} = \mathbf{Y}\mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1}$$

.... repeat until convergence

"pseudoinverse"

Simple and flexible but, a few drawbacks:

- computationally expensive to invert coefficient matrix \mathbf{X}
- since keeping coefficients in \mathbf{X} fixed during dictionary update, slow convergence

~~$\mathbf{Q}\mathbf{X} = \mathbf{Y}$~~
 $\mathbf{Q}\mathbf{X}\mathbf{X}^T = \mathbf{Y}\mathbf{X}^T$
 $\mathbf{Q} = \mathbf{Y}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}$

K-SVD algorithm

K-SVD [Aharon 2006]: Learn optimal dictionary for sparse representation of data

$$\min_{\mathbf{Q}} \left\{ \min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 \text{ subject to } \forall m, \|\mathbf{x}_m\|_0 \leq T \right\}$$

K-SVD algorithm:

1. Solve for coefficients $\mathbf{X}=[\mathbf{x}_1 \dots \mathbf{x}_i]$ for fixed \mathbf{Q} using OMP
2. Solve (1) for dictionary $\mathbf{Q}=[\mathbf{q}_1 \dots \mathbf{q}_i]$, updating both \mathbf{Q} and \mathbf{X} from the SVD of representation error

$$\|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F = \left\| \left(\mathbf{Y} - \sum_{j \neq k} \mathbf{q}_j \mathbf{x}_T^j \right) - \mathbf{q}_k \mathbf{x}_T^k \right\|_F$$

SVD $\rightarrow \|\mathbf{E}_k\|_F = \|\mathbf{q}_k \mathbf{x}_T^k\|_F$

update $\mathbf{q}_k, \mathbf{x}_k$ by SVD

$$\mathbf{E}_k^e = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

$$\mathbf{q}_k = \mathbf{U}(:, 1), \mathbf{x}_T^k = \mathbf{V}(:, 1)\mathbf{S}(1, 1)$$

.... repeat until convergence

2D example

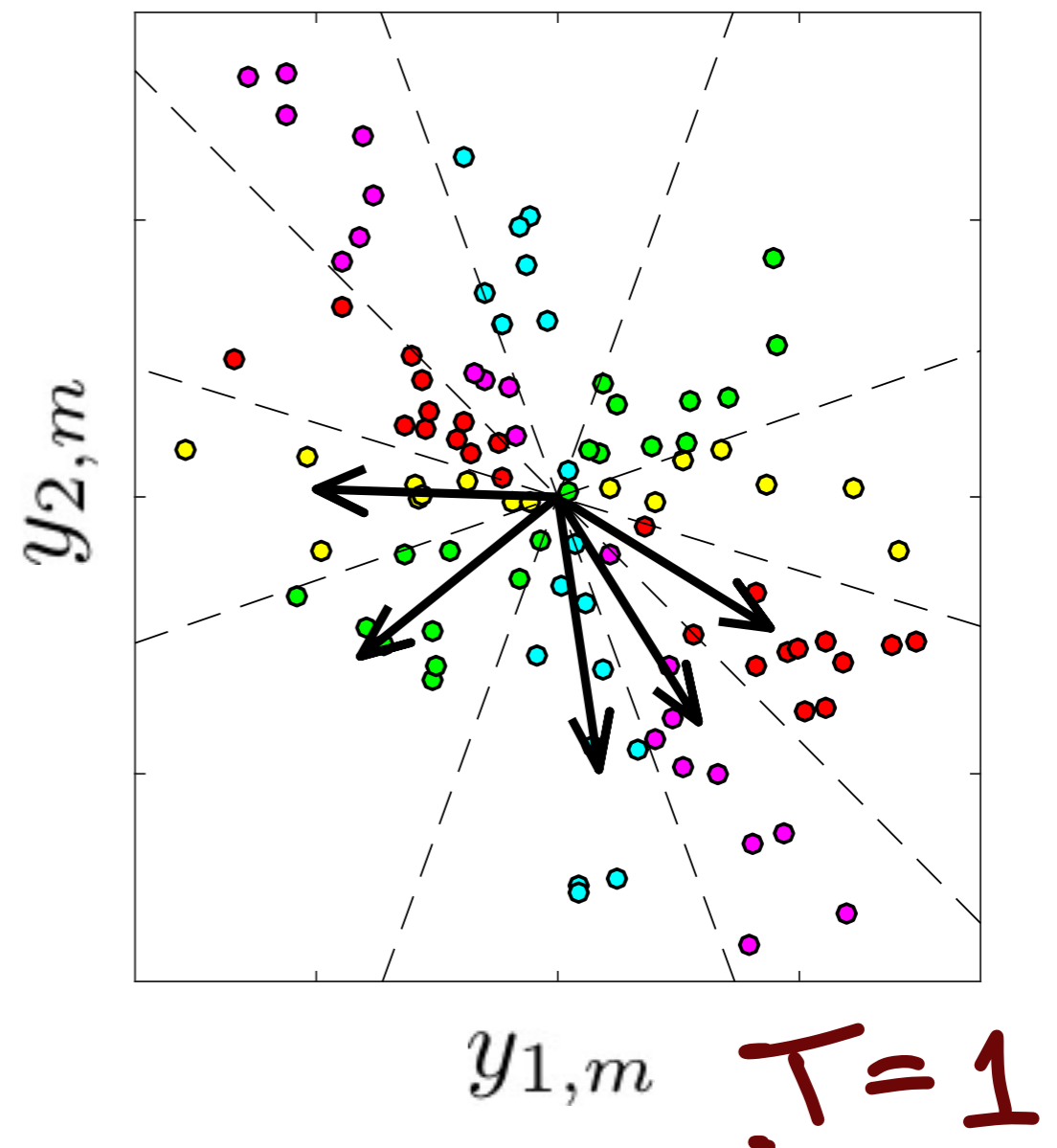


Image restoration tasks

Denoising

Noisy Image (22.1307 dB, $\sigma=20$)



Denosed Image Using Adaptive Dictionary (30.8295 dB)



x_i ~~64x64~~

256 Flad 2006

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}_i\|_0 \leq s$$

Inpainting (a.k.a. matrix completion)



Mairal 2009

256
~60k

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{M}_i(\mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i)\|_2^2 + \lambda \psi(\boldsymbol{\alpha}_i),$$

//

Image restoration tasks

Denoising: learning from noisy image patches for specific image



Noisy Image (22.1307 dB, $\sigma=20$)



Denoised Image Using Adaptive Dictionary (30.8295 dB)

Elad 2006

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}_i\|_0 \leq s$$

$$\{\hat{\boldsymbol{\alpha}}_{ij}; \hat{\mathbf{X}}, \hat{\mathbf{D}}\}$$

$$= \min \|\mathbf{X} - \mathbf{Y}\|_2^2 + \sum_i \|\mathbf{R}_i \mathbf{X} - \mathbf{D}\boldsymbol{\alpha}_i\|_2 + \lambda \|\boldsymbol{\alpha}_i\|_0$$

Solved using block-coordinate descent algorithm (also two steps):

60K (1) $\hat{\boldsymbol{\alpha}}_{ij} = \arg \min_{\boldsymbol{\alpha}} \mu_{ij} \|\boldsymbol{\alpha}\|_0 + \|\mathbf{D}\boldsymbol{\alpha} - \mathbf{x}_{ij}\|_2^2$

$$\Rightarrow \{\hat{\mathbf{D}}, \hat{\boldsymbol{\alpha}}_{ij}\} = \min_{\mathbf{D}} \left\{ \min_{\boldsymbol{\alpha}} \dots \right\}$$

(2) $\hat{\mathbf{X}} = \arg \min_{\mathbf{X}} \lambda \|\mathbf{X} - \mathbf{Y}\|_2^2 + \sum_{ij} \|\mathbf{D}\hat{\boldsymbol{\alpha}}_{ij} - \mathbf{R}_{ij}\mathbf{X}\|_2^2$

$$\hat{\mathbf{X}} = \left(\lambda \mathbf{I} + \sum_{ij} \mathbf{R}_{ij}^T \mathbf{R}_{ij} \right)^{-1} \left(\lambda \mathbf{Y} + \sum_{ij} \mathbf{R}_{ij}^T \mathbf{D}\hat{\boldsymbol{\alpha}}_{ij} \right)$$

Why not just use neural networks?

Burger 2012: Multi-layer perceptron competes with state of art denoising algorithms, using 362 million training samples (~one month of GPU time)

... at least in geoscience (seimsics, ocean acoustics) we rarely have this much training data

Adaptive image denoising-like

→ Handcrafted



Adhere to existing algorithm architecture, few learnable parameters

Pros:

- Likely more generalizable
- Less training data needed
- Natural initialization (from standard algorithm settings)

Cons:

- Reduced chance for optimal performance

MLP-like

Blackbox ←

Deviate from existing algorithm design, many learnable parameters

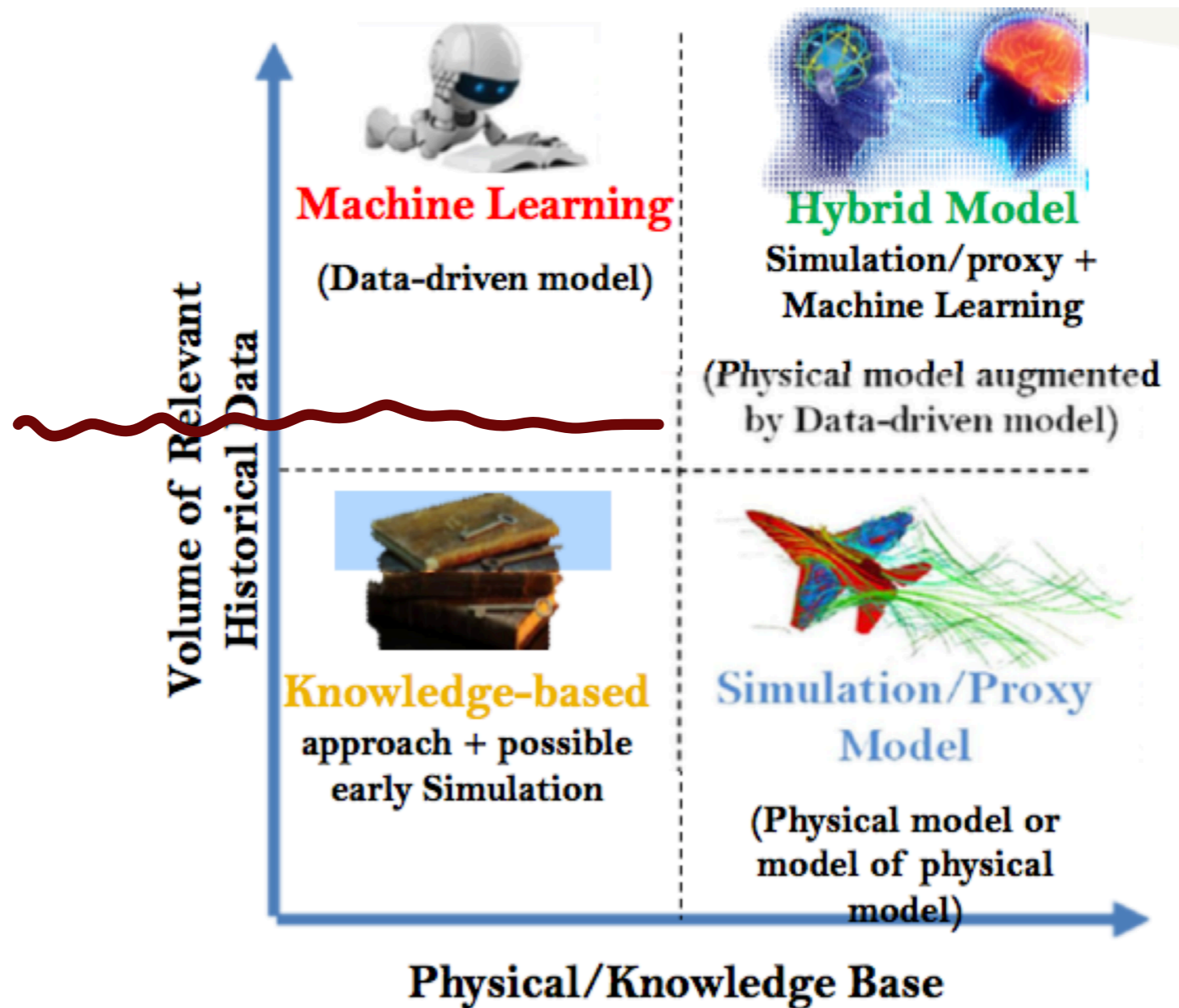
Pros:

- Increased chance of optimal performance given sufficient training data

Cons:

- Maybe less generalizable
- More training data needed
- No algorithm to potentially guide initialization

