

Tutorial-10: Dictionary learning, Deep Image Prior

Sanketh Vedula
sanketh@cs.technion.ac.il

Agenda

- Dictionary learning: K-SVD
 - Problem setup
 - Sparse coding
 - Dictionary update
 - Implementation
- Deep Image Prior
 - Standard image priors
 - Parametrised image priors
 - Deep Image Prior: training strategy
 - Results

Sparseland signal synthesis

$$\mathbf{x} \in \mathbb{R}^N$$

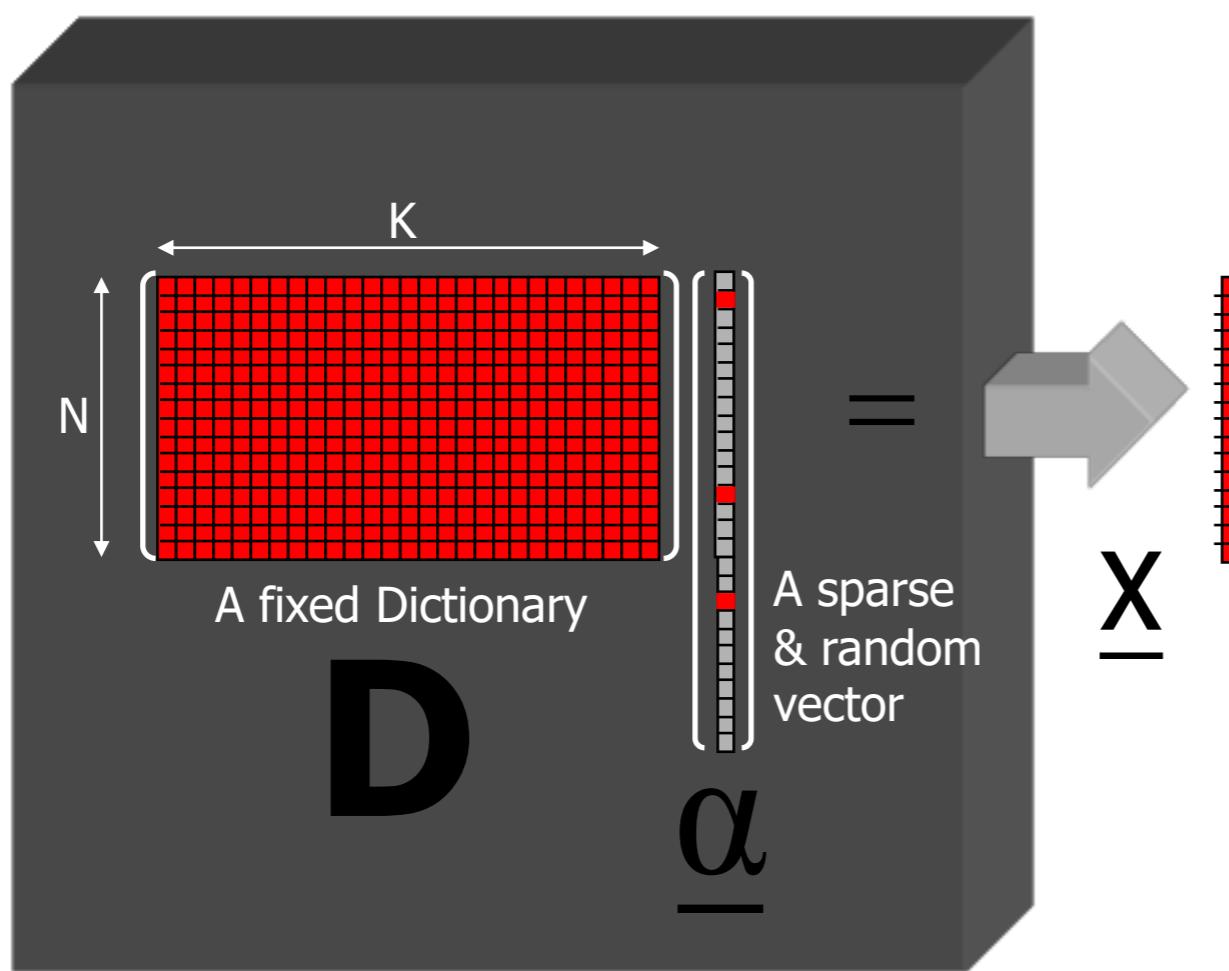
signal

$$\mathbf{D} \in \mathbb{R}^{N \times K}$$

dictionary

$$\alpha \in \mathbb{R}^K$$

representation



- Every column in \mathbf{D} (dictionary) is a prototype signal (Atom).

- The vector $\underline{\alpha}$ is generated randomly with few non-zeros in random locations and random values.

Sparse coding

$$\mathbf{x} \in \mathbb{R}^N \quad \mathbf{D} \in \mathbb{R}^{N \times K} \quad \boldsymbol{\alpha} \in \mathbb{R}^K$$

signal dictionary representation

$$(P_0^\epsilon) : \min \|\boldsymbol{\alpha}\|_0, \text{ such that } \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|^2 < \epsilon$$

$$(P_0^\lambda) : \arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|^2 + \lambda \|\boldsymbol{\alpha}\|_0$$

NP-hard!

Approximate solution: Orthogonal Matching Pursuit



Greedily pick atoms that minimise the residual error

Sparse coding

$$(P_0^\lambda) : \arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{D}\alpha\|^2 + \lambda \|\alpha\|_0$$

Orthogonal Matching Pursuit:

init: support = {}, residual (r) = x

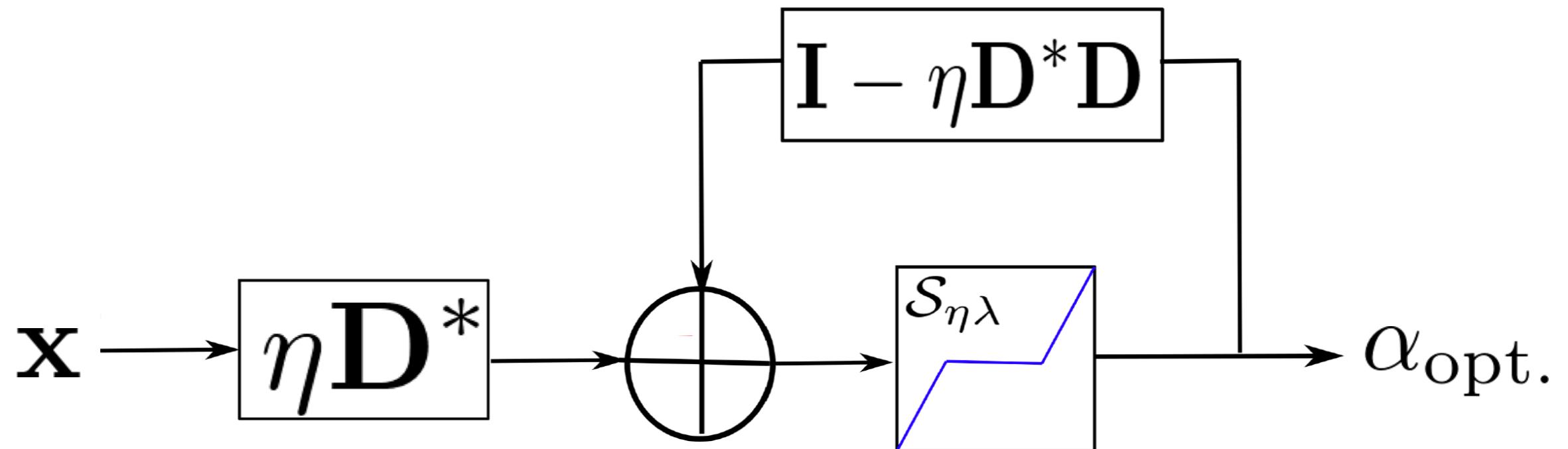
Until reaching the desired level of sparsity:

- Choose an unselected atom/column from D that minimises $\| r - D \alpha \|$
- Add the chosen atom (d_i) into the support => support = support + $\{d_i\}$
- Update the residual (r) => $r = \| x - D \alpha \{d_i\} \|$
- Repeat until reaching desired level of sparsity or when r is close to 0

Approximate Sparse coding

$$(P_1^\lambda) : \arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{D}\alpha\|^2 + \lambda \|\alpha\|_1$$

Solve this? ISTA, Basis Pursuit



Where do dictionaries come from?

