# Sparsity in Convolutional Neural Networks
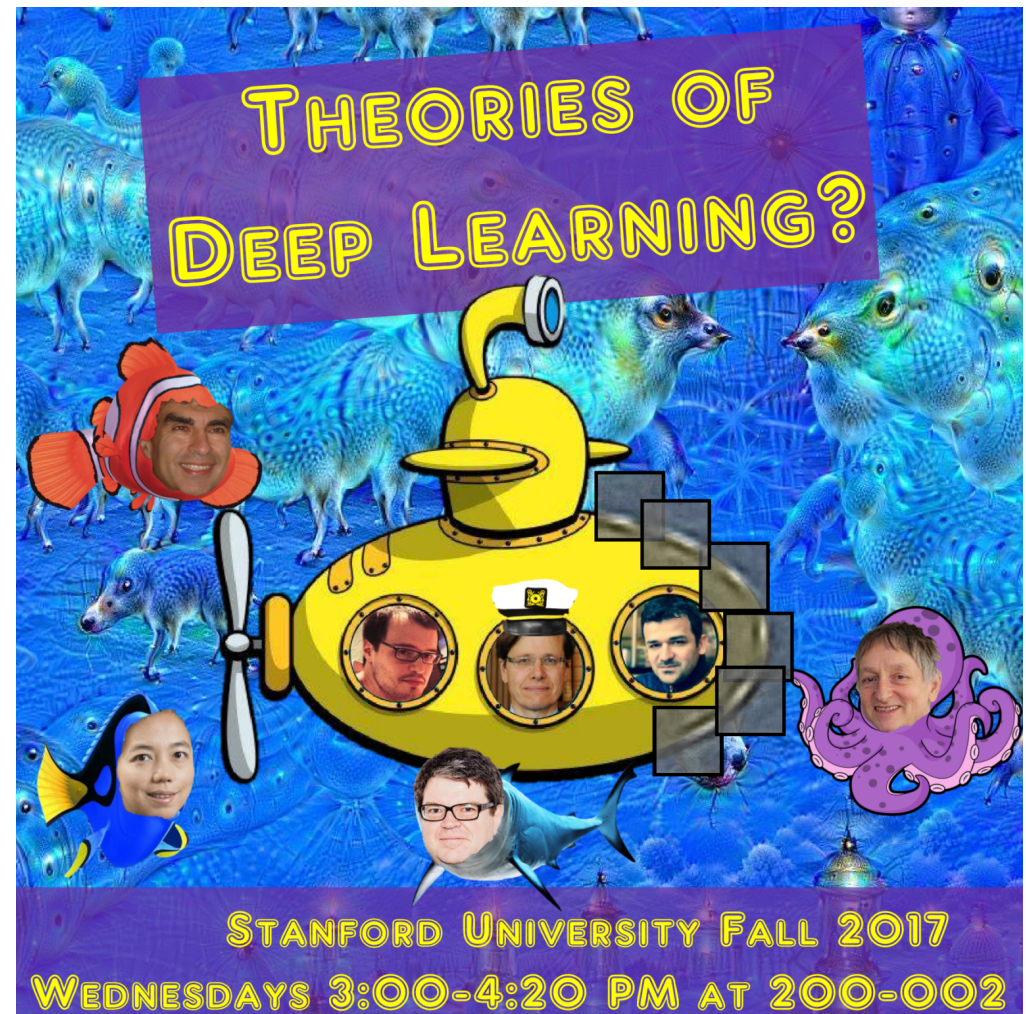
Speaker: Qingyun Sun
Math PhD @ Stanford

Acknowledgement:
Stats 385 @ Stanford
https://stats385.github
.io/

The talk is based on:
Convolutional Neural Networks in View of Sparse Coding,
Vardan Papyan @ Stats 385, Stanford

Based on work of:
Vardan Papyan, Jeremias Sulam, Yaniv Romano,
Michael Elad

# Sparsity: Central idea in Stats

Compressive Sensing:

$$\mathbf{X} = \mathbf{D}\boldsymbol{\Gamma}$$

$$\hat{\boldsymbol{\Gamma}} = \arg\min_{\boldsymbol{\Gamma}} \ \|\boldsymbol{\Gamma}\|_1 \ \ \text{s.t.} \ \ \mathbf{X} = \mathbf{D}\boldsymbol{\Gamma}$$

# Sparsity: Central idea in Stats

Lasso:

$$\mathbf{Y} = \mathbf{D}\boldsymbol{\Gamma} + \mathbf{E}$$

$$\hat{\boldsymbol{\Gamma}} = \arg\min_{\boldsymbol{\Gamma}} \frac{1}{2}\|\mathbf{Y} - \mathbf{D}\boldsymbol{\Gamma}\|_2^2 + \lambda\|\boldsymbol{\Gamma}\|_1$$
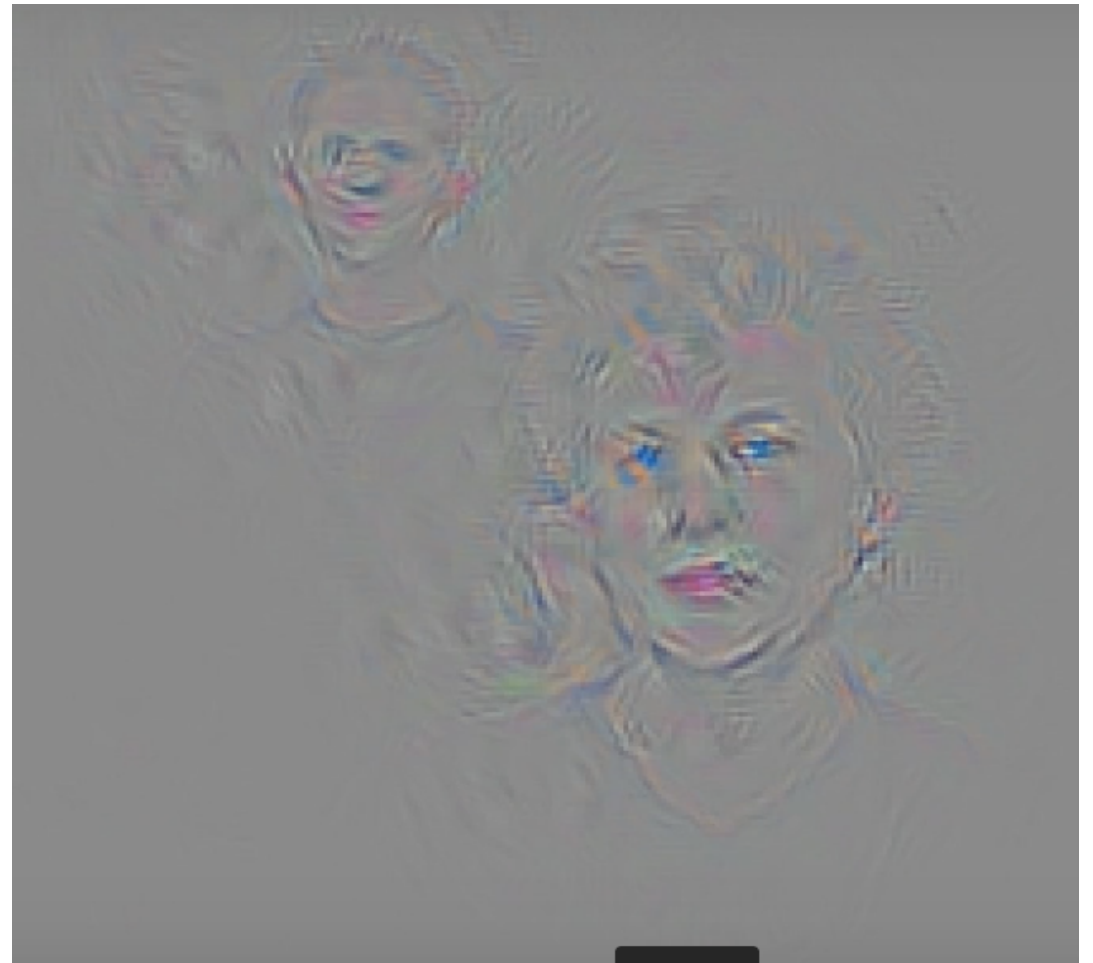
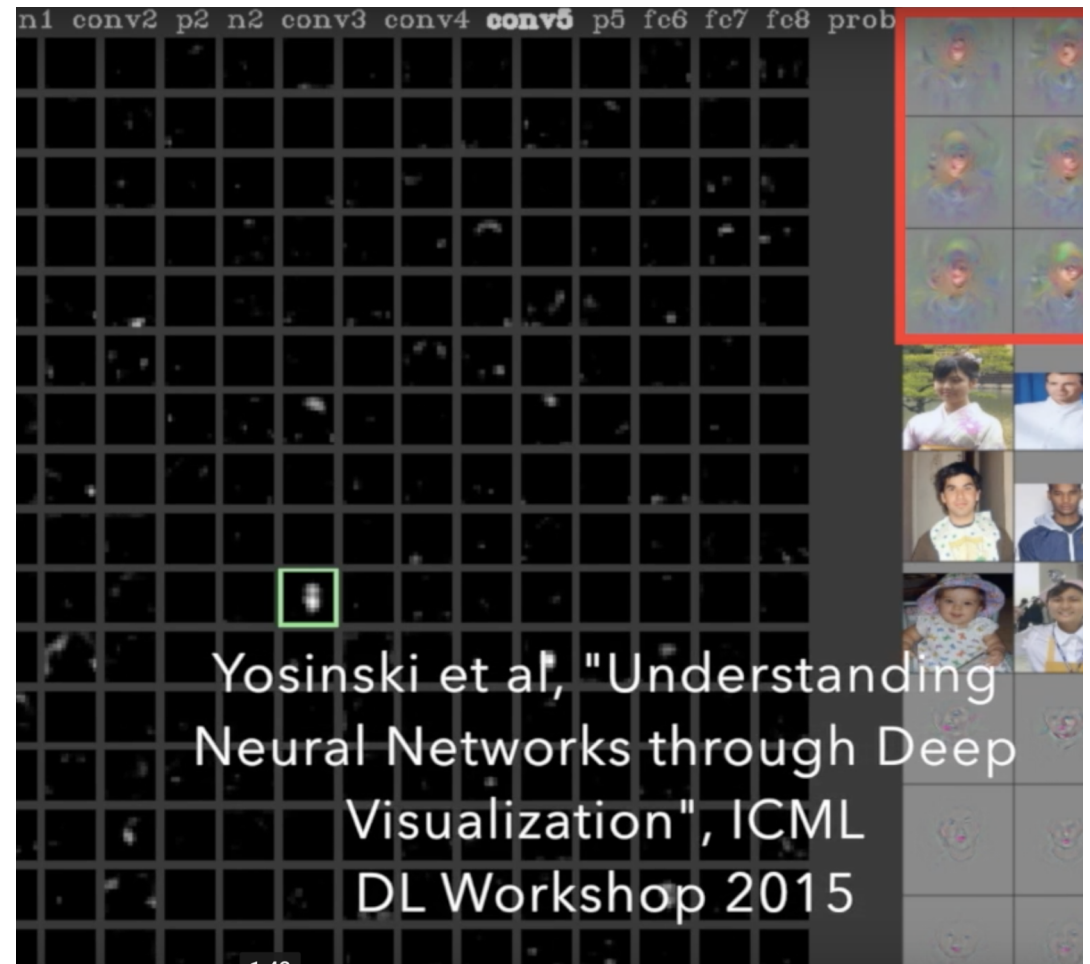# After Deep Revolution, is sparsity still important?

# Sparsity observed in CNN

# Sparsity in Practice

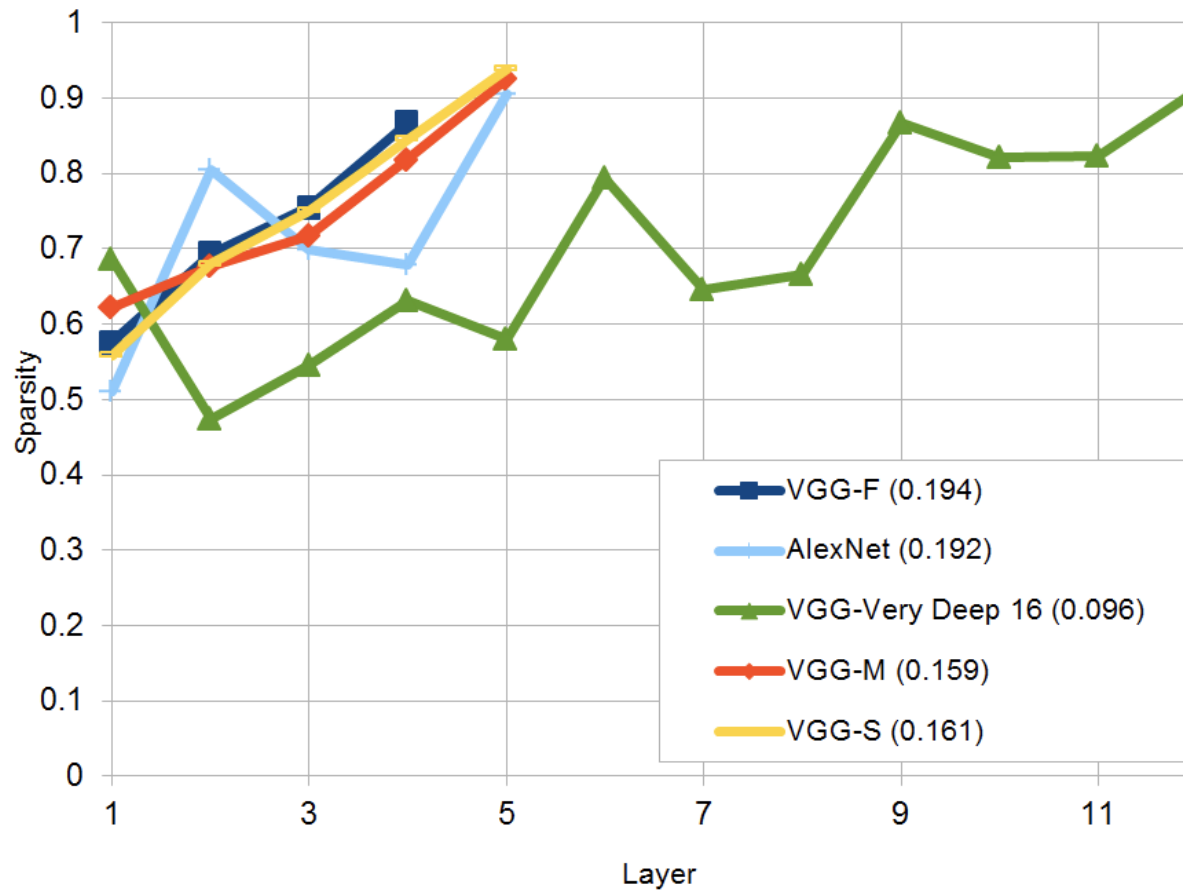The activation of RELU layer is sparse.

# Sparsity in Practice

The activation of RELU layer is sparse.