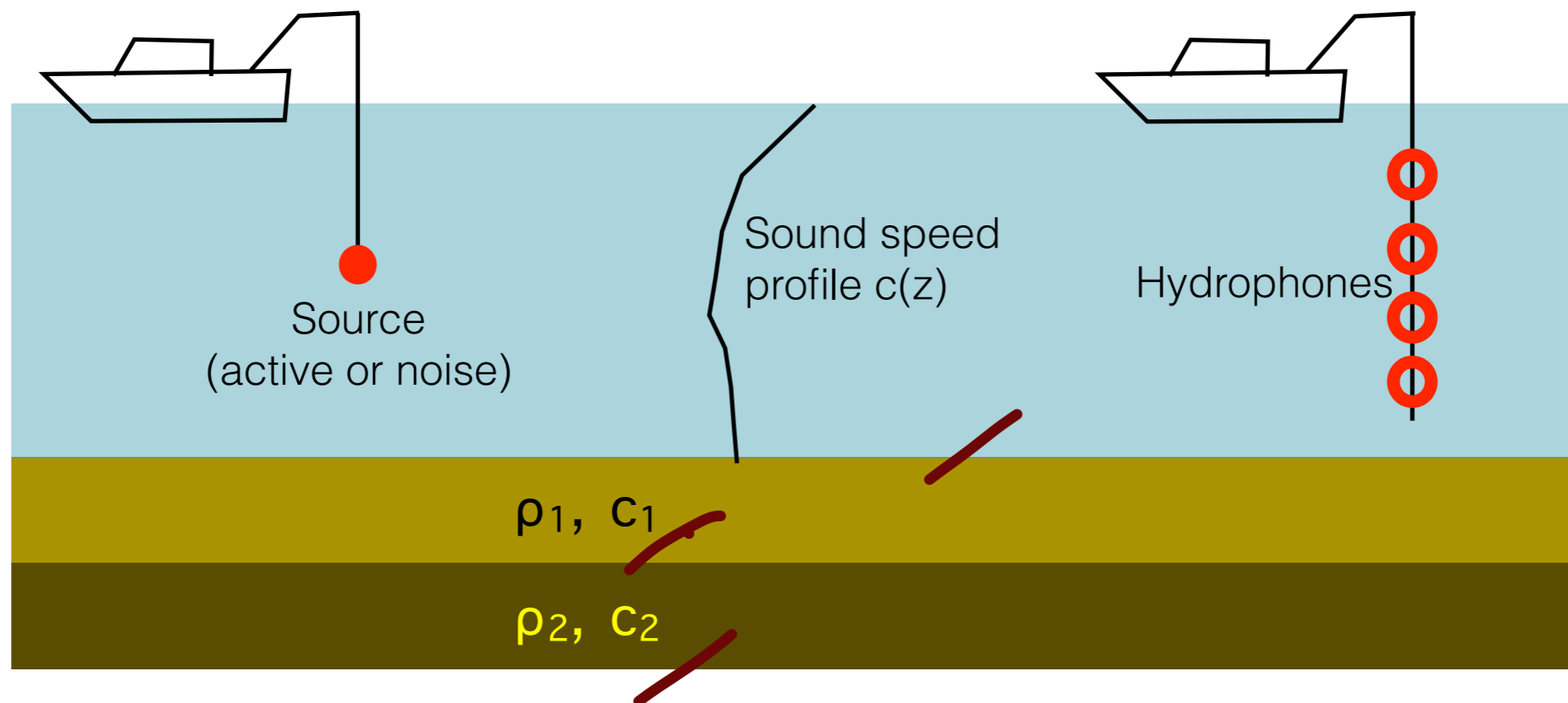
Why not just use neural networks? (cont'd)



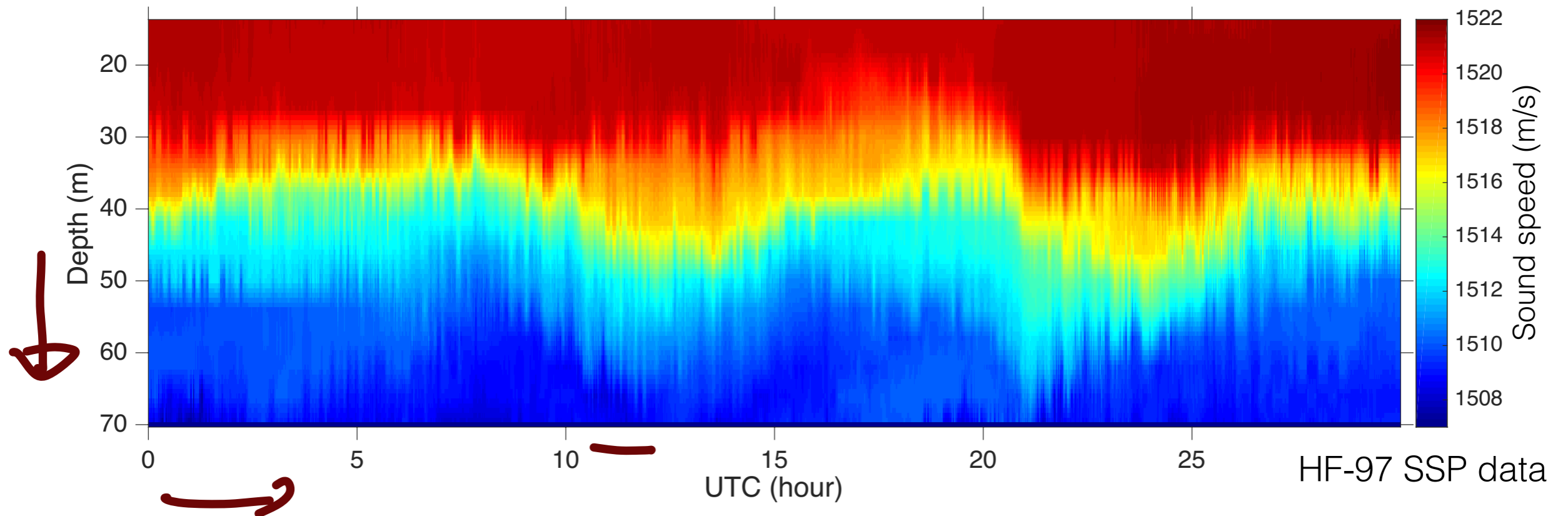
Dictionary learning of ocean sound speed profiles

Bianco and Gerstoft 2017

- Acoustic observations from ocean contain information about ocean environment
- The inversion of environment parameters is limited by physics and signal processing assumptions



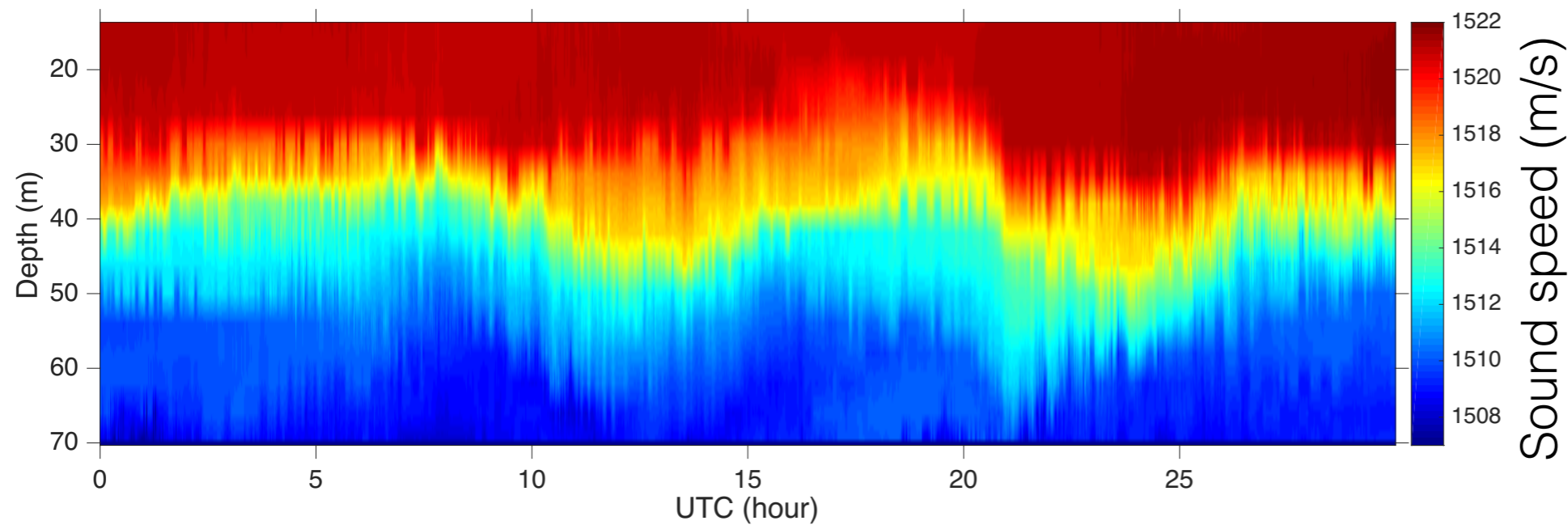
Sound speed profiles



- Sound speed profiles (SSPs) in the ocean are often highly variable with fine scale fluctuations
- Acoustic inversion of SSPs is ill-posed and traditionally regularized using EOFs (=PCA in this case)
- Dictionaries obtained via unsupervised learning may provide better representation of SSP dynamics

Dictionary learning of sound speed profiles

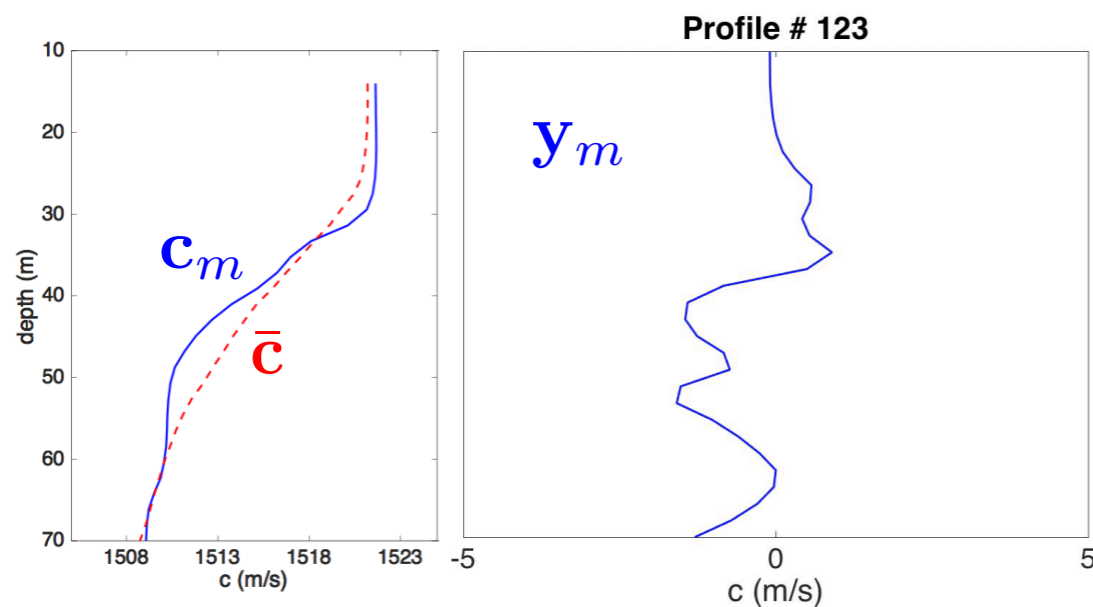
Bianco and Gerstoft JASA 2017 (published)



HF-97 Experiment

- 30 hours of SSP data
- Used 1000 profiles for dictionary learning
- $K = 30$ point SSP's (interpolated from 15 measurements)

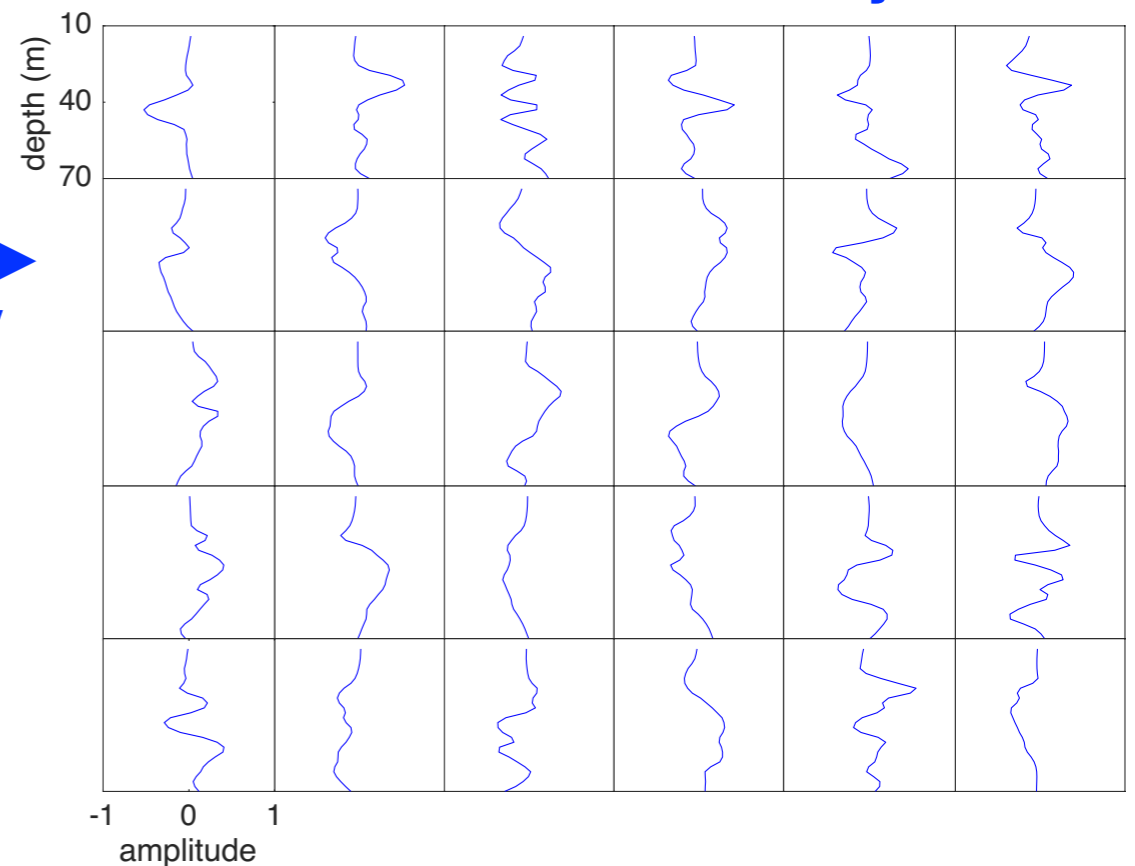
SSP Variability



$$\mathbf{y}_m = \mathbf{c}_m - \bar{\mathbf{c}}$$

Dictionary Learning

'Learned Dictionary'