Carefully incorporate the information we know about the data within these dictionaries:

Modeling?

For example:

If we are working with natural images, then D could be over-complete DCT basis, Wavelet basis, etc.

Pros:

Fast transforms

Proven optimality

Cons:

How do you adapt to different signals?

Not all dictionaries would allow sparsity?

Maybe optimal for large signal families – but are not “signal-specific”.

Where do dictionaries come from?

Learning?

Train the dictionary directly on the given examples, optimizing w.r.to. sparsity and other desired properties (normalized atoms etc.)

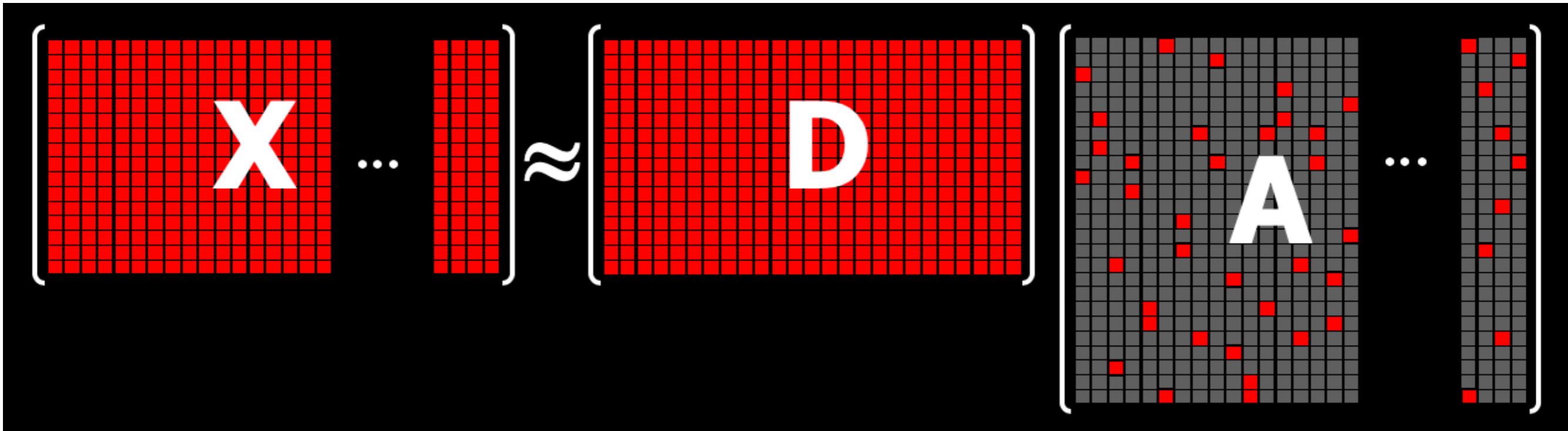
Pros:

Fits the data at hand
Signal-specific or “signal family”
specific modeling

Cons:

No fast transforms anymore
Training takes time
Not scalable?