Yosinski et al, "Understanding Neural Networks through Deep Visualization", ICML DL Workshop 2015
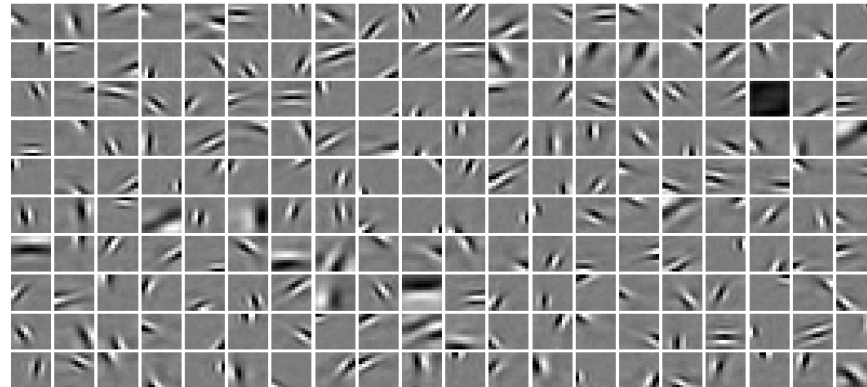
# Sparsity in Practice

# Olshausen & Field and AlexNet

Olshausen & Field

explicit sparsity



AlexNet

implicit sparsity



Credit to: Vardan, Stats385@Stanford

# Theory of sparsity in CNN?

# Breiman's "Two Cultures"
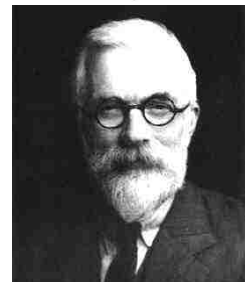
Generative modeling



Gauss    Wiener    Laplace    Bernoulli    Fisher
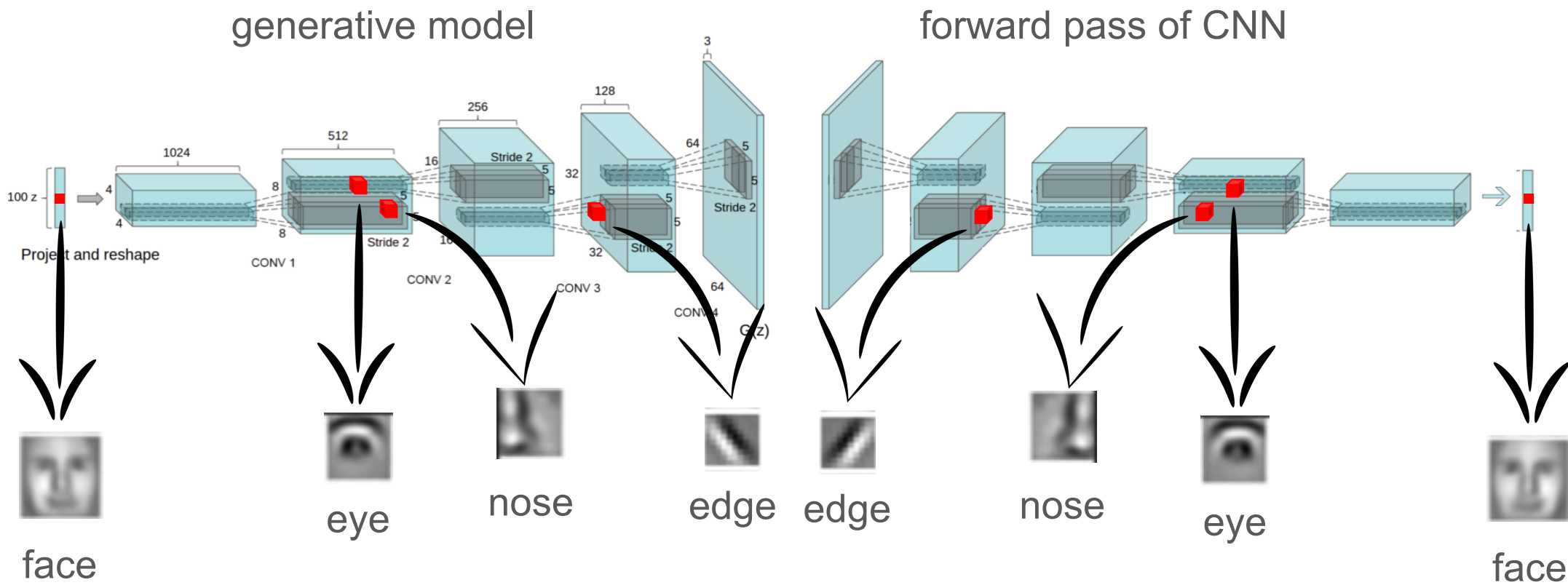
Predictive modeling

# Generative modeling

Seeks to develop stochastic models which fit the data, and then make inferences about the data-generating mechanism based on the structure of those models.

# Predictive modeling

Predictive modeling is effectively silent about the underlying mechanism generating the data, and allows for many different predictive algorithms, preferring to discuss only accuracy of prediction made by different algorithm on various datasets.

# Generative Modeling



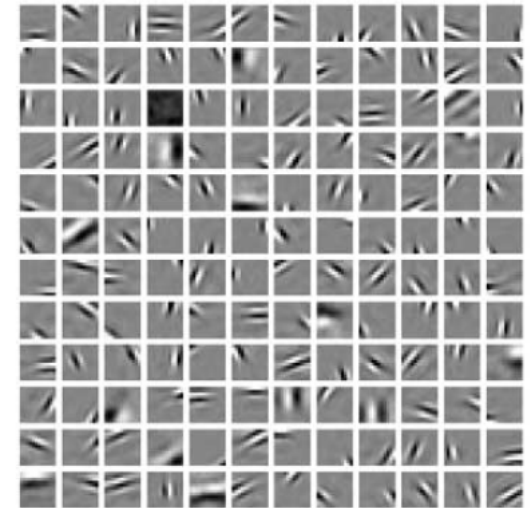generative model

forward pass of CNN
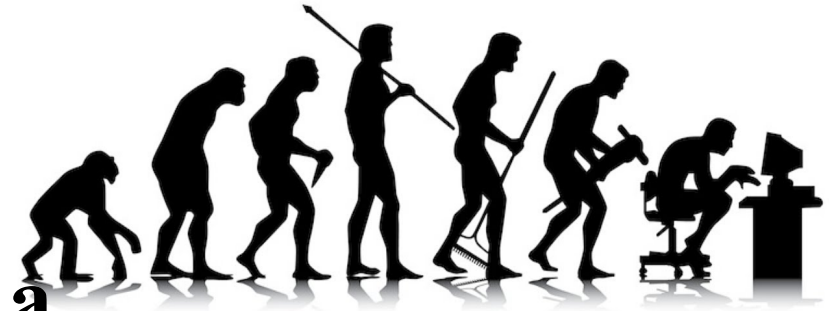
# Sparse Representation Generative Model

- Receptive fields in visual cortex are spatially localized, oriented and bandpass

- Coding natural images while promoting sparse solutions results in a set of filters satisfying these properties [Olshausen and Field 1996]

- Two decades later…
  - vast theoretical study
  - different inference algorithms
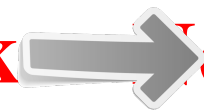  - different ways to train the model