Example: Denoising alphabet with K-SVD algorithm

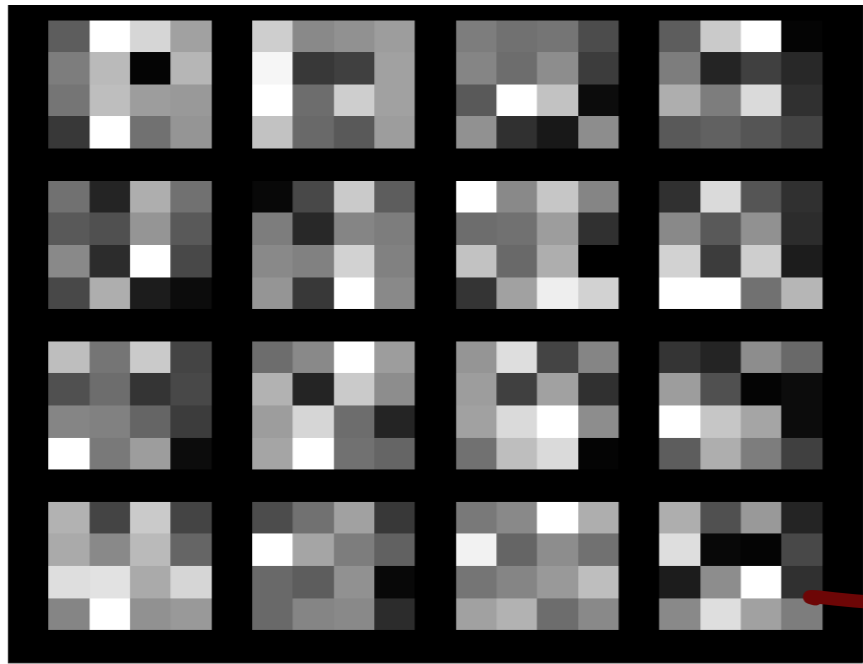
True alphabet



Recovered alphabet (no noise, K-SVD)



Recovered alphabet (no noise, PCA)



Recovered alphabet (noise std = .5, K-SVD)

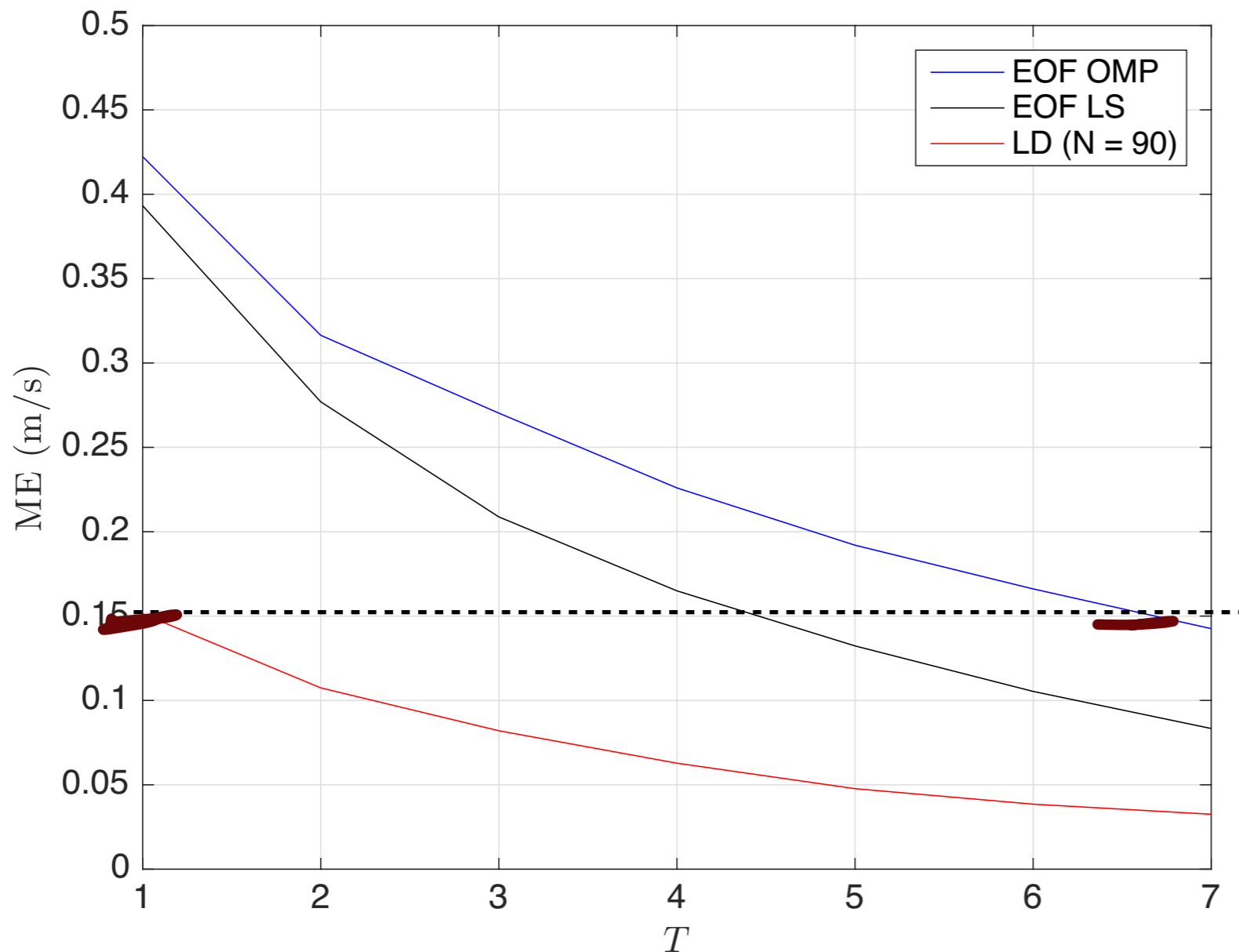


$$Y = [y_1 \dots y_n]$$

$$\|Y - Q \times X\|_2^2 \quad \& \quad \|X_i\| = T \quad \&$$

SSP reconstruction error using Dictionary Learning

Based on 1000 profiles from HF-97



LS: Least squares
OMP: Sparse processor

Mean Error (ME):

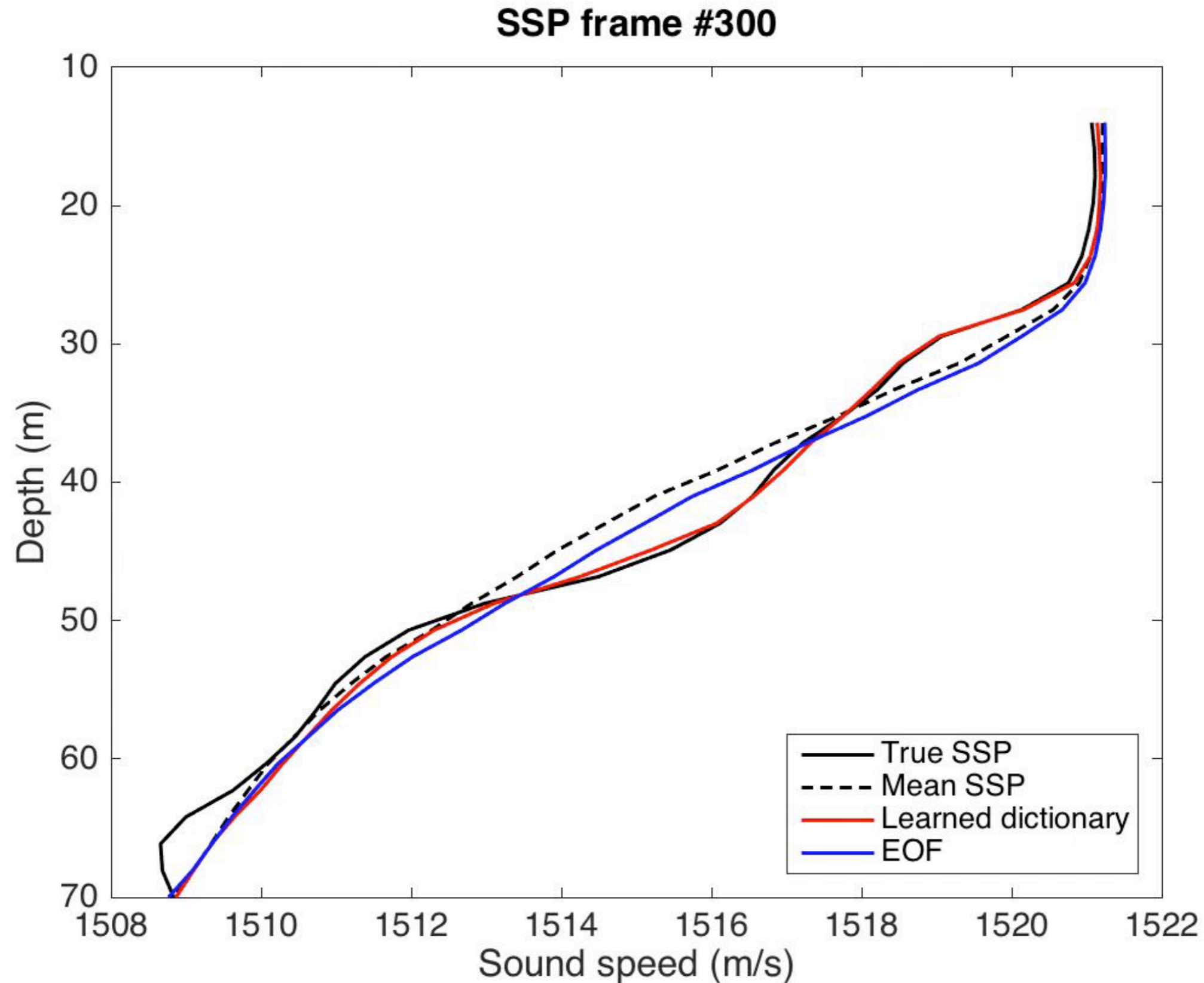
$$ME = \frac{1}{KM} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_1$$

"Compression tech."

- One entry from Learned Dictionary fits SSP data better than 6 EOFs
- Learned dictionary (LD) reconstruction error less than 50% of EOF error

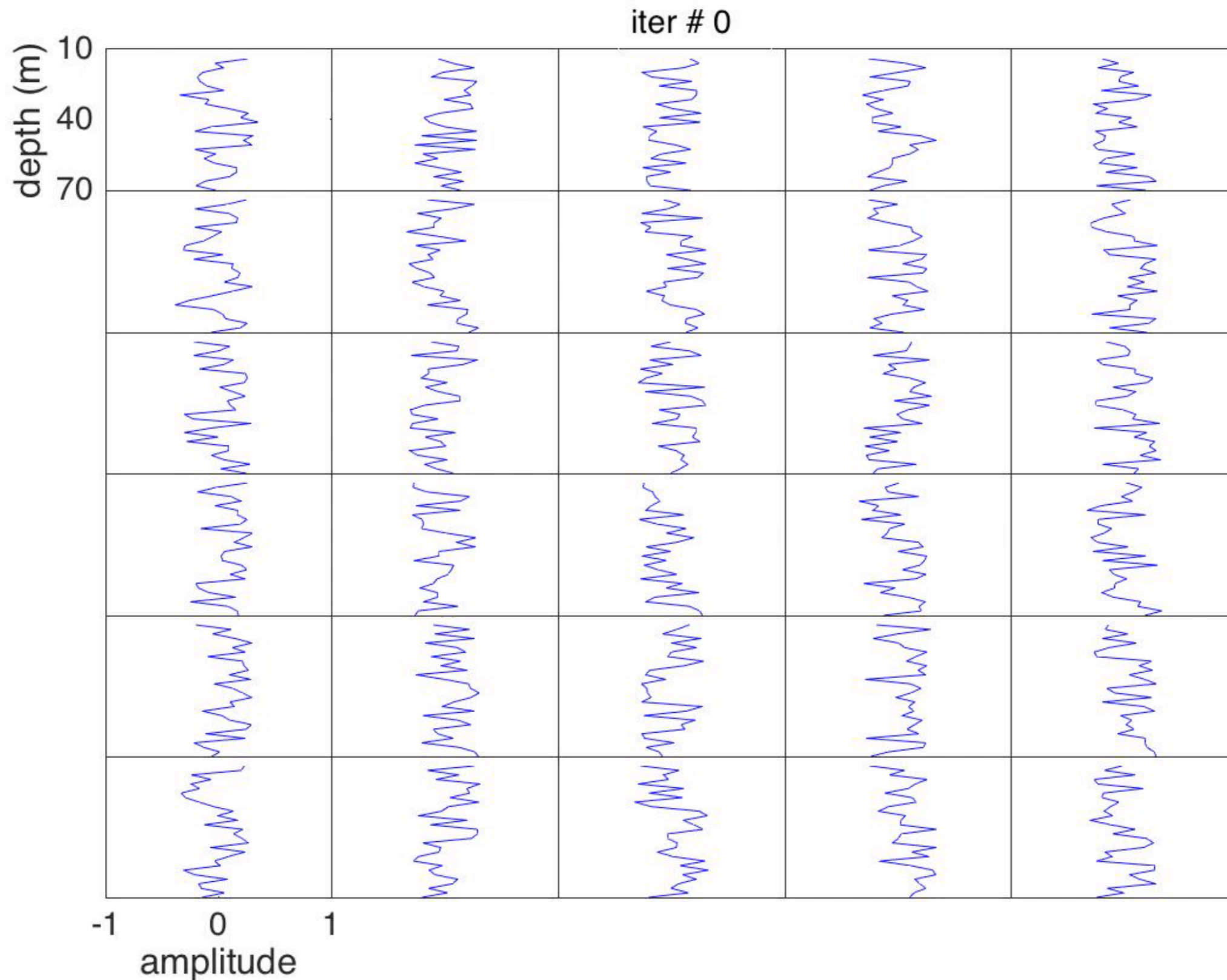
SSP reconstruction using Dictionary Learning

HF-97: One coefficient from Learned Dictionary vs. One EOF coefficient



Learning dictionary from HF-97 SSP variation

Q random initialized, converges within 15 iterations



LD solution space much smaller than EOFs

Inversion for SSP:

Assuming a potentially non-linear mapping:

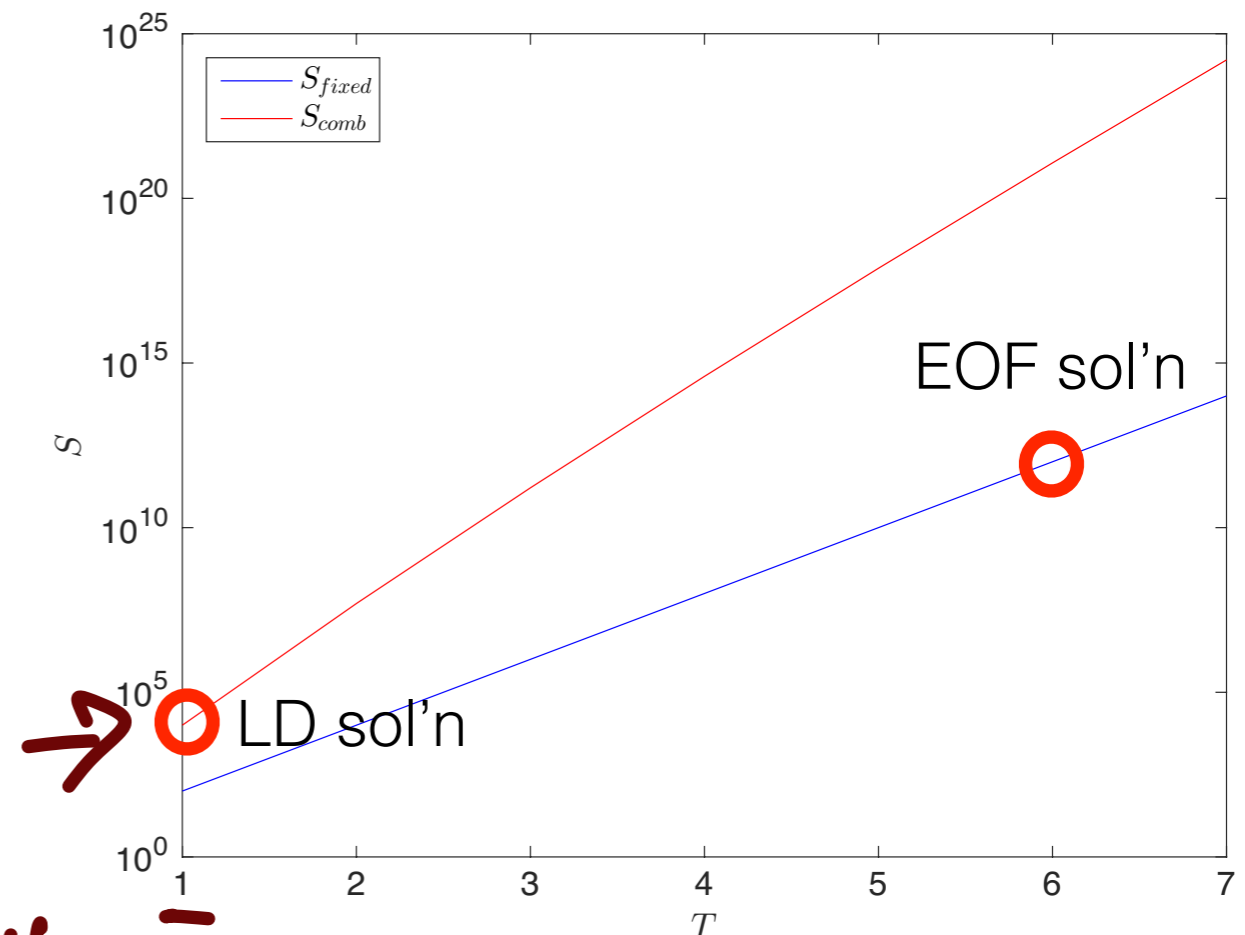
- EOF solution: T leading order coefficients (fixed indices)

$$S_{\text{fixed}} = H^T \quad \leftarrow 100^6$$

- LD solution: T -non-zero coefficients (combinatorial indices)

$$S_{\text{comb}} = H^T \frac{N!}{T!(N-T)!}$$

$11 \times 11_0$ - combi-indices \rightarrow



- Since 6 EOFs or 1 LD entry required, if coefficients discretized in $H=100$ coefficients number of possible solutions are

EOFs: $S_{\text{fixed}} = 10^{12}$ solutions

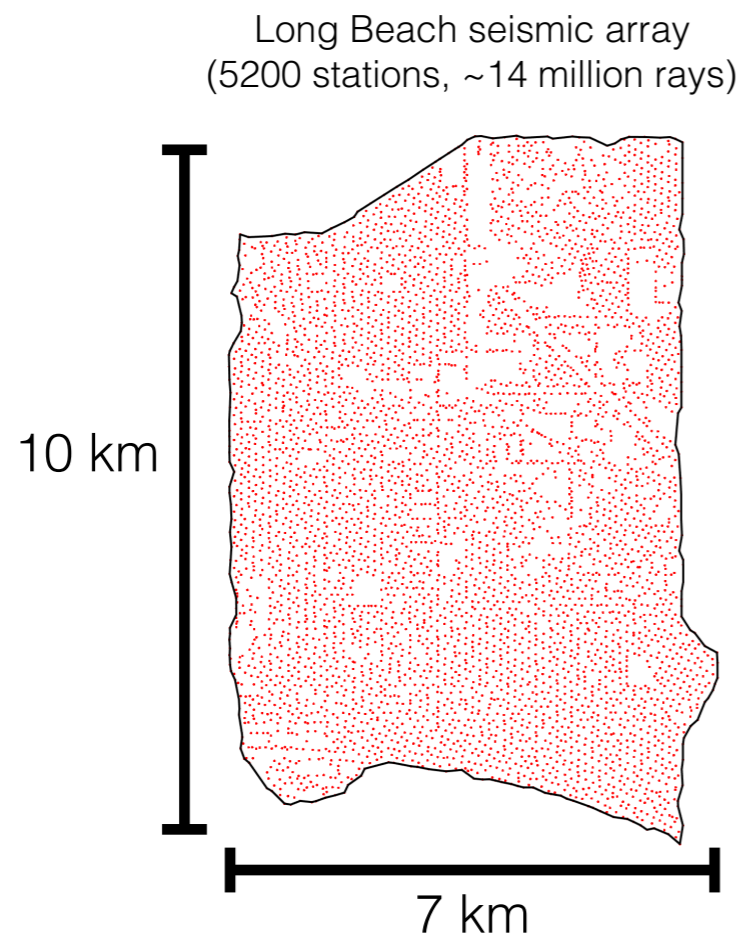
LD: $S_{\text{comb}} = 10^4$ solutions

10^{12}

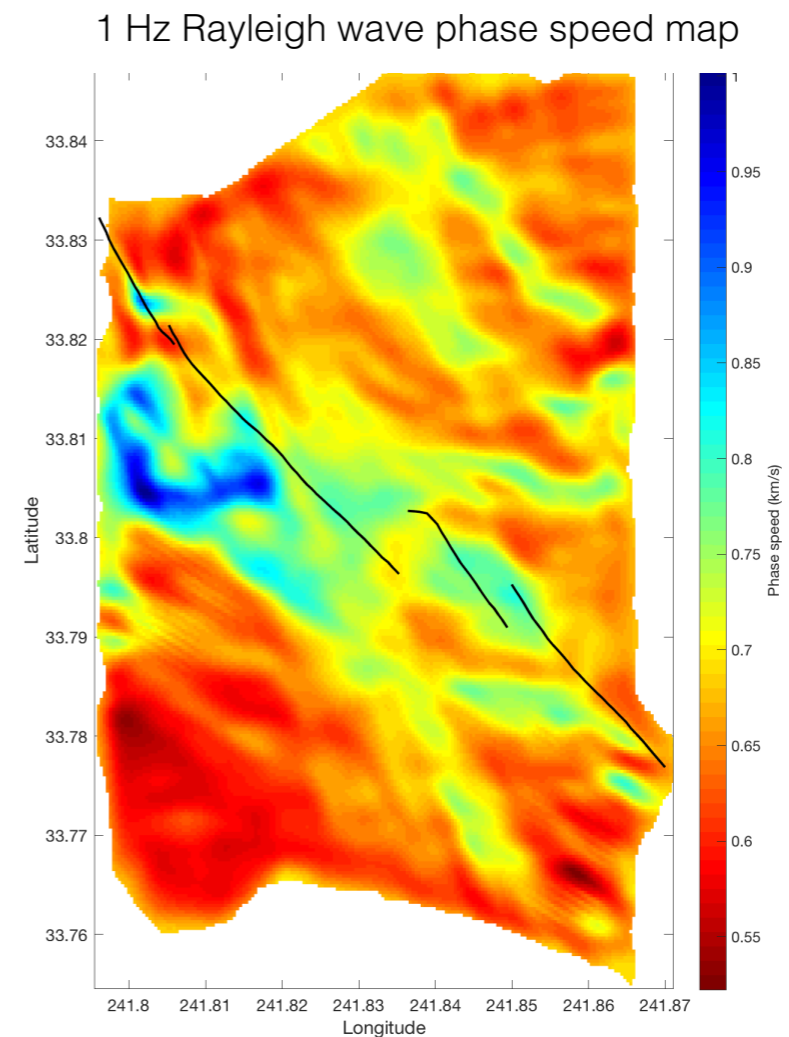
Dictionary learning in travel time tomography

Bianco and Gerstoft 2018

- The Earth contains both smooth and discontinuous variations in slowness (e.g. Moho, faults) at multiple spatial scales
- Most existing travel time inversion methods are ad hoc: regularize inversion assuming exclusively smooth or discontinuous slownesses
- Propose locally-sparse 2D travel time tomography (LST) method with three main ingredients:
 - Sparsity constraint on slowness patches
 - Dictionary learning (unsupervised machine learning)
 - Damped least squares regularization on overall slowness map



LST



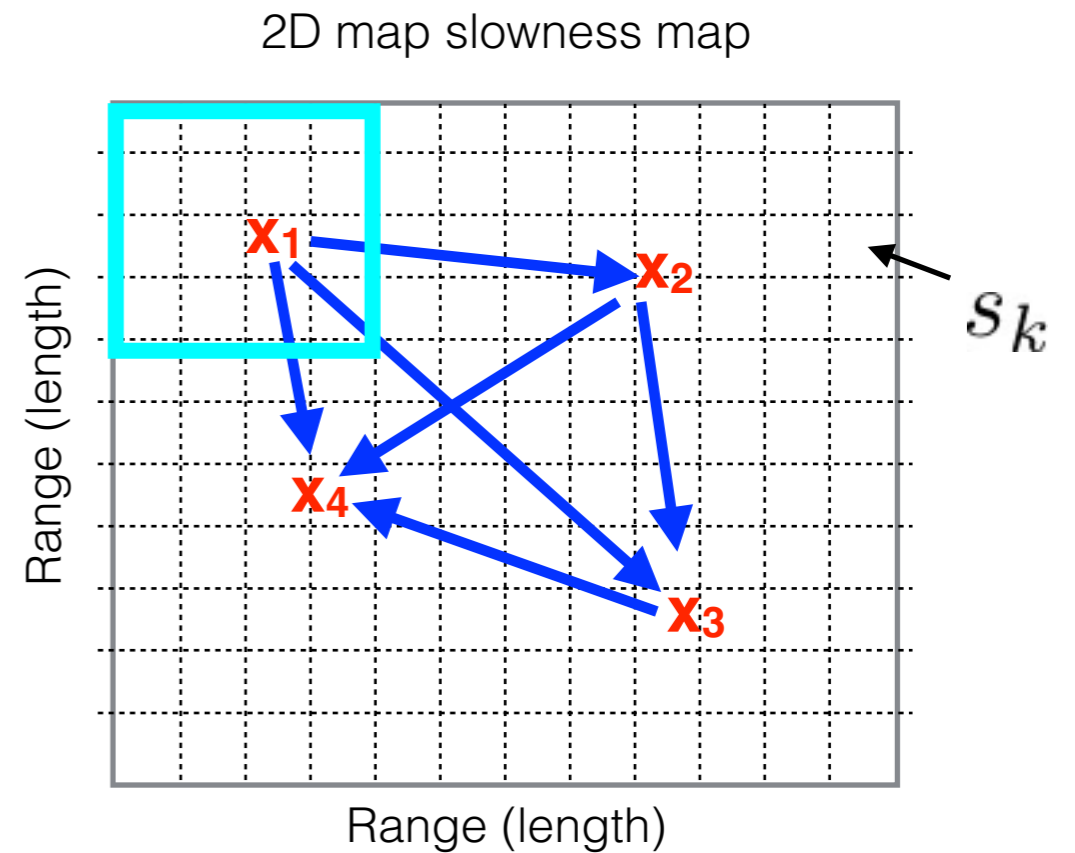
Consider simple travel time model

For slowness field, get travel time:

$$t = \frac{d}{c} = ds$$

c = wave speed

s = slowness



For straight-rays get simple formulation:

$$\begin{bmatrix} t_{12} \\ \vdots \\ t_{ij} \end{bmatrix} = \begin{bmatrix} \delta r_{12,1} & \dots & \delta r_{12,k} \\ \vdots & \ddots & \vdots \\ \delta r_{12,k} & \dots & \delta r_{ij,k} \end{bmatrix} \begin{bmatrix} s_1 \\ \vdots \\ s_k \end{bmatrix} \longrightarrow \boxed{\mathbf{t} = \mathbf{A}\mathbf{s}}$$