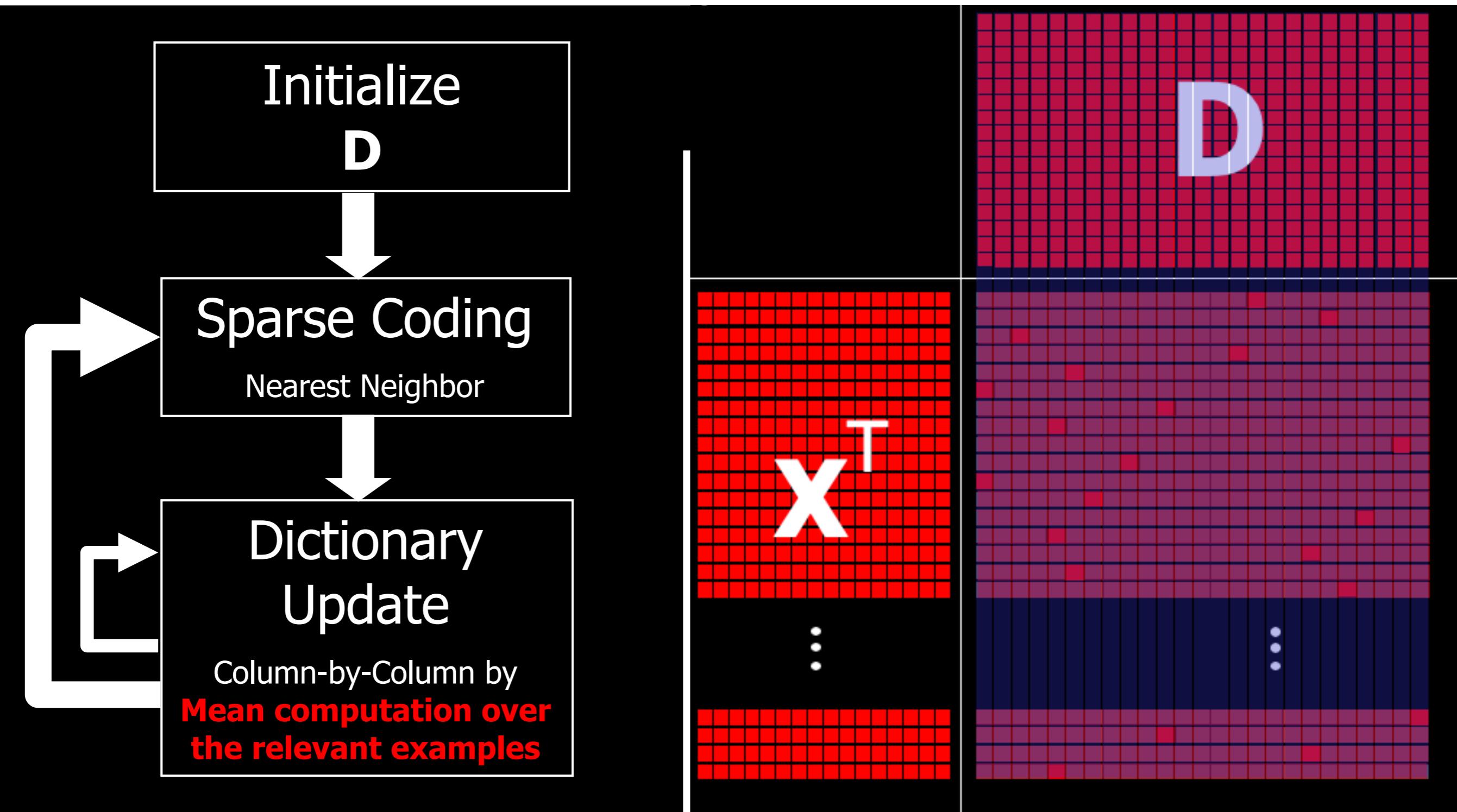
Dictionary Learning



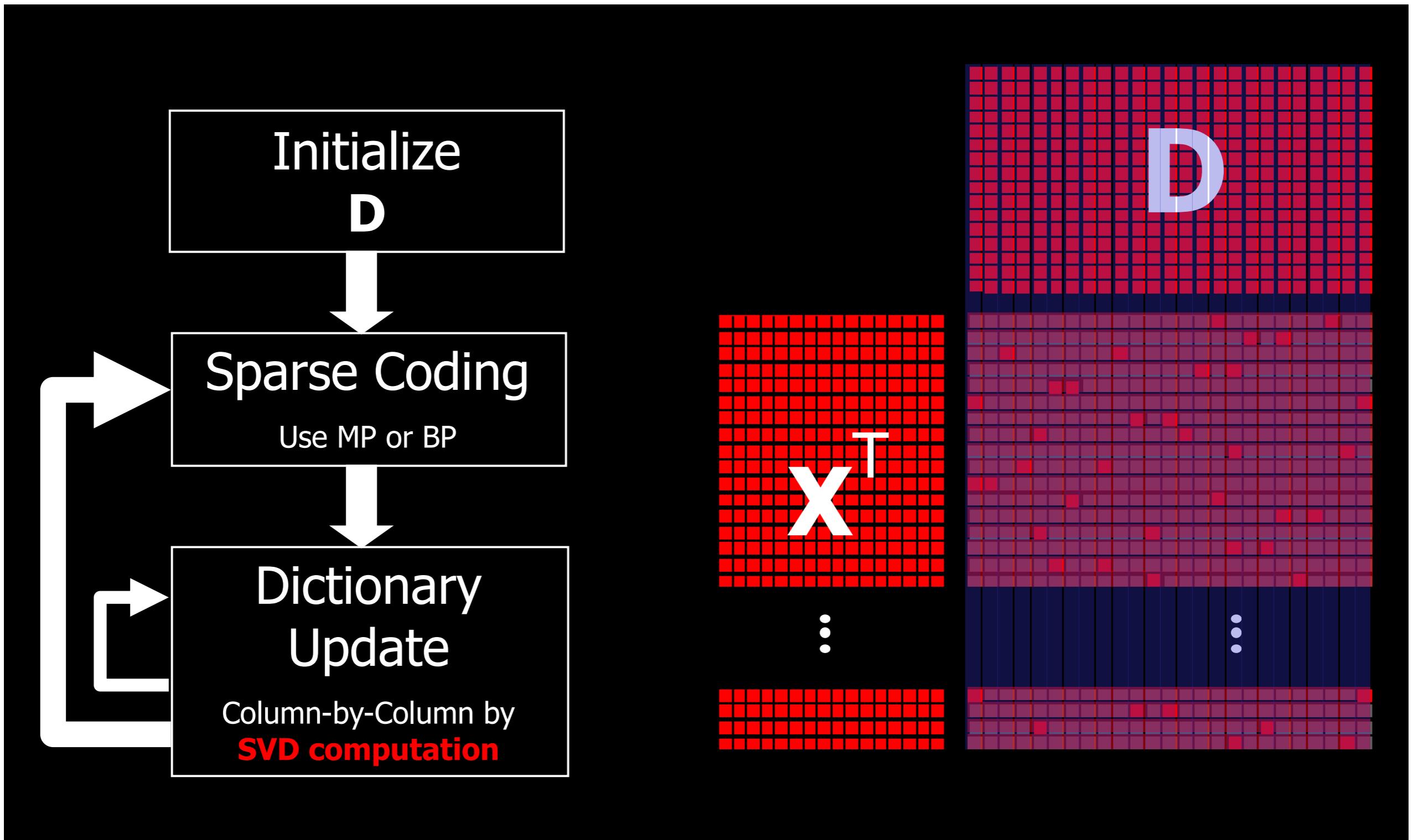
$$\underset{\mathbf{D}, \mathbf{A}}{\text{Min}} \sum_{j=1}^P \left\| \mathbf{D} \underline{\alpha}_j - \underline{\mathbf{x}}_j \right\|_2^2 \quad \text{s.t. } \forall j, \left\| \underline{\alpha}_j \right\|_0 \leq L$$

Objective: Update dictionary and representations such that the representation error is minimal while still maintaining sparsity.

K-Means



K-SVD

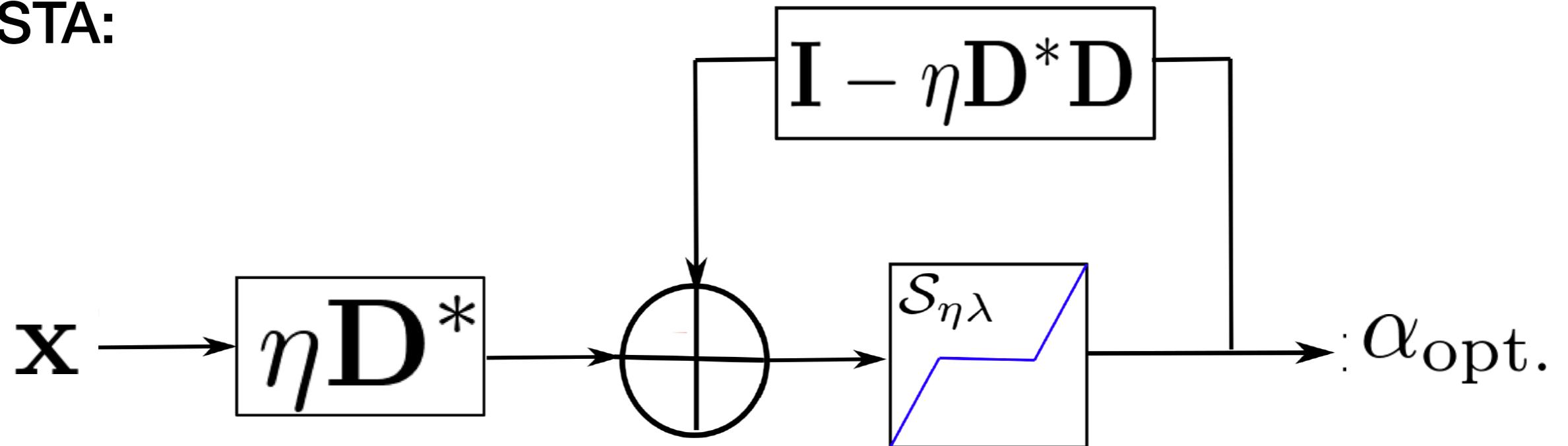


K-SVD: Sparse coding

Goal: We know D , we need to get the sparse code.

Pursuit methods: OMP, Basis Pursuit

ISTA:



K-SVD: dictionary update

$$\|\mathbf{X} - \mathbf{DA}\|_F^2 = \left\| \mathbf{X} - \sum_{j=1}^K d_j \alpha_j \right\|_F^2 = \left\| \left(\mathbf{X} - \sum_{j \neq k} d_j \alpha_j \right) - d_k \alpha_k \right\|_F^2 = \|\mathbf{E}_k - d_k \alpha_k\|_F^2$$

↑
SVD?