# Evolution of Models



**Multi-Layered Convolutional Neural Network**

➡️

**First Layer of a Convolutional Neural Network**

➡️

**First Layer of a Neural Network**

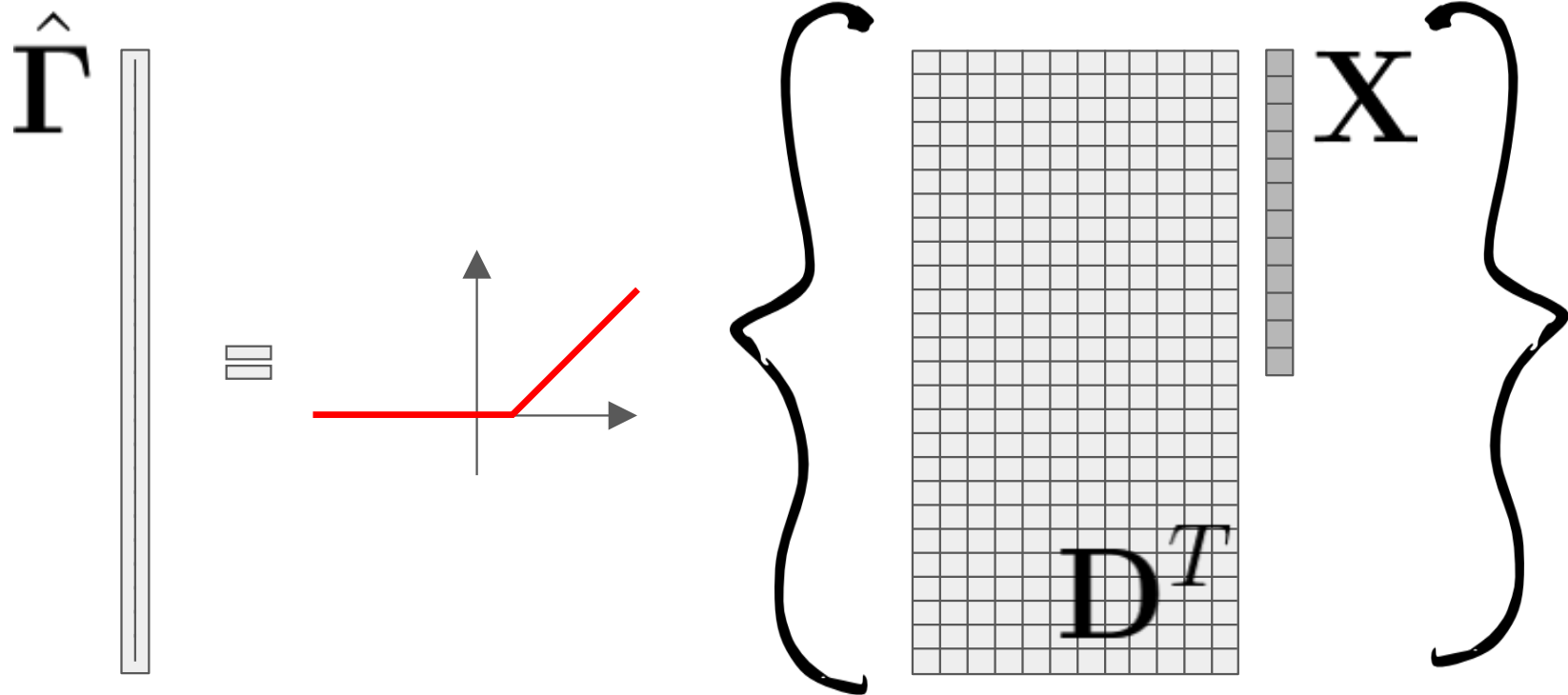**Multi-Layered Convolutional Sparse Representation**

⬅️

**Convolutional sparse representation**

⬅️

**Sparse representations**

# First Layer of a Neural Network

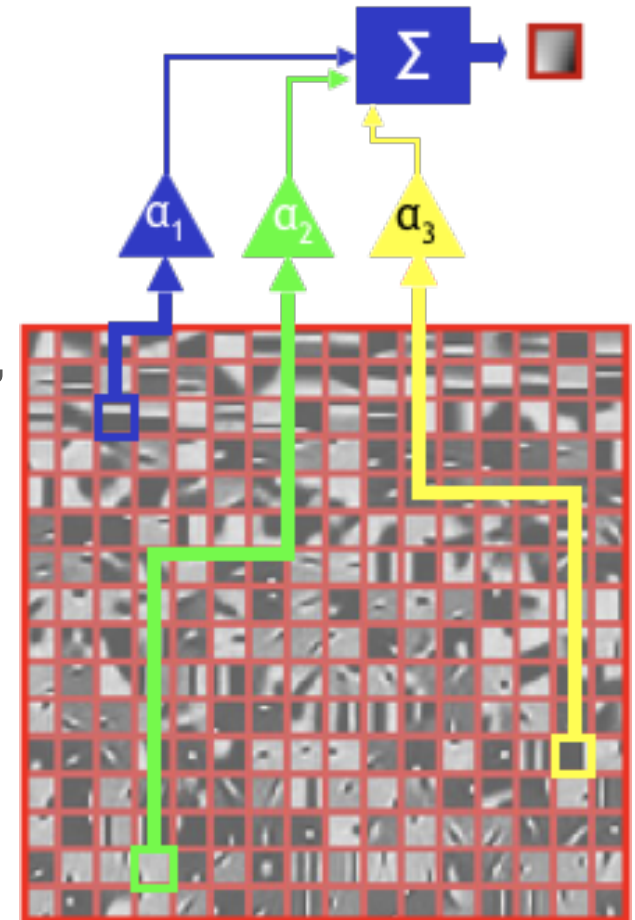$$\hat{\Gamma} = \left\{ \mathbf{D}^T \mathbf{X} \right\}$$

# Sparse Modeling

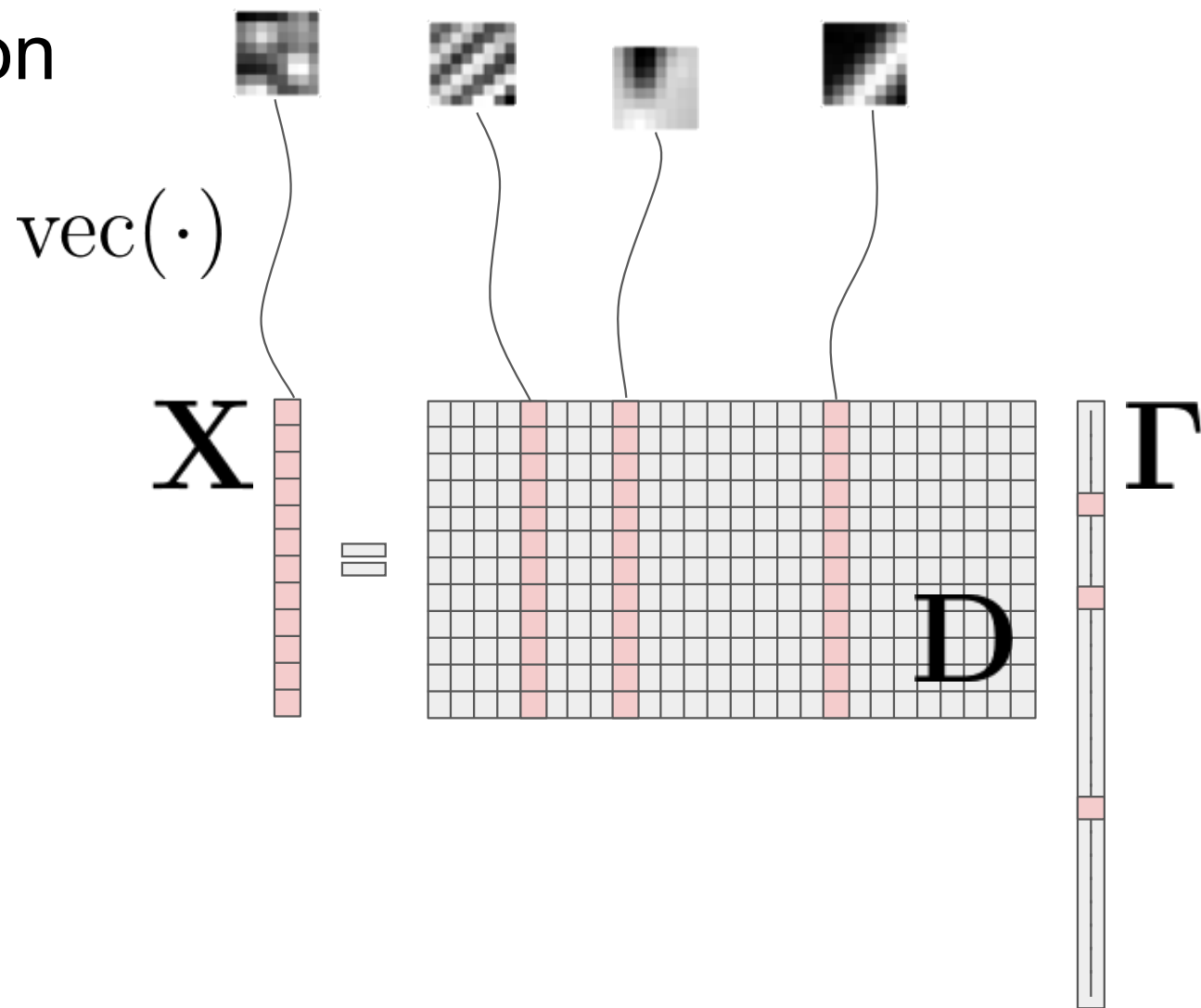Task: model image patches of size 8x8 pixels