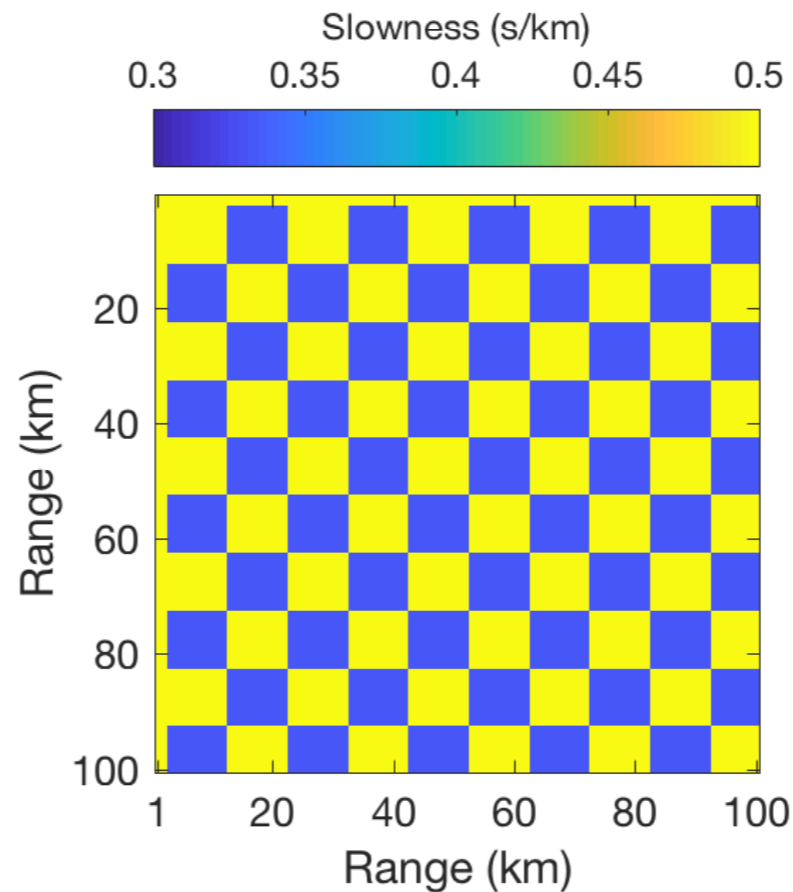
"tomography matrix" \mathbf{A}

- Propose LST tomography ingredients:

- Sparsity constraint on slowness patches
- Dictionary learning (unsupervised machine learning)
- Damped least squares regularization on overall slowness map

Proposed locally-sparse tomography (LST) basics

Synthetic "checkerboard" slowness example



LST approach three ingredients: classified as **local** and **global** models

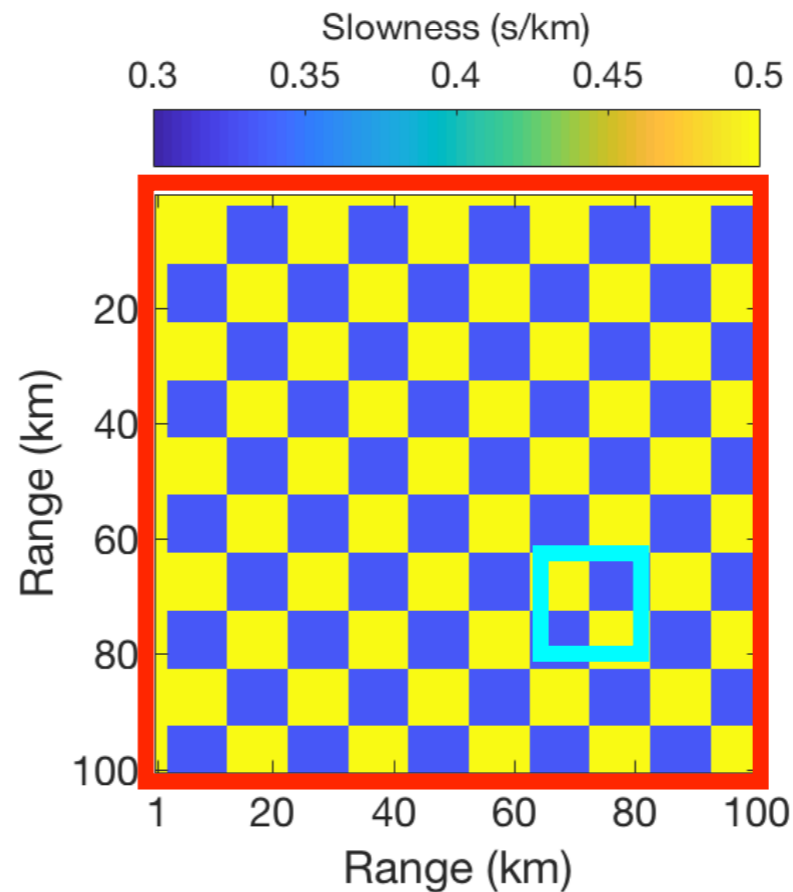
- | | |
|-----------------------|--|
| “Local” model | 1. Sparsity constraint on slowness patches |
| | 2. Dictionary learning (unsupervised machine learning) |
| “Global” model | 3. Damped least squares regularization on overall slowness map |

“Local” model: Models small-scale features as patches

“Global” model: Models larger-scale features with damped least squares

Proposed locally-sparse tomography (LST) basics

Synthetic "checkerboard" slowness example



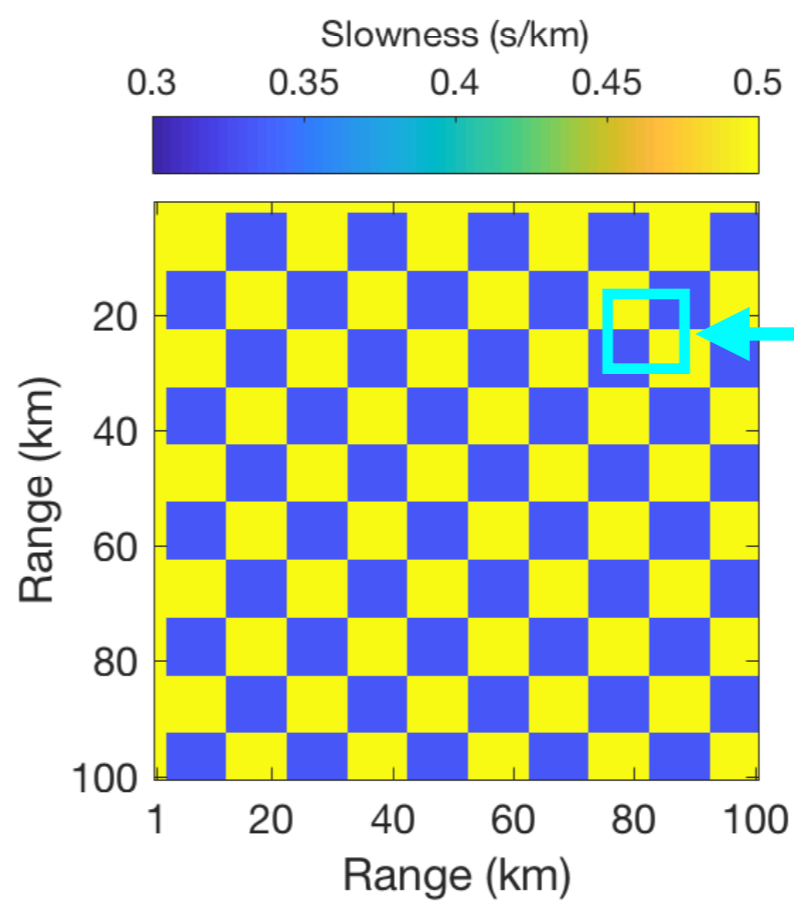
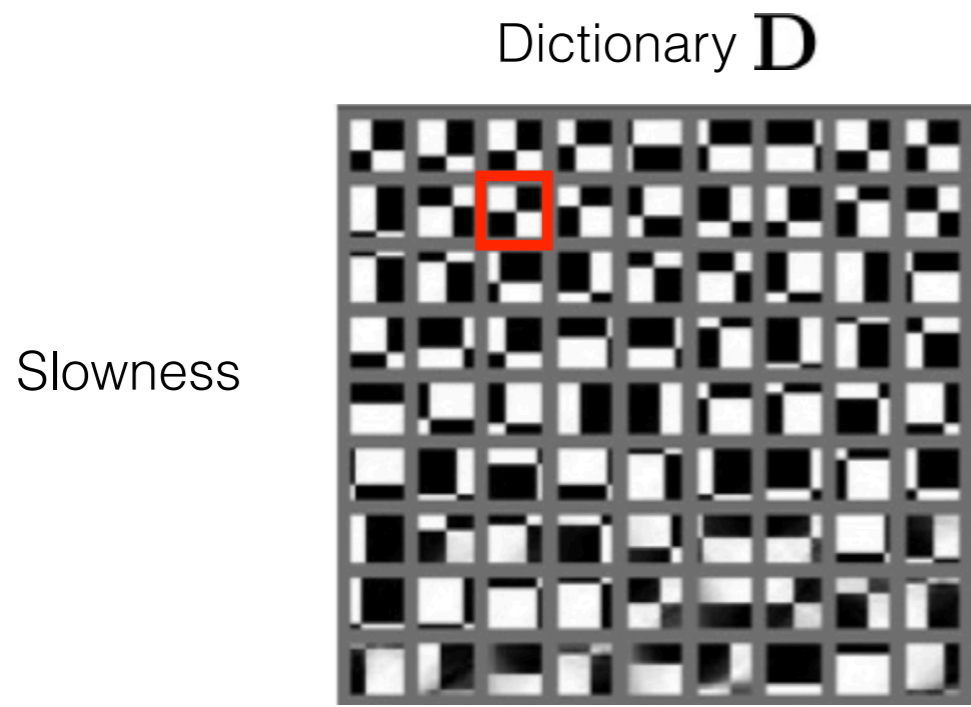
LST approach three ingredients: classified as **local** and **global** models

- | | |
|-----------------------|--|
| “Local” model | 1. Sparsity constraint on slowness patches |
| | 2. Dictionary learning (unsupervised machine learning) |
| “Global” model | 3. Damped least squares regularization on overall slowness map |

“Local” model: Models small-scale features as patches

“Global” model: Models larger-scale features with damped least squares

Local model: slowness patches related to dictionary entries



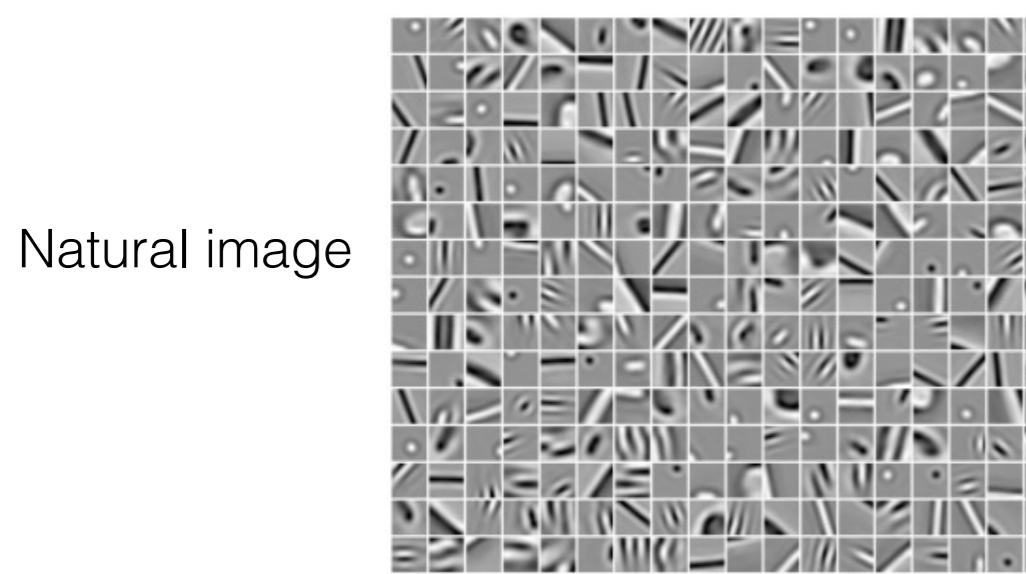
$$\mathbf{y} = \mathbf{R}_i \mathbf{s} = \mathbf{D} \mathbf{x}_i$$

$$\mathbf{R}_i \mathbf{s} = \begin{bmatrix} \blacksquare & \square \\ \square & \blacksquare \end{bmatrix} \mathbf{x}_i$$

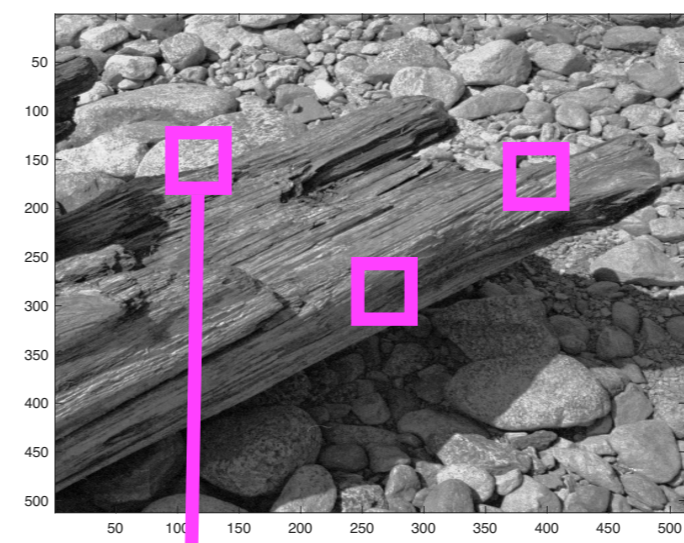
10x10 pixel patches

$$\mathbf{D} \in \mathbb{R}^{n \times Q}$$

$$\mathbf{R}_i \in \{0, 1\}^{n \times N}$$



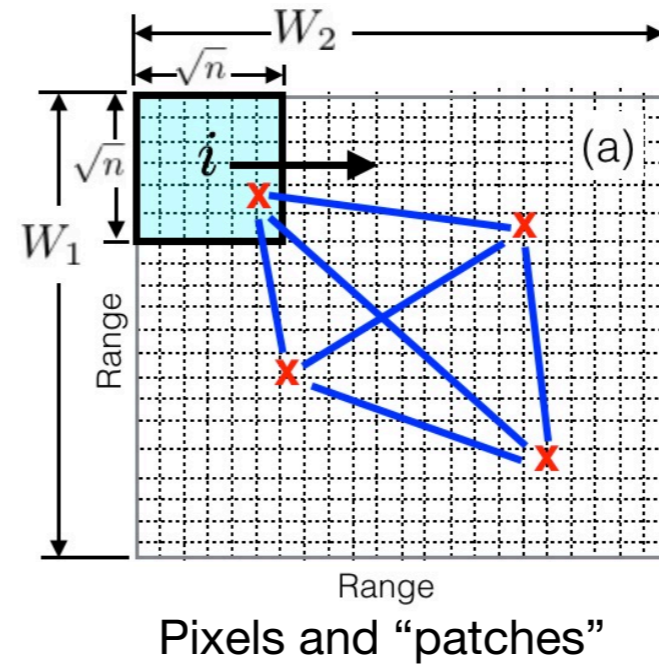
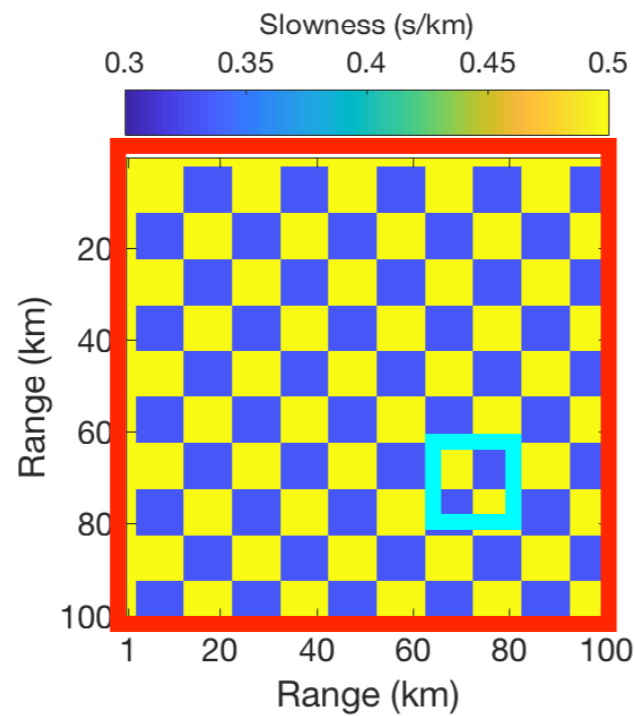
Olshausen 2009