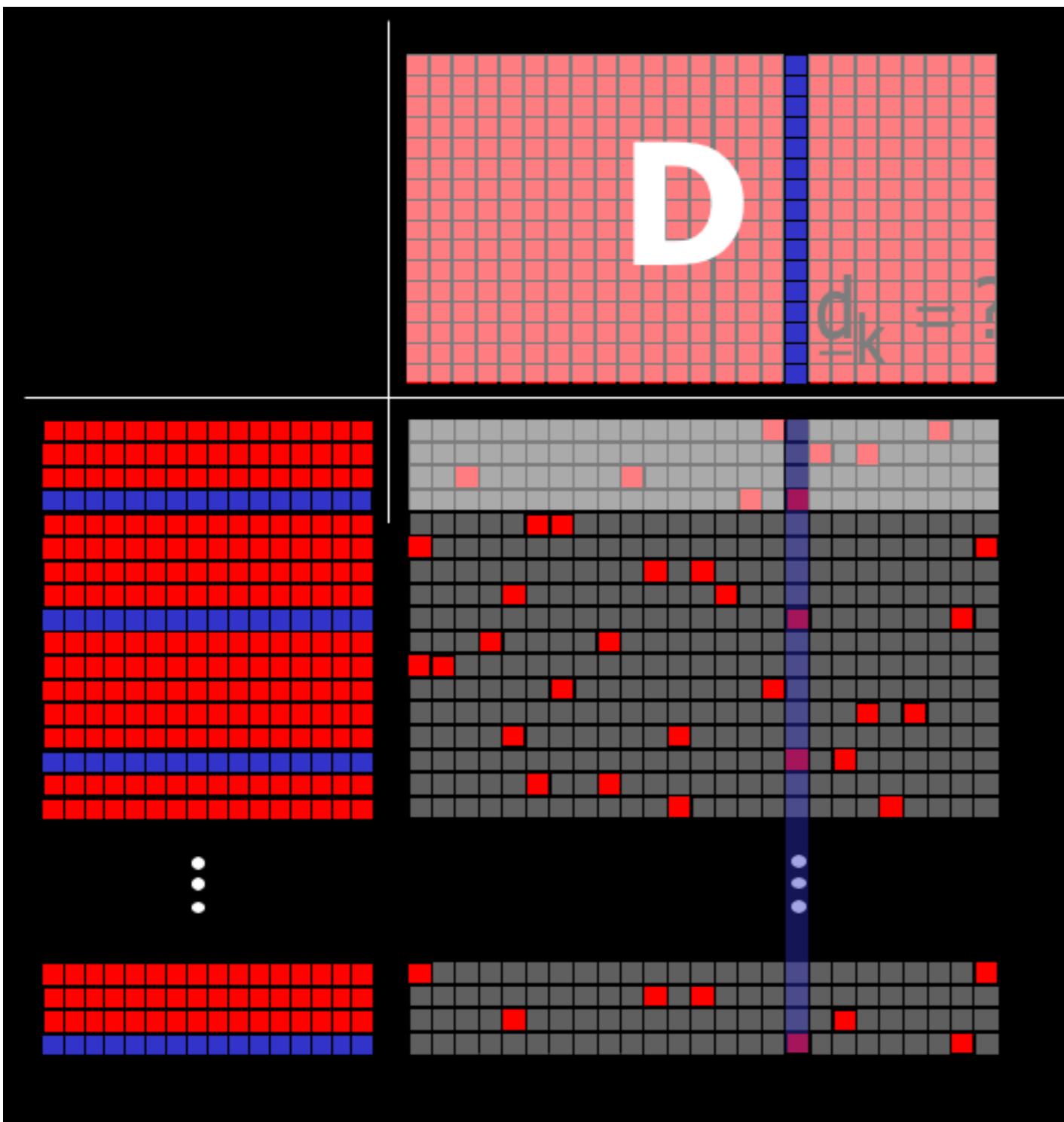
We want to update kth column in D.
How to do it?

Possible solution: For each column in D,
let us treat the update step as a matrix factorisation
problem – just do singular value decomposition
and update the D and A.

Caveat?

This might be destroying the sparsity constraint.
Because, there's good chance that the resulting sparse
code is not sparse.

K-SVD: Dictionary update



Do the following steps for each dictionary column d_k :



Find the signals that use d_k



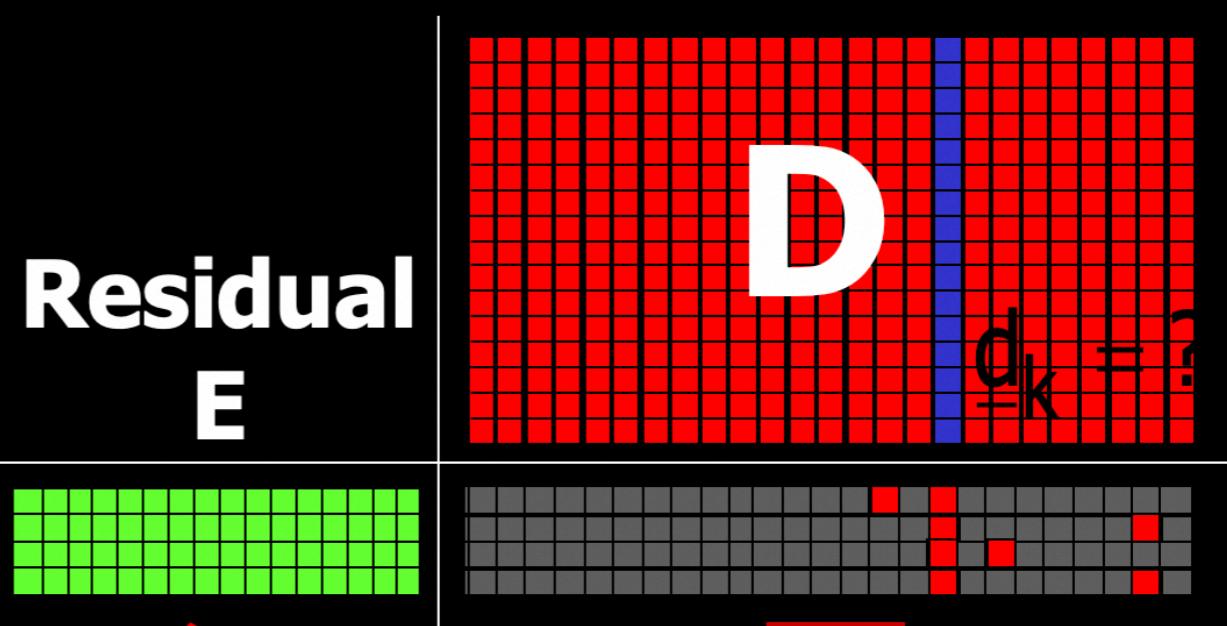
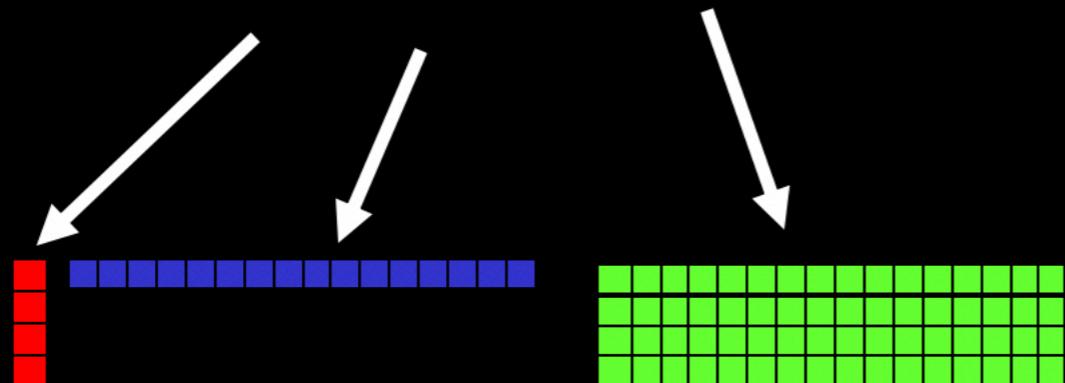
Fix the rest of the A and D apart from the kth column and seek both d_k and kth column in A to better fit the residual