We assume a dictionary of such image patches is given, containing 256 atoms

Assumption: every patch can be described as a linear combination of a few atoms

Key properties: sparsity and redundancy

# Matrix Notation

$$\text{vec}(\cdot)$$

# Sparse Coding

Given a signal, we would like to find its sparse representation

$$\min_{\boldsymbol{\Gamma}} \|\boldsymbol{\Gamma}\|_0 \quad \text{s.t.} \quad \mathbf{X} = \mathbf{D}\boldsymbol{\Gamma}$$

Convexify

$$\min_{\boldsymbol{\Gamma}} \|\boldsymbol{\Gamma}\|_1 \quad \text{s.t.} \quad \mathbf{X} = \mathbf{D}\boldsymbol{\Gamma}$$

# Sparse Coding

Given a signal, we would like to find its sparse representation

$$\min_{\boldsymbol{\Gamma}} \ \|\boldsymbol{\Gamma}\|_0 \ \ \text{s.t.} \ \ \mathbf{X} = \mathbf{D}\boldsymbol{\Gamma}$$

Convexify

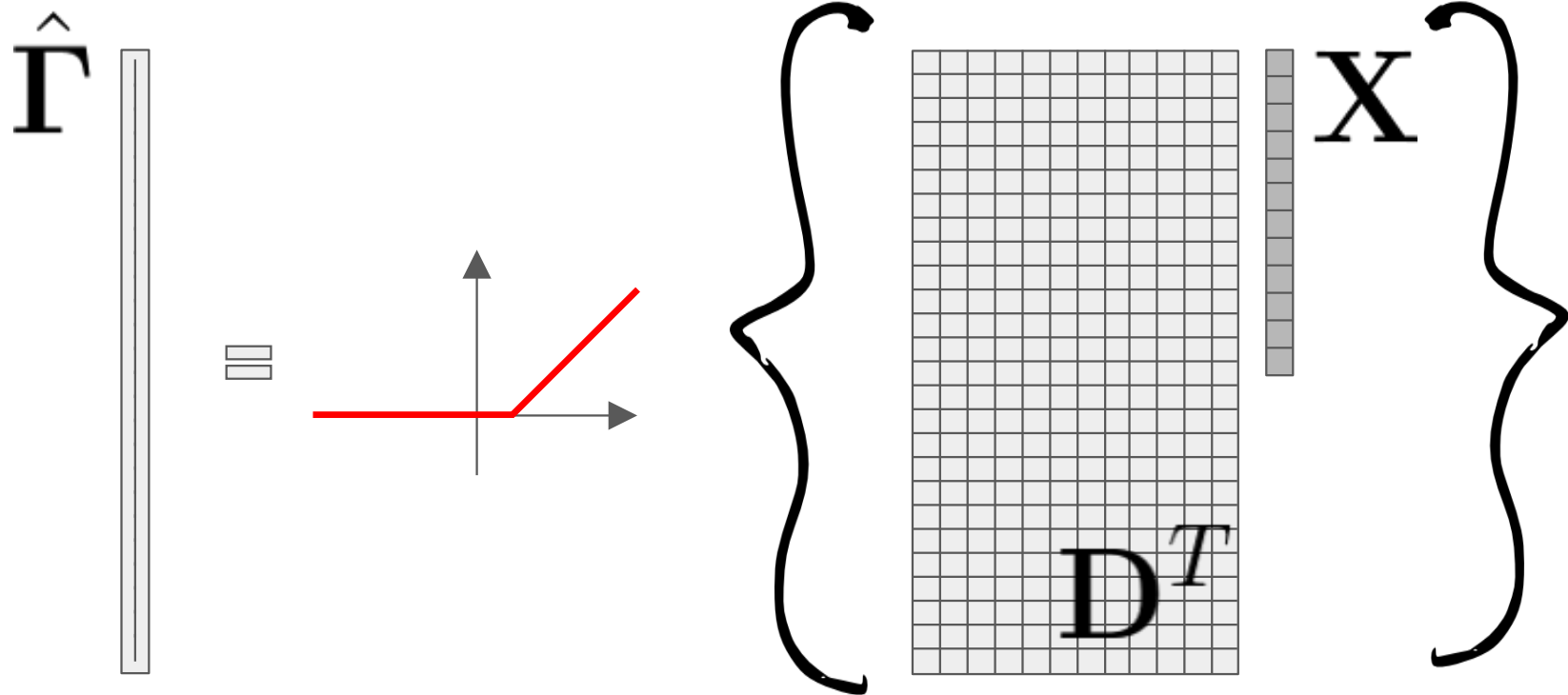$$\min_{\boldsymbol{\Gamma}} \ \|\boldsymbol{\Gamma}\|_1 \ \ \text{s.t.} \ \ \mathbf{X} = \mathbf{D}\boldsymbol{\Gamma}$$

Crude approximation

$$\mathcal{S}_{\beta}\{\mathbf{D}^T\mathbf{X}\}$$

# Thresholding Algorithm

# First Layer of a Neural Network

$$\hat{\Gamma} = \quad \left\{ \mathbf{D}^T \right\} \mathbf{X}$$

# ReLU = Soft Nonnegative Thresholding



ReLU is equivalent to soft nonnegative thresholding

# First layer of a **Convolutional** Neural Network



$$\hat{\boldsymbol{\Gamma}} = \quad \Bigg\{ \mathbf{X} * \mathbf{D} \Bigg\}$$

width

height

filters

# Convolutional Sparse Modeling

# Convolutional Sparse Modeling

$$\mathbf{X} = \sum_{k=1}^{K} \mathbf{d}_k * \mathbf{z}_k$$

# Convolutional Sparse Modeling

$$\mathbf{X} = \text{deconv} \left\{ \boldsymbol{\Gamma} , \ \mathbf{D} \right\}$$

# Thresholding Algorithm

# First layer of a Convolutional Neural Network

$$\hat{\Gamma} = \left\{ \mathbf{D}^T \right\} \mathbf{X}$$

# First layer of a Convolutional Neural Network

**X**

width

height

# Convolutional Neural Network

$$\text{ReLU}\left(\text{conv}(\mathbf{X}, \mathbf{D}_1) + \beta_1\right)$$

$\hat{\boldsymbol{\Gamma}}_1$

width

height

filters

$$\text{ReLU}\left(\text{conv}(\hat{\boldsymbol{\Gamma}}_1, \mathbf{D}_2) + \beta_2\right)$$