$$\hat{\mathbf{x}}_i = \arg \min_{\mathbf{x}_i} \|\mathbf{y}_i - \mathbf{D} \mathbf{x}_i\|_2 \quad \text{subject to} \quad \|\mathbf{x}_i\|_0 \leq T$$

$$\mathbf{y} = \begin{bmatrix} \text{patch} \end{bmatrix} = \begin{bmatrix} \text{patch} \end{bmatrix} x_1 + \begin{bmatrix} \text{patch} \end{bmatrix} x_2 + \dots$$

LST slowness image and sampling



x's are stations

Slowness map and sampling:

- Discrete slowness map $N=W_1 \times W_2$ pixels
- I overlapping $\sqrt{n} \times \sqrt{n}$ pixel patches
- M straight-ray paths

Tomography matrix
(straight ray)

$$\mathbf{A} \in \mathbb{R}^{M \times N}$$

Slowness dictionary

$$\mathbf{D} \in \mathbb{R}^{n \times Q}$$

$$Q \ll I$$

“Local” model

$$\hat{\mathbf{x}}_i = \arg \min_{\mathbf{x}_i} \|\mathbf{R}_i \mathbf{s}_s - \mathbf{D} \mathbf{x}_i\|_2^2 \text{ subject to } \|\mathbf{x}_i\|_0 = T$$

“Global” model

$$\mathbf{t} = \mathbf{A} \mathbf{s}_g + \epsilon, \quad \hat{\mathbf{s}}_g = \arg \min_{\mathbf{s}_g} \|\mathbf{t} - \mathbf{A} \mathbf{s}_g\|_2^2 + \lambda_1 \|\mathbf{s}_g - \mathbf{s}_s\|_2^2,$$

Formulation of LST and algorithm

Bayesian MAP objective:

$$\{\hat{\mathbf{s}}_g, \hat{\mathbf{s}}_s, \hat{\mathbf{X}}\} = \arg \min_{\mathbf{s}_g, \mathbf{s}_s, \mathbf{X}} \left\{ \frac{1}{\sigma_\epsilon^2} \|\mathbf{t} - \mathbf{A}\mathbf{s}_g\|_2^2 + \frac{1}{\sigma_g^2} \|\mathbf{s}_g - \mathbf{s}_s\|_2^2 + \frac{1}{\sigma_{p,i}^2} \sum_i \|\mathbf{D}\mathbf{x}_i - \mathbf{R}_i\mathbf{s}_s\|_2^2 \right\}$$

subject to $\|\mathbf{x}_i\|_0 = T \quad \forall i.$

Solution via block-coordinate descent

- **Global model:** the global slowness is solved as

$$\hat{\mathbf{s}}_g = \arg \min_{\mathbf{s}_g} \|\mathbf{t} - \mathbf{A}\mathbf{s}_g\|_2^2 + \lambda_1 \|\mathbf{s}_g - \mathbf{s}_s\|_2^2, \quad \lambda_1 = (\sigma_\epsilon/\sigma_g)^2$$

- **Local model:** sparse coding and dictionary learning, decoupled from MAP objective

$$\hat{\mathbf{x}}_i = \arg \min_{\mathbf{x}_i} \|\mathbf{D}\mathbf{x}_i - \mathbf{R}_i\hat{\mathbf{s}}_g\|_2^2 \quad \text{subject to } \|\mathbf{x}_i\|_0 = T. \quad (\mathbf{s}_s = \hat{\mathbf{s}}_g)$$

Dictionary learning by iterative thresholding and signed K-means (ITKM) algorithm, Schnass 2015

$$\max_{\mathbf{D}} \sum_i \max_{|K|=T} \|\mathbf{D}_K^T \mathbf{y}_i\|_1,$$

- The sparse slowness is then solved from

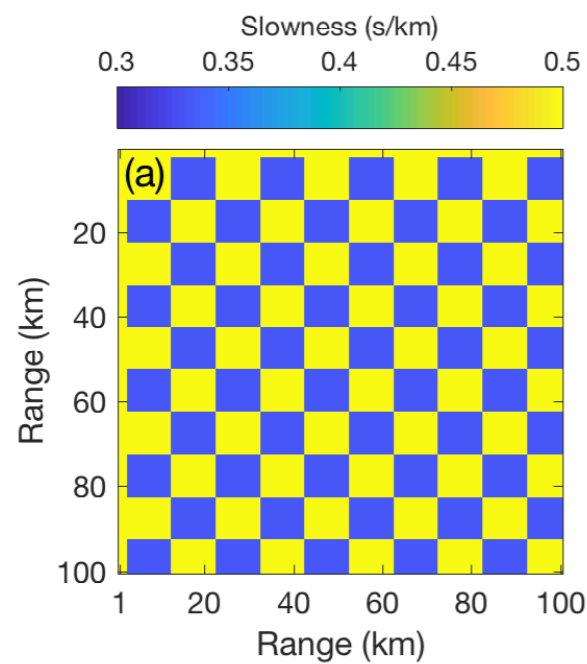
$$\hat{\mathbf{s}}_s = \arg \min_{\mathbf{s}_s} \lambda_2 \|\hat{\mathbf{s}}_g - \mathbf{s}_s\|_2^2 + \sum_i \|\mathbf{D}\hat{\mathbf{x}}_i - \mathbf{R}_i\mathbf{s}_s\|_2^2, \quad \lambda_2 = (\sigma_p/\sigma_g)^2$$

"Slowness at pixel n"

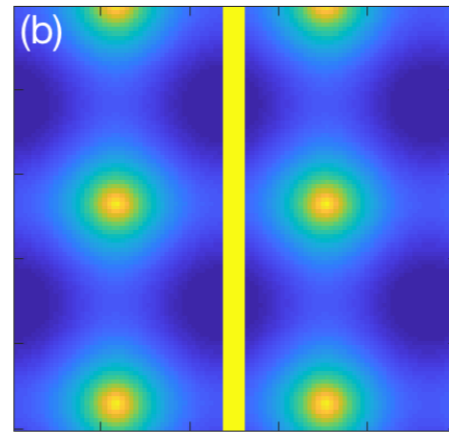
$$\hat{\mathbf{s}}_s = \left(\lambda_2 \mathbf{I} + \sum_i \mathbf{R}_i^T \mathbf{R}_i \right)^{-1} \left(\lambda_2 \hat{\mathbf{s}}_g + \sum_i \mathbf{R}_i^T \mathbf{D} \hat{\mathbf{x}}_i \right) \longrightarrow \hat{s}_{s,n} = \frac{\lambda_2 \hat{s}_{g,n} + b_n s_{p,n}}{\lambda_2 + b_n}$$

$$\mathbf{b} = \text{diag} \left(\sum_i \mathbf{R}_i^T \mathbf{R}_i \right) \in \mathbb{Z}^N \quad s_{p,n} = p_n/b_n$$

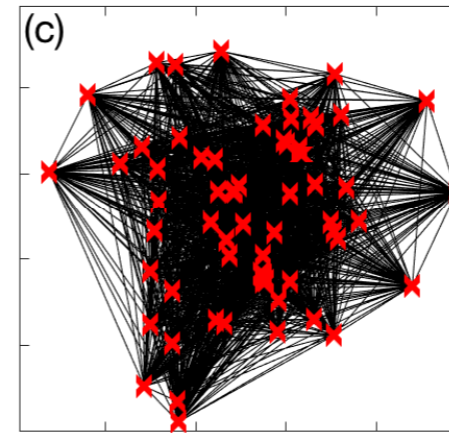
Synthetic slownesses and dictionaries



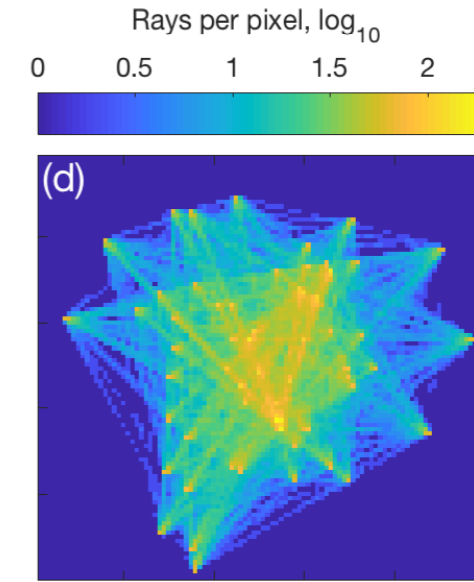
Checkerboard



"Fault" profile

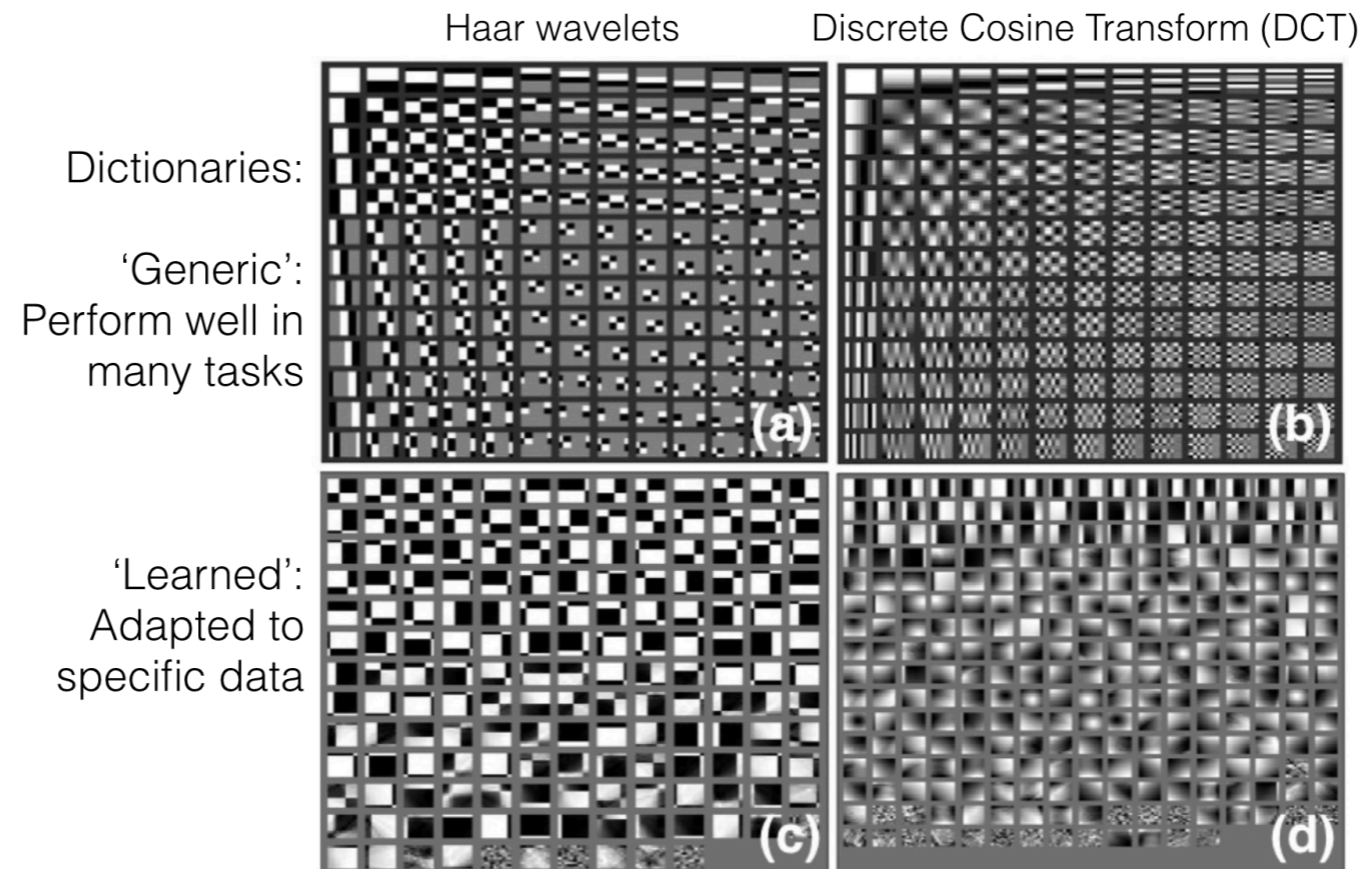


Ray sampling (64 stations, 2016 rays)