$$\|\Omega \alpha_k d_k - \Omega E_k\|_F^2$$

K-SVD: Dictionary update

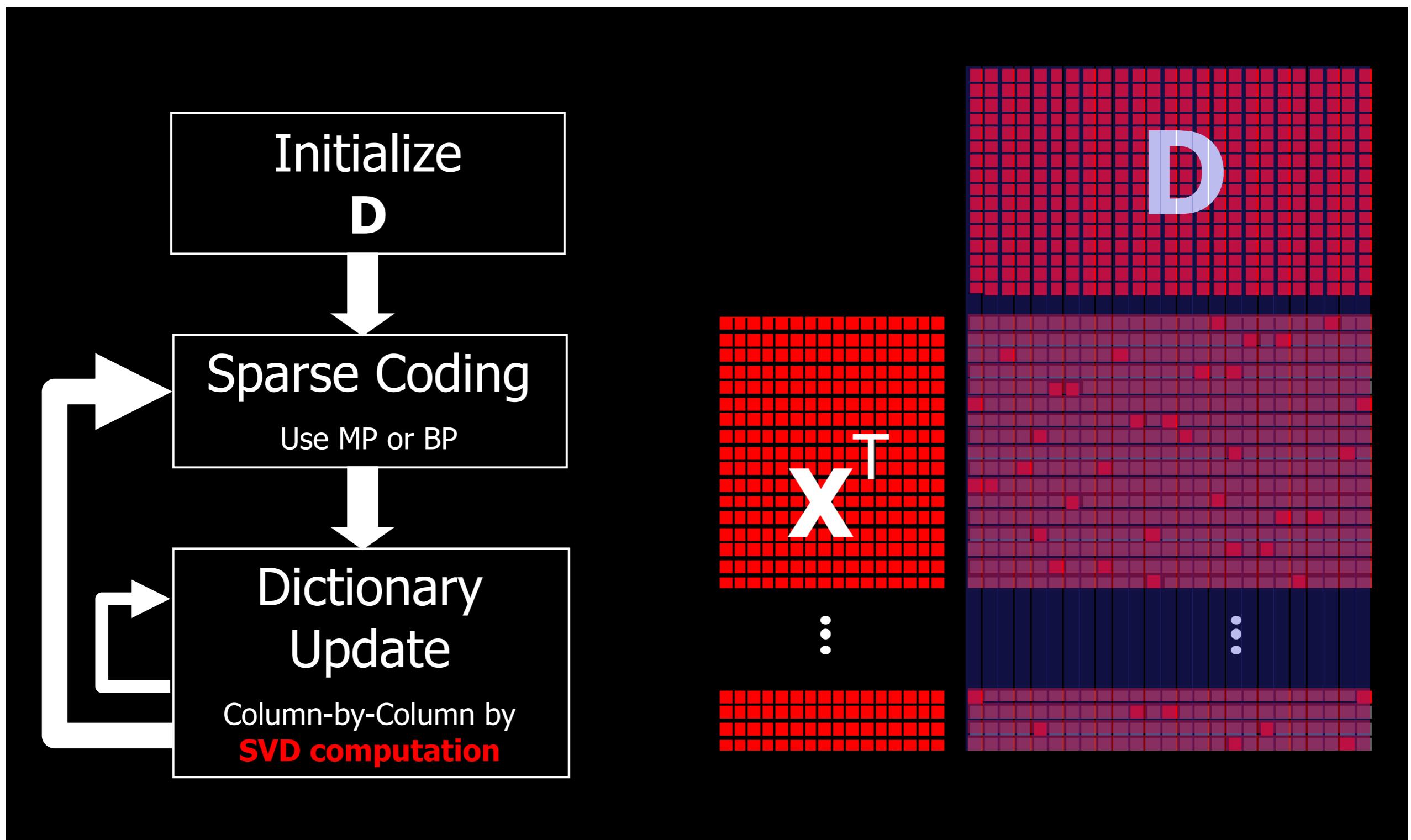
Don't update the zeros
Just pick the non-zeros
and do the matrix factorization.
Less complexity + sparsity preserved!

$$\|\Omega \alpha_k d_k - \Omega E_k\|_F^2$$



Omega represents the set
of rows corresponding
to the non-zero elements

K-SVD - summarised



Unrolled networks

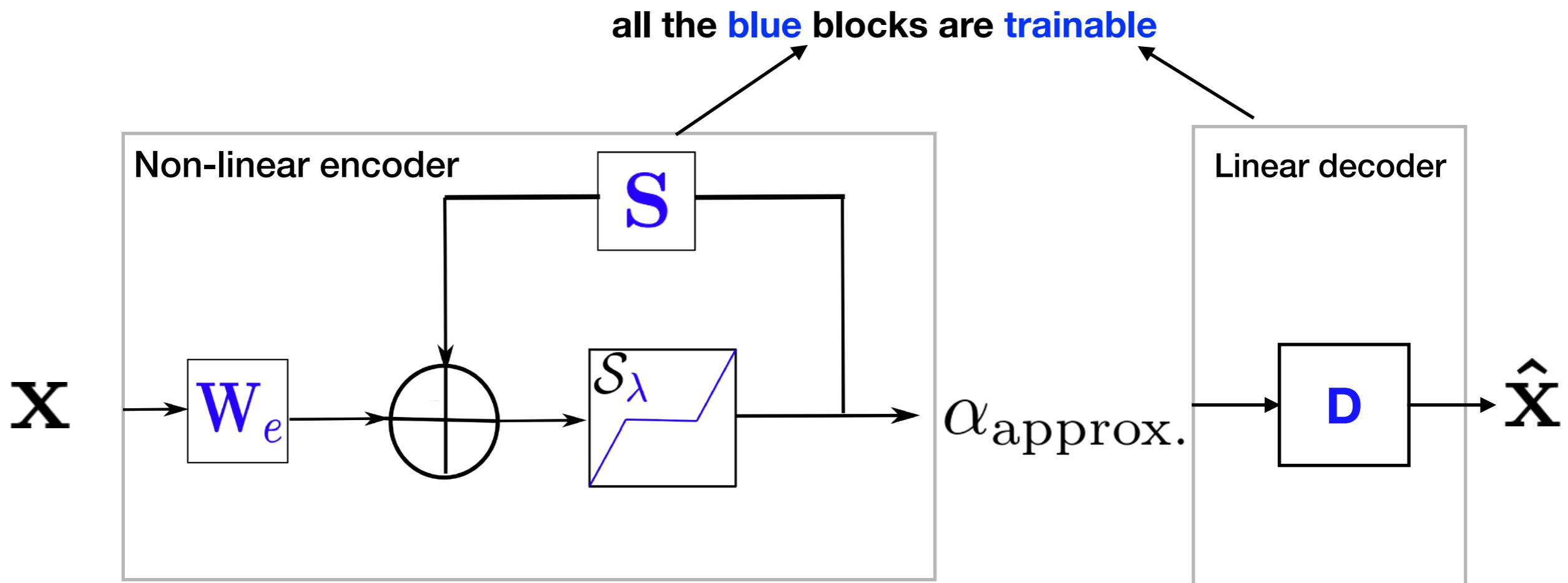
Denoising:

x - noisy input

\hat{x} - denoised

$$\mathcal{L}(x, \hat{x}) = \|\hat{x} - \hat{x}\|^2$$

optimized with gradient descent
and backpropagation



Notice the resemblance with neural networks – if the proximal operator
is non-negative and if it has a non-linear decoder

Image Priors

denoising,
inpainting, blur etc.