$\hat{\boldsymbol{\Gamma}}_2$

width

height

filters

# Multi-layered Convolutional Sparse Modeling

# Multi-layered Convolutional Sparse Modeling

$$\mathbf{X}$$



$$\mathbf{\Gamma}_2 \qquad\qquad \mathbf{\Gamma}_1$$

$$\mathrm{deconv}(\mathbf{\Gamma}_2, \mathbf{D}_2) \qquad \mathrm{deconv}(\mathbf{\Gamma}_1, \mathbf{D}_1)$$

width

height

filters

filters

height

width

height

# Layered Thresholding

$$\hat{\Gamma}_2 = \quad \mathbf{D}_2^T \quad \hat{\Gamma}_1$$

$$\hat{\Gamma}_1 = \quad \mathbf{D}_1^T \quad \mathbf{X}$$

Convolutional Neural Network

$$\hat{\Gamma}_1 = \text{ReLU}\left( \mathbf{D}_1^T \mathbf{X} \right)$$

$$\hat{\Gamma}_2 = \text{ReLU}\left( \mathbf{D}_2^T \hat{\Gamma}_1 \right)$$

**X**

width

height

# Convolutional Neural Network

$\hat{\Gamma}_1$

width

height

$\text{ReLU}\left(\text{conv}(\mathbf{X}, \mathbf{D}_1) + \beta_1\right)$

filters

$\hat{\Gamma}_2$

width

height

$\text{ReLU}\left(\text{conv}(\hat{\Gamma}_1, \mathbf{D}_2) + \beta_2\right)$

filters

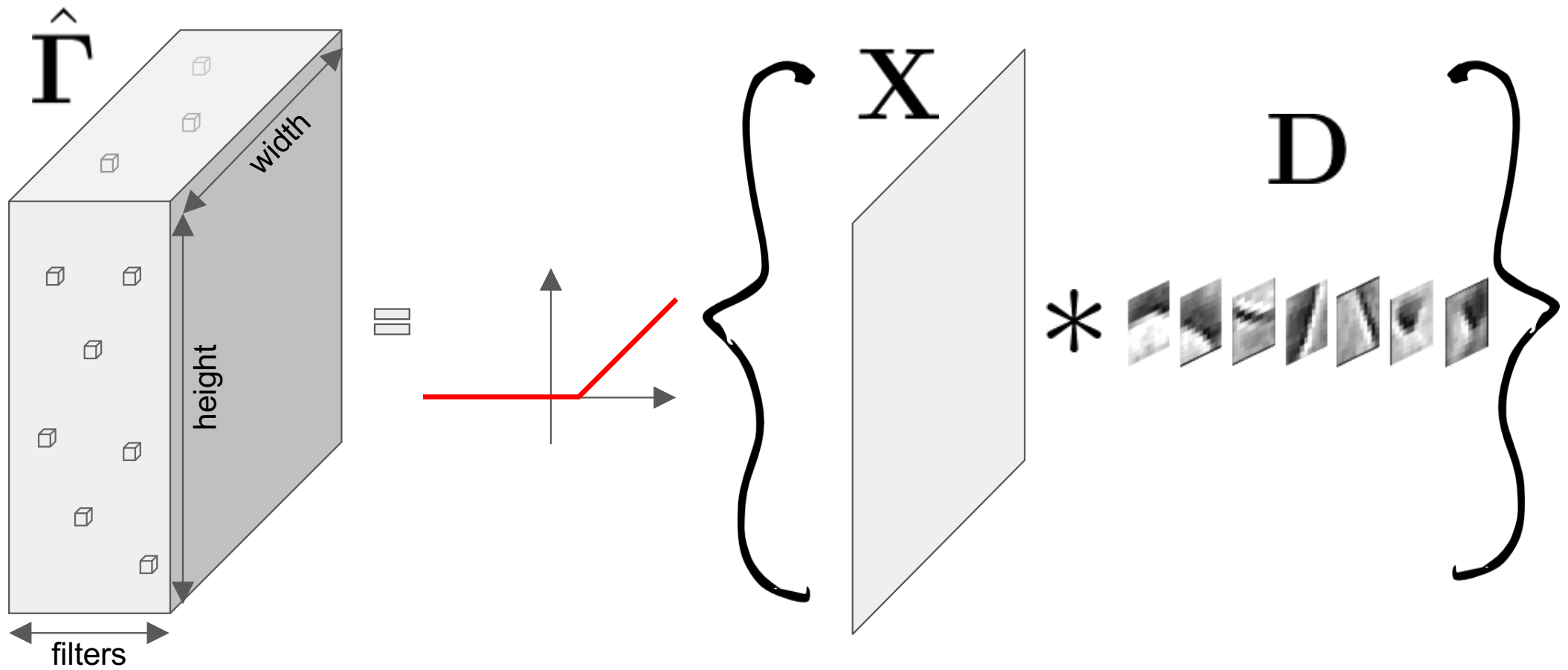# Theories of Deep Learning ?

Evolution of Models



**Multi-Layered Convolutional Neural Network**

→

**First Layer of a Convolutional Neural Network**

→

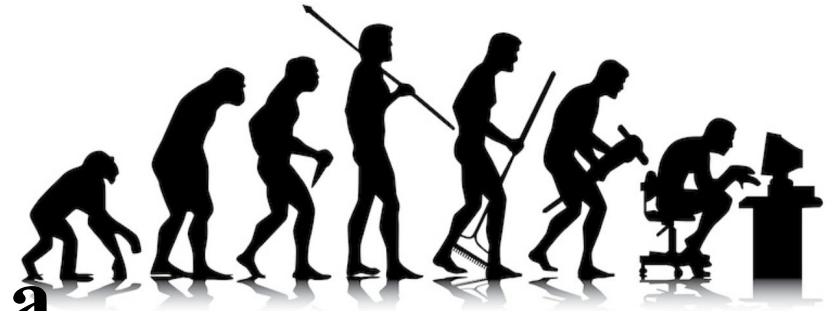**First Layer of a Neural Network**

**Multi-Layered Convolutional Sparse Representation**

←

**Convolutional sparse representation**

←

**Sparse representations**

# Sparse Modeling



$$\mathbf{X} = \mathbf{D}\,\boldsymbol{\Gamma}$$

# Classic Sparse Theory

$$\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$$

$$\hat{\mathbf{\Gamma}} = \arg\min_{\mathbf{\Gamma}} \ \|\mathbf{\Gamma}\|_1 \ \ \text{s.t.} \ \ \mathbf{X} = \mathbf{D}\mathbf{\Gamma}$$

**Theorem:** [Donoho and Elad, 2003]

Basis pursuit is guaranteed to recover the true sparse vector assuming that

$$\|\mathbf{\Gamma}\|_0 < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D})}\right)$$

Mutual Coherence:    $\mu\left(\mathbf{D}\right) = \max_{i \neq j}\left|\left(\mathbf{D}^T\mathbf{D}\right)_{i,j}\right|$

# Convolutional Sparse Modeling

# Classic Sparse Theory for Convolutional Case

**Theorem:** [Donoho and Elad, 2003]

Basis pursuit is guaranteed to recover the true sparse vector assuming that

$$\|\mathbf{\Gamma}\|_0 < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D})}\right)$$

Assuming 2 atoms of length 64    $\mu(\mathbf{D}) \geq 0.063$    [Welch, 1974]

Success guaranteed when   $\|\mathbf{\Gamma}\|_0 < 8.43$   **Very pessimistic!**

# Local Sparsity

$$\|\mathbf{\Gamma}\|_{0,\infty}$$

maximal number of non-zeroes
in a local neighborhood

$$\min_{\mathbf{\Gamma}} \|\mathbf{\Gamma}\|_{0,\infty} \quad \text{s.t.} \quad \mathbf{X} = \mathbf{D\Gamma}$$
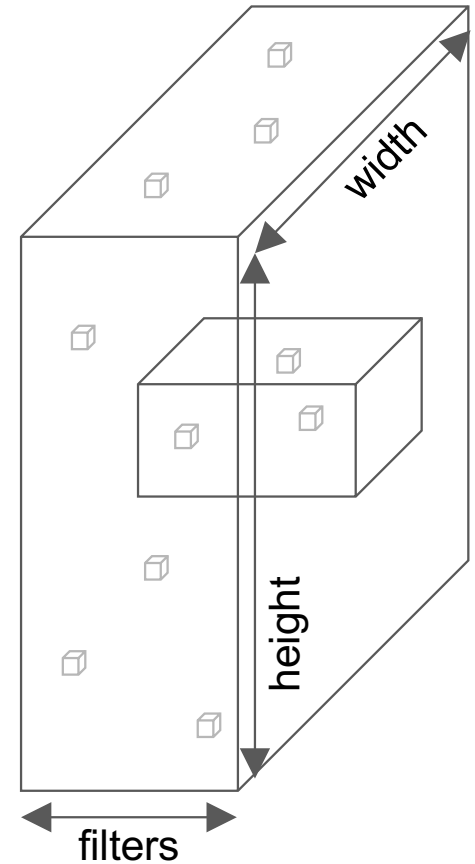


width

height

filters

# Success of Basis Pursuit

$$\mathbf{Y} = \mathbf{D}\boldsymbol{\Gamma} + \mathbf{E}$$

$$\hat{\boldsymbol{\Gamma}} = \arg\min_{\boldsymbol{\Gamma}} \ \frac{1}{2}\|\mathbf{Y} - \mathbf{D}\boldsymbol{\Gamma}\|_2^2 + \lambda\|\boldsymbol{\Gamma}\|_1$$