Ray density

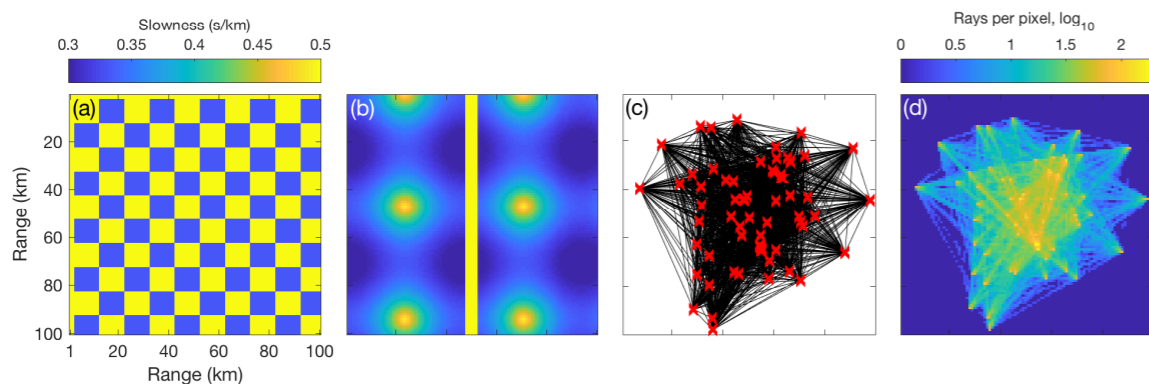
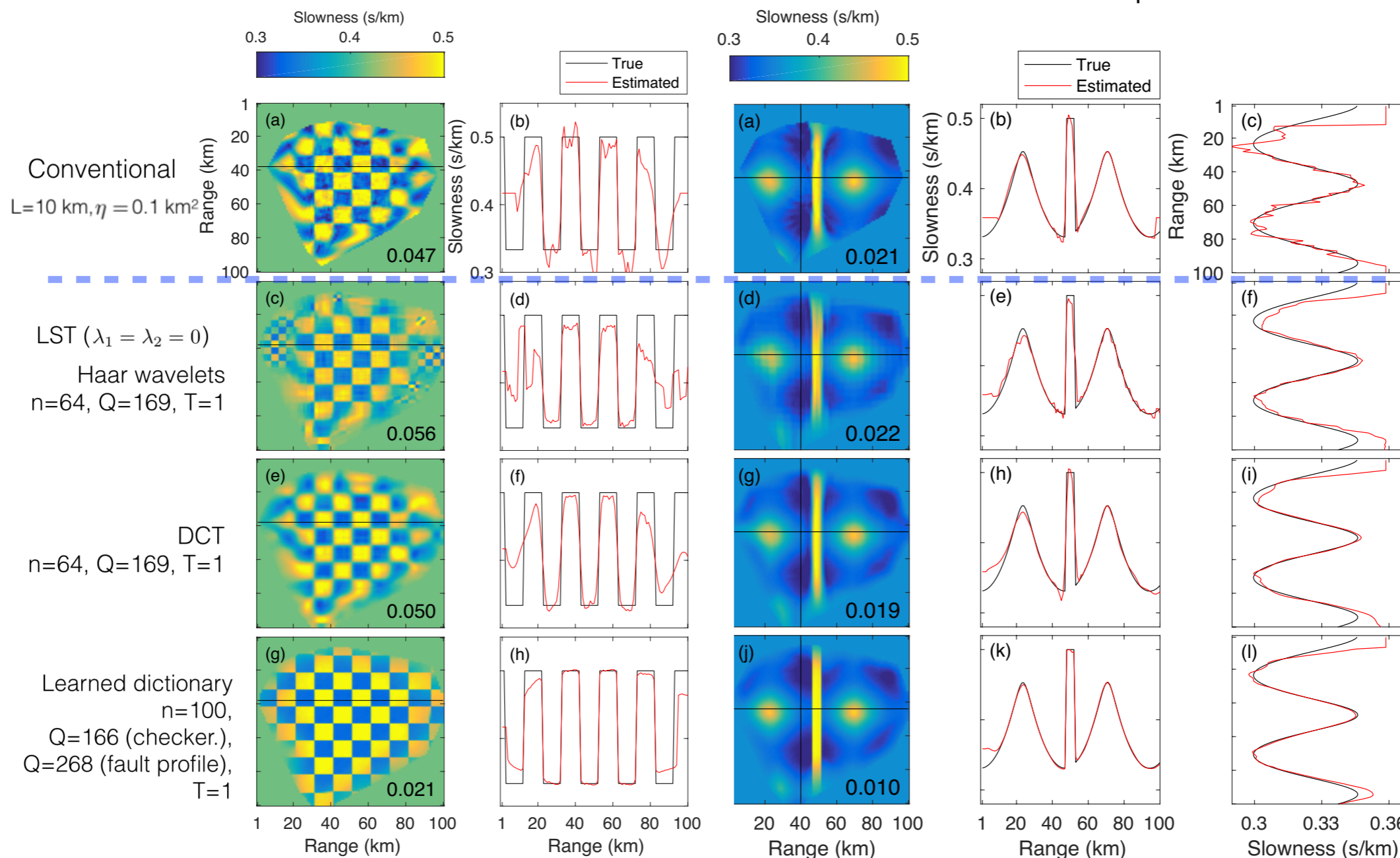
Dictionaries: Prescribed and Learned



Checkerboard (Q=166, T=1) Fault profile (Q=268, T=1)

LST vs. conventional method: synthetic inversions without noise

Each example took ~5 min on MacBook Pro



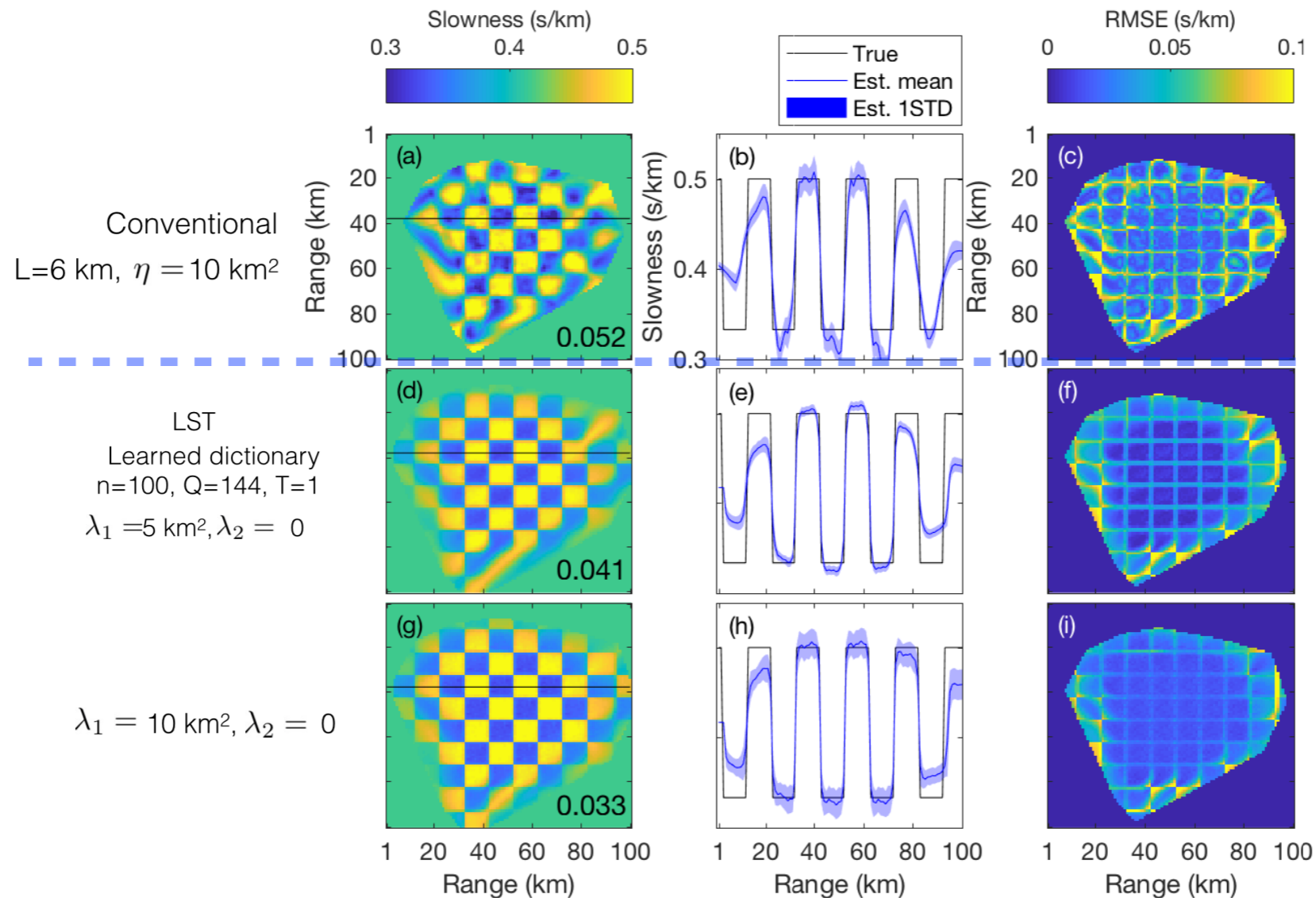
Conventional tomography method
 (Rodgers 2000)

$$\hat{\mathbf{s}}_g = (\mathbf{A}^T \mathbf{A} + \eta \Sigma_L^{-1})^{-1} \mathbf{A}^T \mathbf{t}, \quad \Sigma_L(i, j) = \exp(-D_{i,j}/L)$$

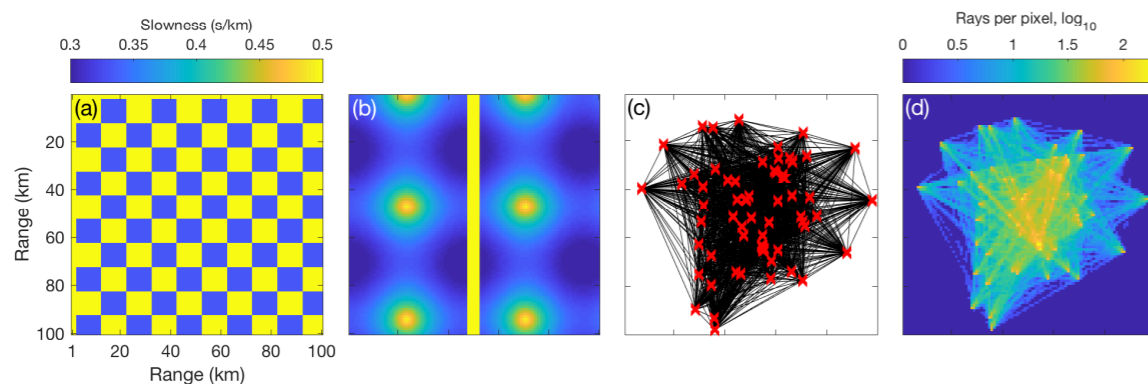
with $\eta = (\sigma_\epsilon / \sigma_g)^2$

LST vs. conventional method: synthetic inversions with travel time noise

Checkerboard



- Slowness RMSE (s/km) written on 2D estimates
- Noise is Gaussian with STD 2% of mean travel time



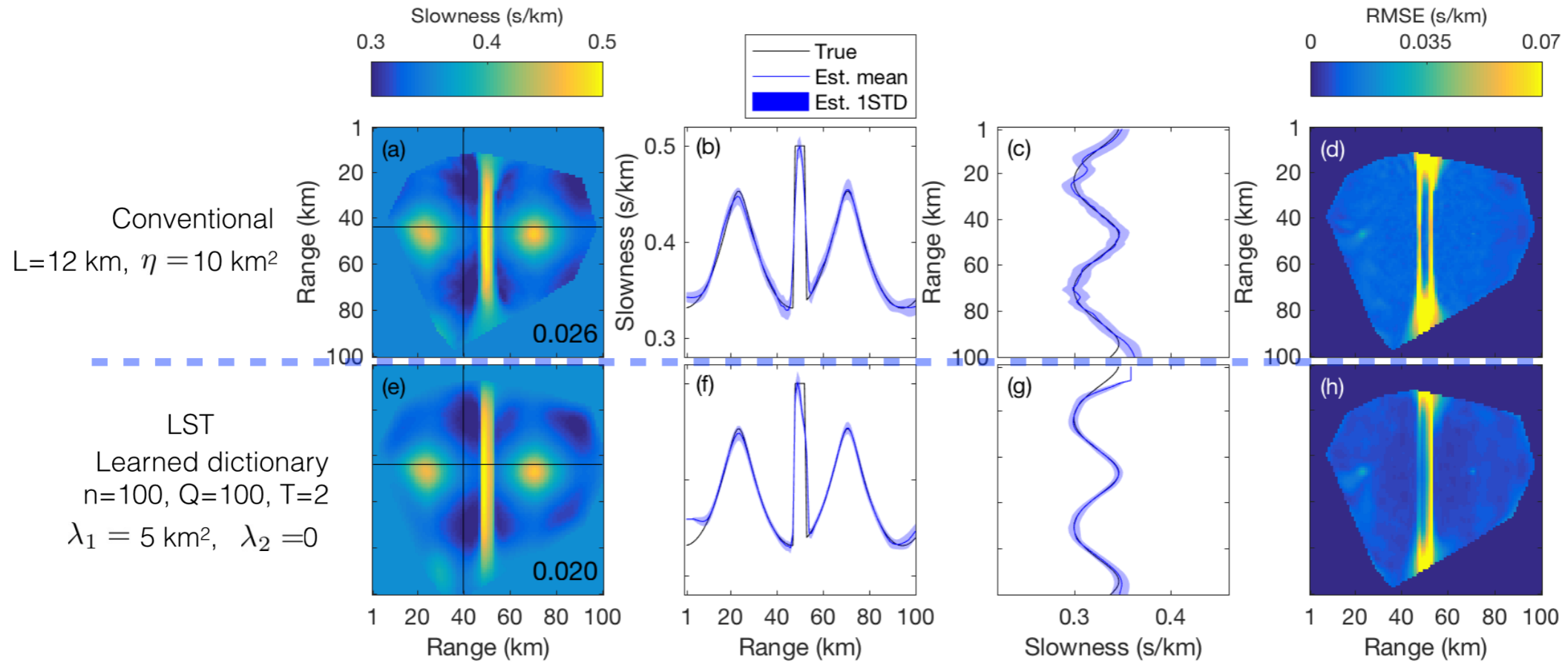
Conventional tomography method
 (Rodgers 2000)

$$\hat{\mathbf{s}}_g = (\mathbf{A}^T \mathbf{A} + \eta \mathbf{\Sigma}_L^{-1})^{-1} \mathbf{A}^T \mathbf{t}, \quad \mathbf{\Sigma}_L(i, j) = \exp(-D_{i, j} / L)$$