MAP: $x^* = \arg \max_{\textcolor{teal}{x}} p(\textcolor{teal}{x}|\hat{x})$

$$p(\textcolor{teal}{x}|\hat{x}) = \frac{p(\hat{x}|\textcolor{teal}{x})p(\textcolor{teal}{x})}{p(\hat{x})} \propto p(\hat{x}|\textcolor{teal}{x})p(\textcolor{teal}{x})$$

Likelihood Prior

Image Priors

- \mathbf{x} - Clean image
- $\hat{\mathbf{x}}$ - Corrupted image (observed)
- \mathbf{x}^* - Restored image

$$\begin{aligned}\mathbf{x}^* &= \arg \max_{\mathbf{x}} p(\hat{\mathbf{x}}|\mathbf{x})p(\mathbf{x}) \\ &= \arg \min_{\mathbf{x}} -\log p(\hat{\mathbf{x}}|\mathbf{x}) - \log p(\mathbf{x}) \\ &= \arg \min_{\mathbf{x}} E(\mathbf{x}; \hat{\mathbf{x}}) + R(\mathbf{x})\end{aligned}$$

We've seen several instances of how to model priors.

Parametrised Image Priors

- \mathbf{x} - Clean image
- $\hat{\mathbf{x}}$ - Corrupted image (observed)
- \mathbf{x}^* - Restored image

Regular:

$$\arg \min_{\mathbf{x}} E(\mathbf{x}; \hat{\mathbf{x}}) + R(\mathbf{x})$$

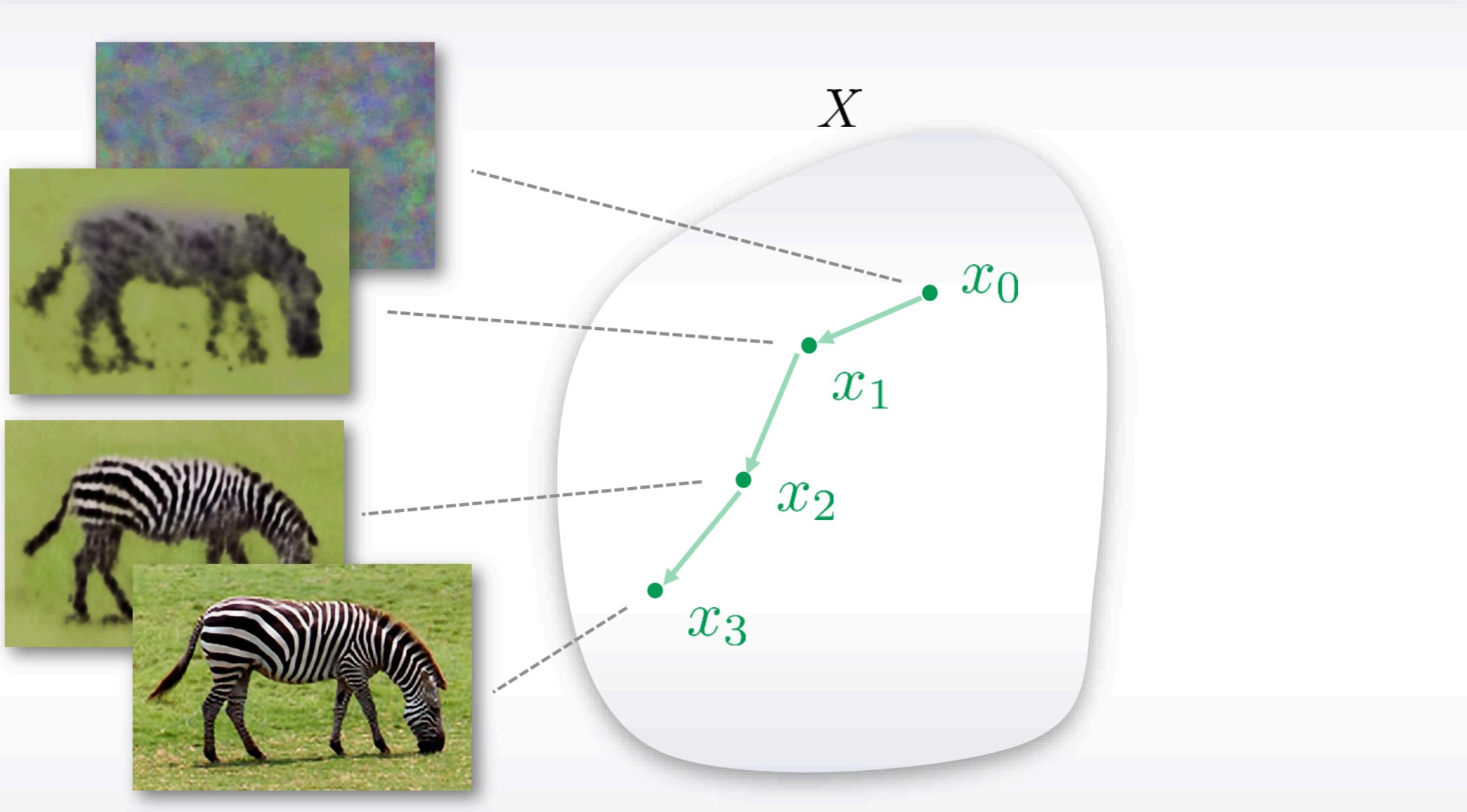
Search in the image space

Parametrized:

$$\arg \min_{\theta} E(g(\theta); \hat{\mathbf{x}}) + R(g(\theta))$$

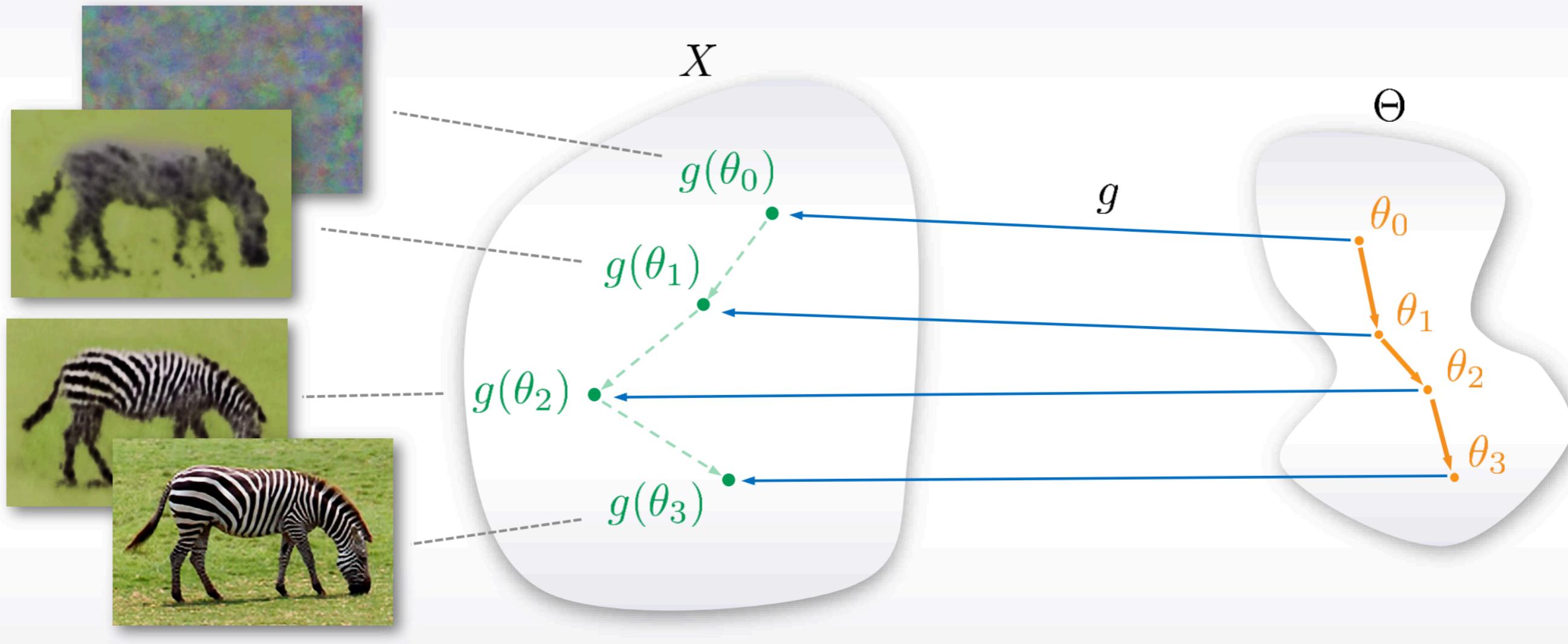
Search in some other space

Regular Image Priors



Regular: $\arg \min_{\mathbf{x}} E(\mathbf{x}; \hat{\mathbf{x}}) + R(\mathbf{x})$