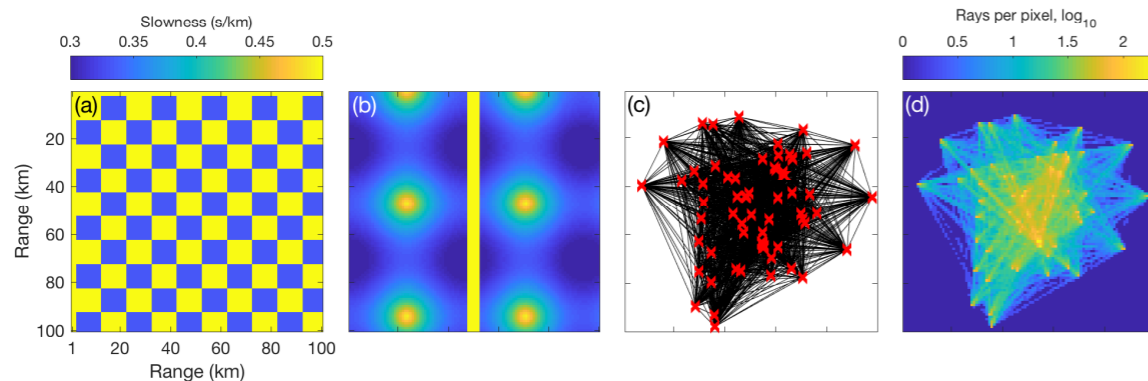
with $\eta = (\sigma_\epsilon / \sigma_g)^2$

LST vs. conventional method: synthetic inversions with travel time noise

Fault profile



- Slowness RMSE (s/km) written on 2D estimates
- Noise is Gaussian with STD 2% of mean travel time



Conventional tomography method
 (Rodgers 2000)

$$\hat{\mathbf{s}}_g = (\mathbf{A}^T \mathbf{A} + \eta \boldsymbol{\Sigma}_L^{-1})^{-1} \mathbf{A}^T \mathbf{t}, \quad \boldsymbol{\Sigma}_L(i, j) = \exp(-D_{i, j} / L)$$

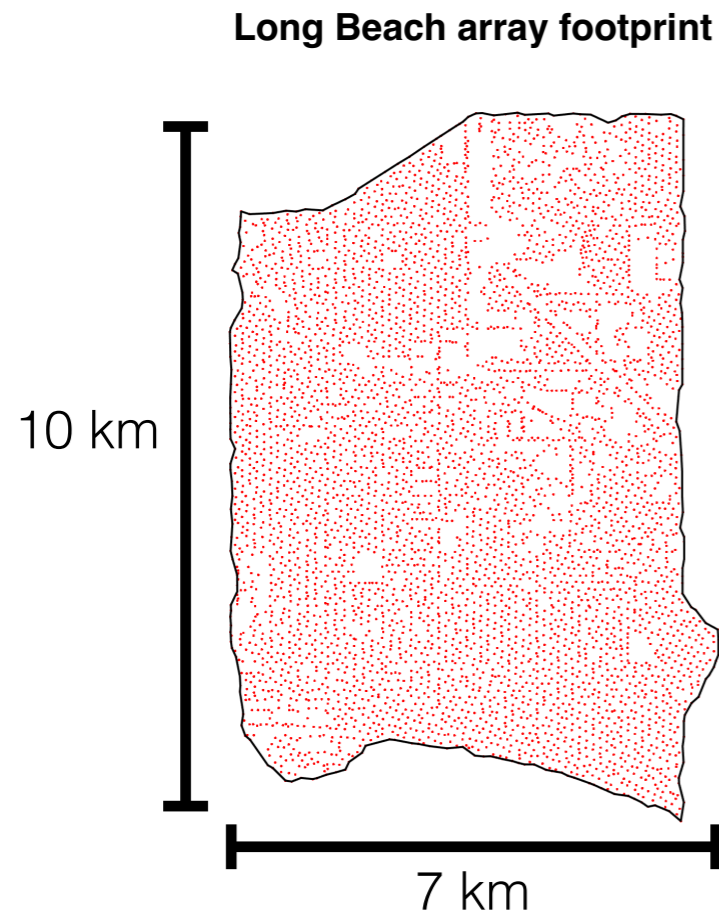
with $\eta = (\sigma_\epsilon / \sigma_g)^2$

Imaging Long Beach, CA using LST: Big Data task

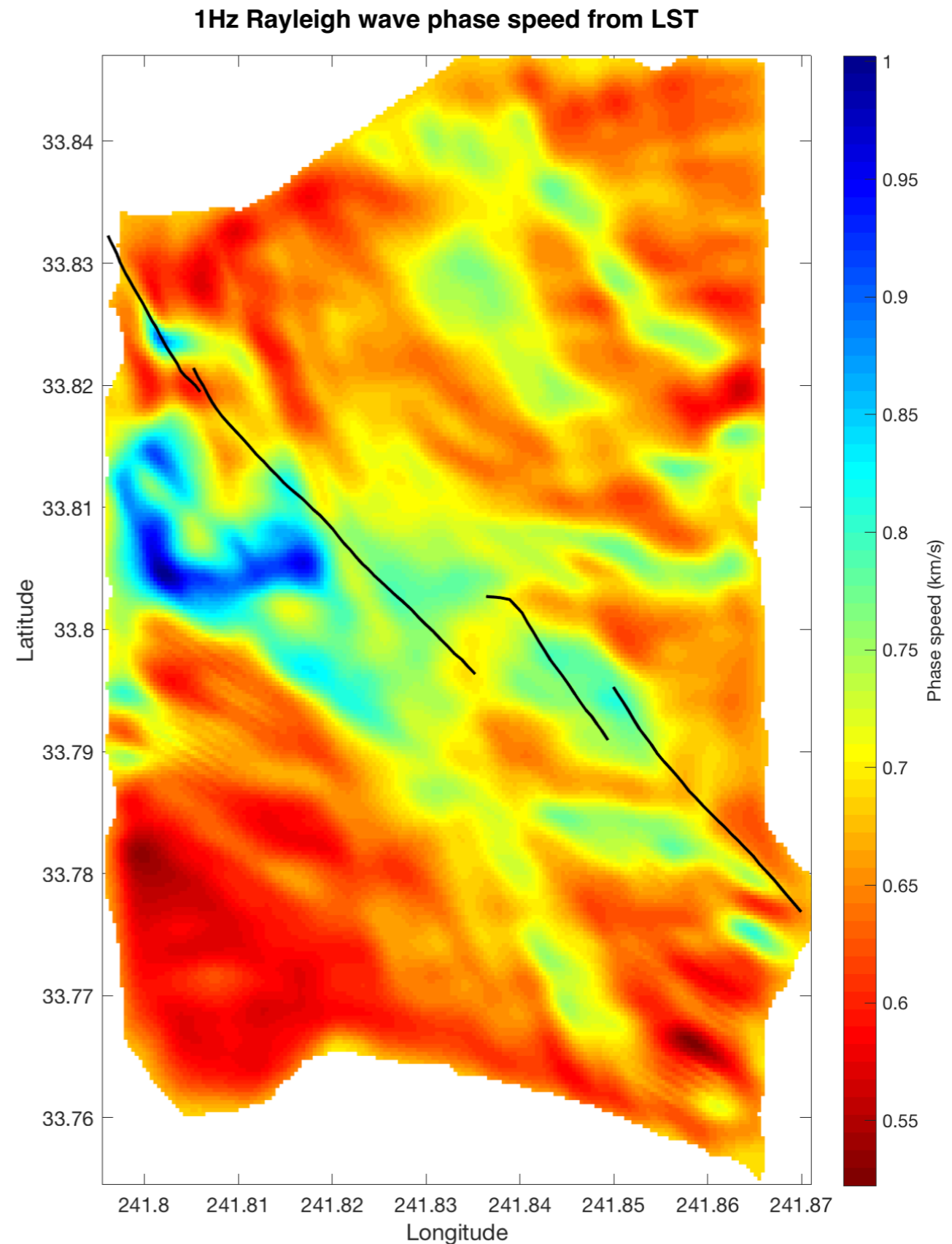


- In March 2011, 5200 seismic stations were deployed in Long Beach, California over 70km² area
- Ambient seismic noise cross-correlations were obtained for all unique virtual source-receiver pairs (~14 million ray paths) using 3 weeks of data
- We consider only the 1Hz vertical component data, corresponding to Rayleigh surface waves (from Lin et al. 2013)
- After quality control there were ~8 million ray paths

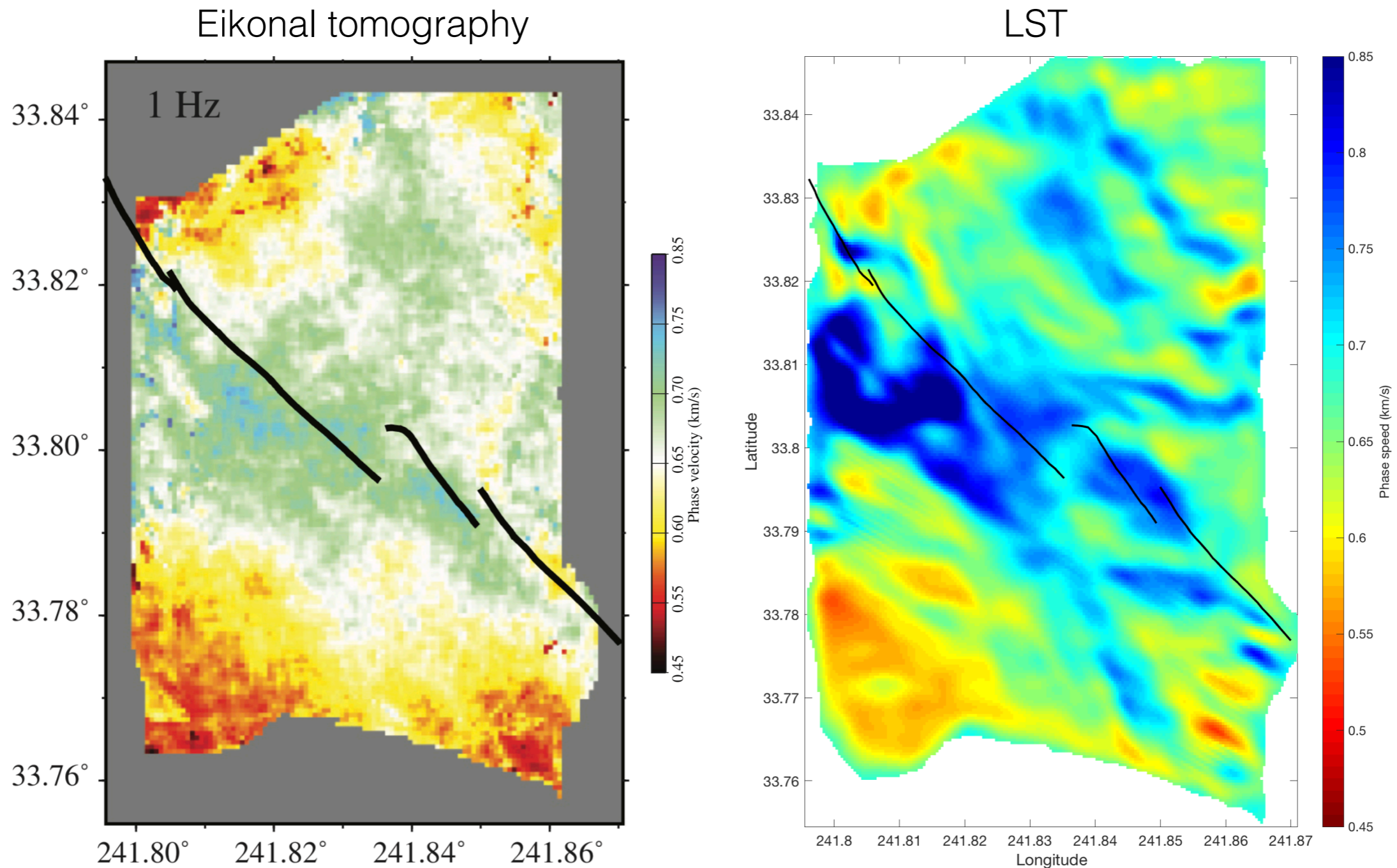
High-resolution LST phase speed map from 8 million cross-correlations



- For LST we generate a 300x200 pixel slowness map with 8 million rays (tomography matrix \mathbf{A} has dimensions $M=8$ million, $N=60000$)
- 10 iterations, used ~2 cpu-hours
- Since we are not imposing global correlations on pixels, can treat \mathbf{A} as sparse matrix, get fast inversion for global model (which is bottleneck)
- Newport-Inglewood fault network shown as black line



LST comparison with eikonal tomography (Lin et al. 2013)



- We observe the same general trends between eikonal and LST
- From LST we have improved contrast along fault lines, for example near Signal Hill
- The LST results are preliminary and they can likely be improved with more careful preprocessing (future work)