Parametrised Image Priors



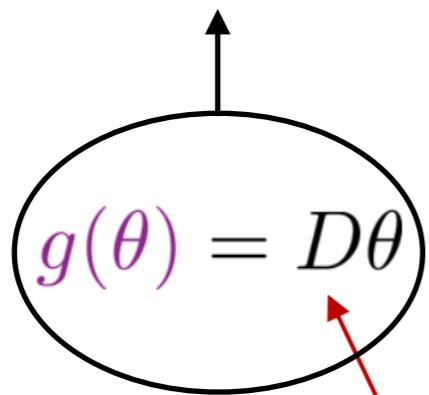
Parametrized: $\arg \min_{\theta} E(g(\theta); \hat{x}) + R(g(\theta))$

Parametrised Image Priors

Example:

Parametrized: $\arg \min_{\theta} E(g(\theta); \hat{x}) + R(g(\theta))$

Theta is sparse
D is a dictionary

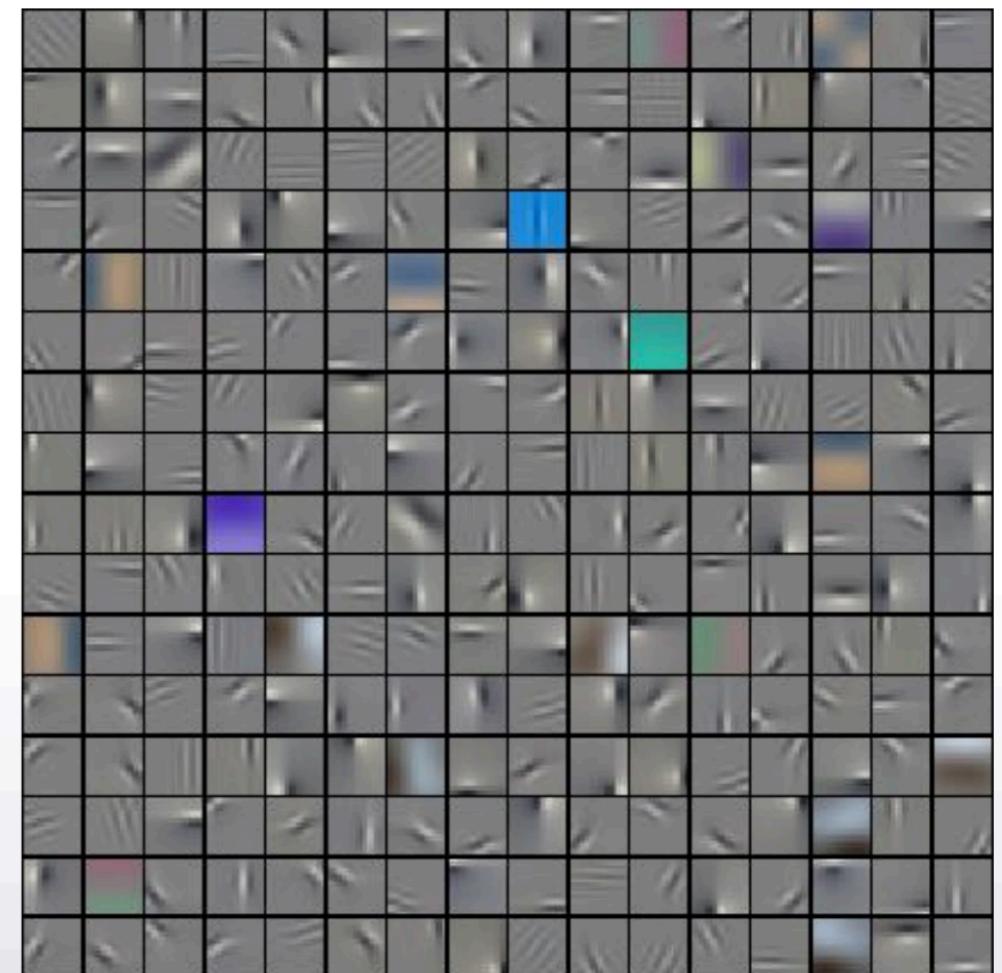


But this model is linear

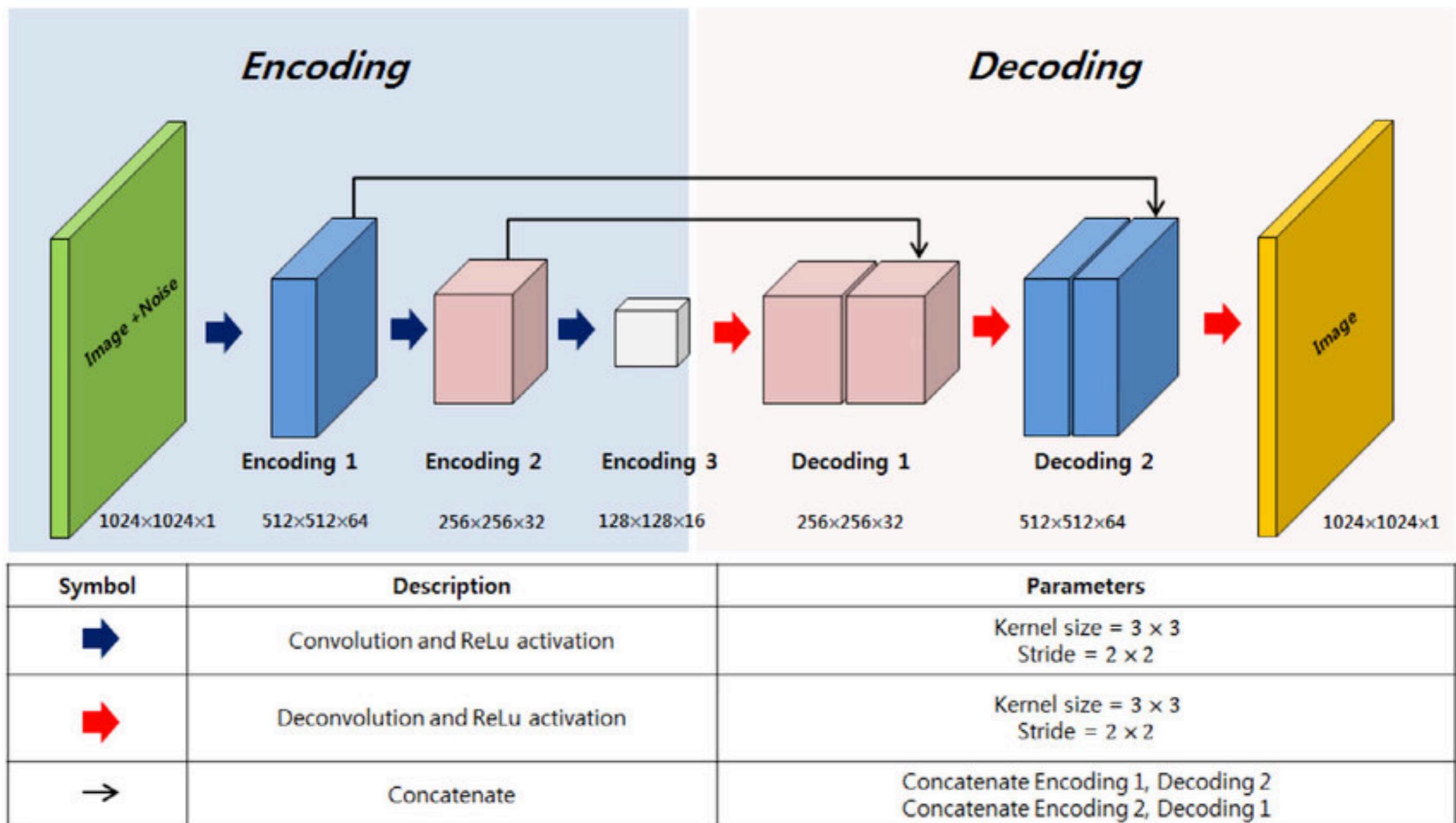
Limited expressivity

Coefficients

Dictionary



Convolutional Neural Networks



Deep Image Prior

We know CNNs are excellent generative models!

Let's just model $g(\Theta)$ using a CNN and "learn" the prior instead of putting it by hand.

Parametrized: $\arg \min_{\theta} E(g(\theta); \hat{x}) + R(g(\theta))$

$$g(\theta) \equiv f_{\theta}(z)$$

Fixed input

Convolutional network with parameters θ

g itself is a prior! And maybe it is sufficient to optimize only data term

Optimize only data term: $\arg \min_{\theta} E(f_{\theta}(z); \hat{x})$

Deep Image Prior: training

- \hat{x} – Corrupted image (observed)

1. Initialize z

- For example fill it with uniform noise $U(-1, 1)$

2. Solve

$$\arg \min_{\theta} E(f_{\theta}(z); \hat{x}))$$

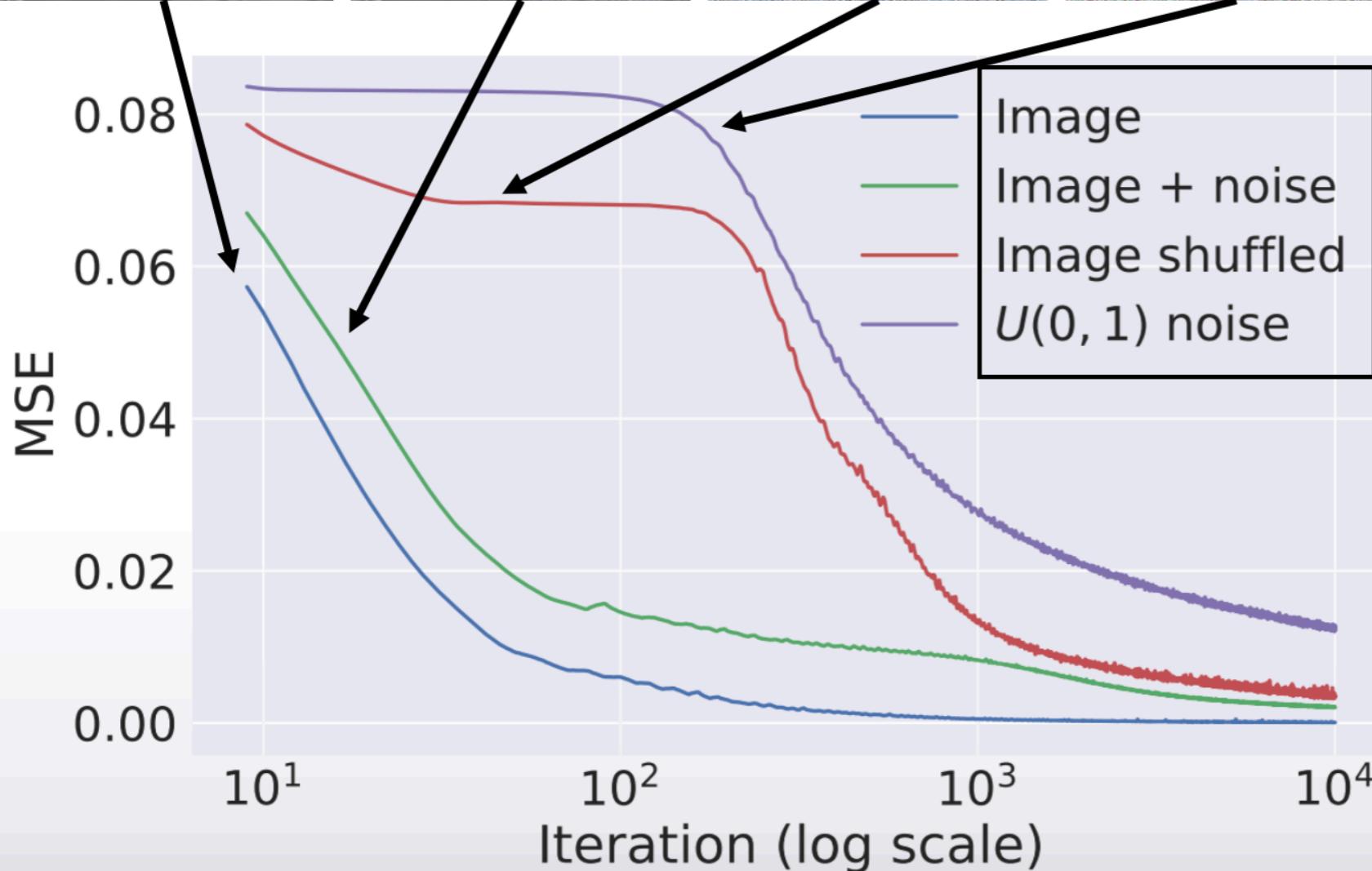
- With your favorite gradient-based method

$$\theta^{k+1} = \theta^k - \alpha \frac{\partial E(f_{\theta}(z); \hat{x})}{\partial \theta}$$

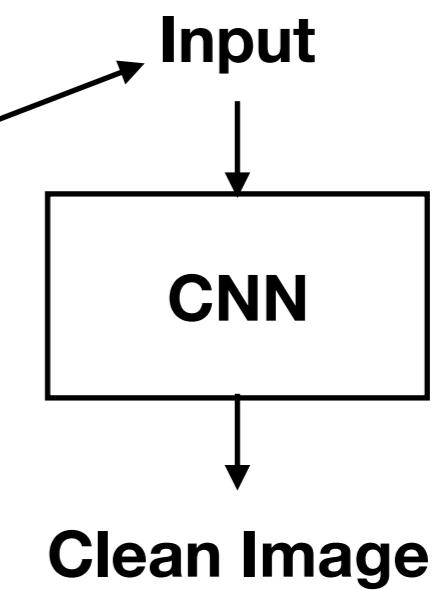
3. Get the solution

$$x^* = f_{\theta^*}(z)$$

Why should it work?



Networks have very strong expressive capabilities



Experiments

- \mathbf{x} - Clean image
- $\hat{\mathbf{x}}$ - Corrupted image (observed)
- m - Binary mask

Objective: $\arg \min_{\theta} E(f_{\theta}(z); \hat{\mathbf{x}}))$

- **Denoising:** $E(\mathbf{x}; \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|^2$ **Needs early stopping!**
- **Inpainting:** $E(\mathbf{x}; \hat{\mathbf{x}}) = \|(\mathbf{x} - \hat{\mathbf{x}}) \cdot m\|^2$
- **Super-Res.:** $E(\mathbf{x}; \hat{\mathbf{x}}) = \|d(\mathbf{x}) - \hat{\mathbf{x}}\|^2$

Results + Code

https://dmitryulyanov.github.io/deep_image_prior

Thanks!