**Theorem:** [Papyan, Sulam and Elad, 2016]

Assume: $\|\boldsymbol{\Gamma}\|_{0,\infty} < \frac{1}{3}\left(1 + \frac{1}{\mu(\mathbf{D})}\right)$

Then: $\|\hat{\boldsymbol{\Gamma}} - \boldsymbol{\Gamma}\|_\infty \le 7.5\|\mathbf{E}\|_{2,\infty}$

Theoretical guarantee for:

- [Zeiler et. al 2010]
- [Wohlberg 2013]
- [Bristow et. al 2013]
- [Fowlkes and Kong 2014]
- [Zhou et. al 2014]
- [Kong and Fowlkes 2014]
- [Zhu and Lucey 2015]
- [Heide et. al 2015]
- [Gu et. al 2015]
- [Wohlberg 2016]
- [Šorel and Šroubek 2016]
- [Serrano et. al 2016]
- [Papyan et. al 2017]
- [Garcia-Cardona and Wohlberg 2017]
- [Wohlberg and Rodriguez 2017]
- ...

# Multi-layered Convolutional Sparse Modeling

## Deep Coding Problem

Given $\mathbf{X}$, find a set of representations satisfying:

$$\mathbf{X} = \mathbf{D}_1 \boldsymbol{\Gamma}_1, \qquad \|\boldsymbol{\Gamma}_1\|_{0,\infty} \leq \lambda_1$$

$$\boldsymbol{\Gamma}_1 = \mathbf{D}_2 \boldsymbol{\Gamma}_2, \qquad \|\boldsymbol{\Gamma}_2\|_{0,\infty} \leq \lambda_2$$

$$\vdots$$

$$\boldsymbol{\Gamma}_{L-1} = \mathbf{D}_L \boldsymbol{\Gamma}_L, \qquad \|\boldsymbol{\Gamma}_L\|_{0,\infty} \leq \lambda_L$$

## Deep Coding Problem

Given $\mathbf{Y}$, find a set of representations satisfying:

$$\|\mathbf{Y} - \mathbf{D}_1\boldsymbol{\Gamma}_1\|_2 \leq \epsilon, \qquad \|\boldsymbol{\Gamma}_1\|_{0,\infty} \leq \lambda_1$$

$$\boldsymbol{\Gamma}_1 = \mathbf{D}_2\boldsymbol{\Gamma}_2, \qquad \|\boldsymbol{\Gamma}_2\|_{0,\infty} \leq \lambda_2$$

$$\vdots$$

$$\boldsymbol{\Gamma}_{L-1} = \mathbf{D}_L\boldsymbol{\Gamma}_L, \qquad \|\boldsymbol{\Gamma}_L\|_{0,\infty} \leq \lambda_L$$

# Uniqueness

## Uniqueness Theorem

$$\|\mathbf{\Gamma}_l\|_{0,\infty} \leq \lambda_l < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D}_l)}\right)$$

$\{\mathbf{\Gamma}_l\}_{l=1}^{L}$ are the unique feature maps of $\mathbf{X}$

# Success of Forward Pass

# Success of Forward Pass Theorem

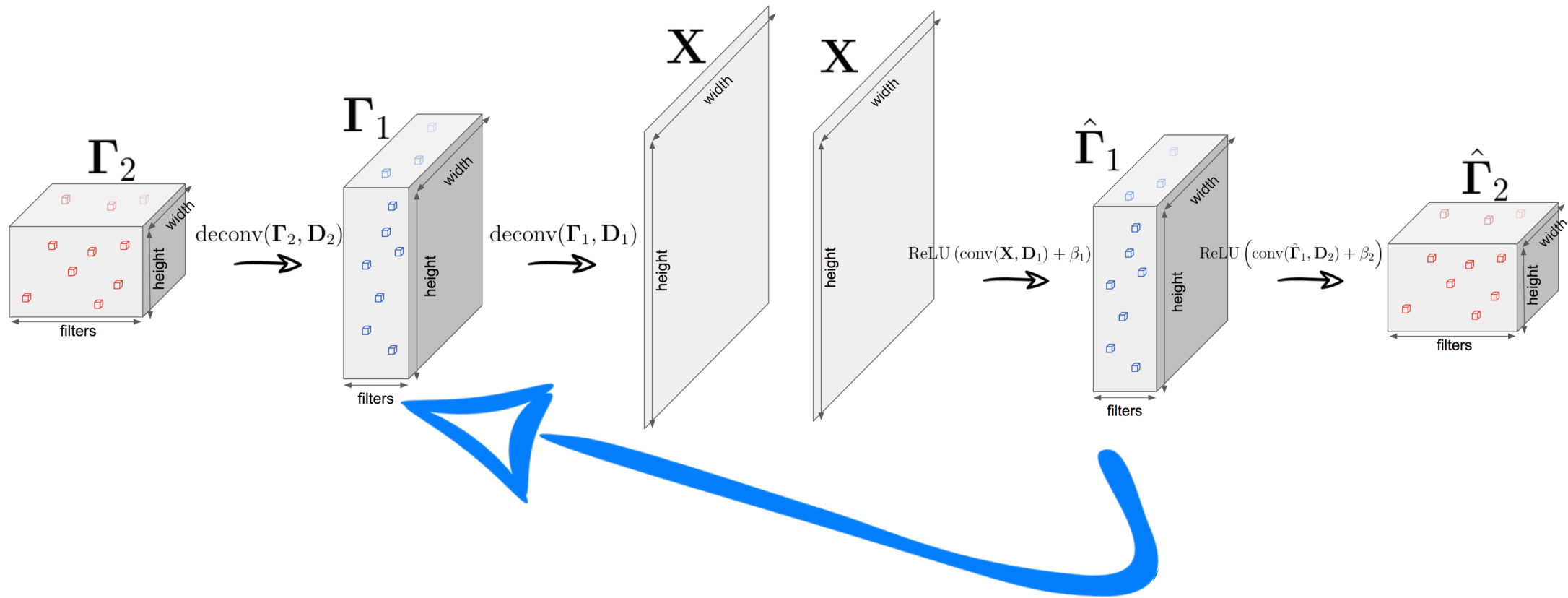$$\|\mathbf{\Gamma}_l\|_{0,\infty} < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D}_l)}\frac{|\Gamma_l^{\min}|}{|\Gamma_l^{\max}|}\right) - \frac{1}{\mu(\mathbf{D}_l)}\frac{\epsilon_{l-1}}{|\Gamma_l^{\max}|}$$
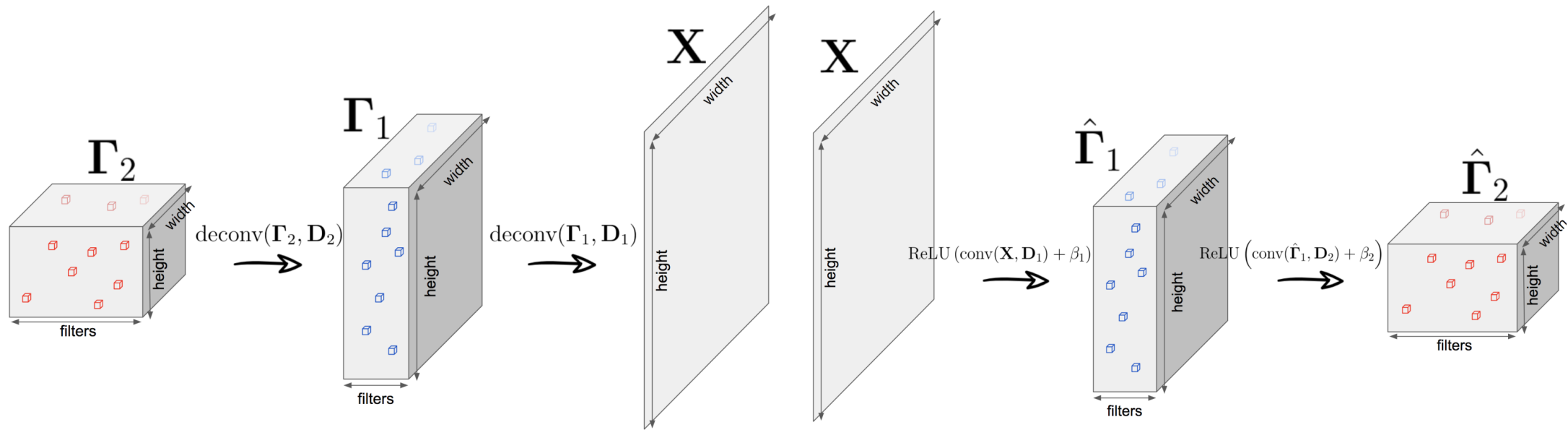
Layered thresholding guaranteed:

1. Find correct places of nonzeros

$$\|\hat{\mathbf{\Gamma}}_l - \mathbf{\Gamma}_l\|_{2,\infty} \le \epsilon_l$$

❌ Forward pass always fails at recovering representations exactly

✗ Success depends on ratio

❌ Distance increases with layer

# Generative Model and Crude Inference

Layered Lasso

# StatsDepartment

$$\hat{\boldsymbol{\Gamma}}_1 = \arg\min_{\boldsymbol{\Gamma}_1} \frac{1}{2}\|\mathbf{Y} - \mathbf{D}_1\boldsymbol{\Gamma}_1\|_2^2 + \alpha_1\|\boldsymbol{\Gamma}_1\|_1$$

$$\hat{\boldsymbol{\Gamma}}_2 = \arg\min_{\boldsymbol{\Gamma}_2} \frac{1}{2}\|\hat{\boldsymbol{\Gamma}}_1 - \mathbf{D}_2\boldsymbol{\Gamma}_2\|_2^2 + \alpha_2\|\boldsymbol{\Gamma}_2\|_1$$

# Success of Layered Lasso

$$\|\mathbf{\Gamma}_l\|_{0,\infty} < \frac{1}{3}\left(1 + \frac{1}{\mu(\mathbf{D}_L)}\right)$$

Layered Lasso guaranteed:

1. Find only correct places of nonzeros
2. Find all coefficients that are big enough

$$\|\hat{\mathbf{\Gamma}}_l - \mathbf{\Gamma}_l\|_{2,\infty} \leq \epsilon_l$$

❌ ~~Forward pass always fails at recovering representations exactly~~

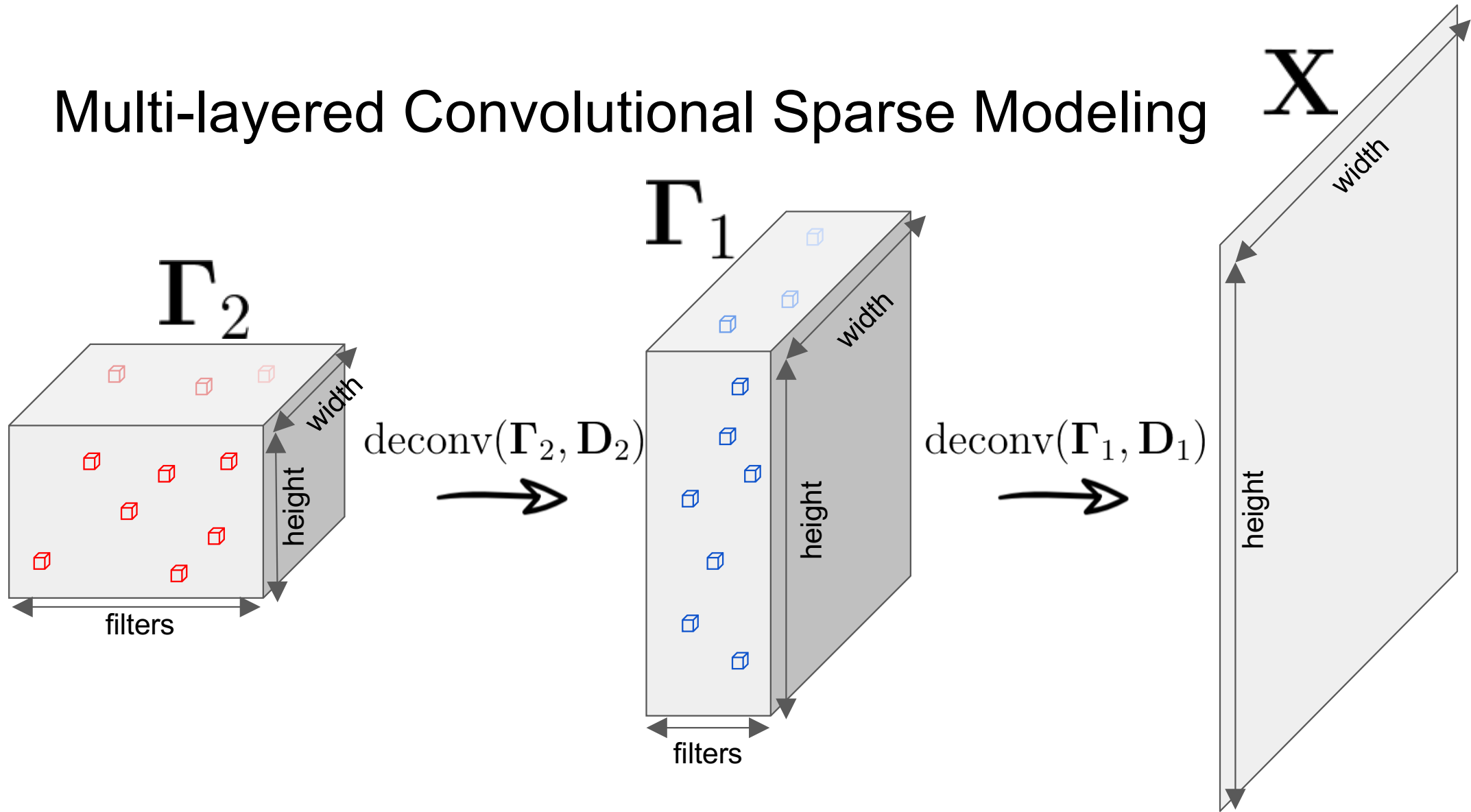❌ ~~Success depends on ratio~~

❌ Distance increases with layer

# Layered Iterative Thresholding

$$\mathbf{\Gamma}_1^t = \mathcal{S}_{\alpha_1} \left( \mathbf{D}_1^T \mathbf{Y} + \left( \mathbf{I} - \mathbf{D}_1^T \mathbf{D}_1 \right) \mathbf{\Gamma}_1^{t-1} \right)$$

$$\mathbf{\Gamma}_2^t = \mathcal{S}_{\alpha_2} \left( \mathbf{D}_2^T \hat{\mathbf{\Gamma}}_1 + \left( \mathbf{I} - \mathbf{D}_2^T \mathbf{D}_2 \right) \mathbf{\Gamma}_2^{t-1} \right)$$

# Multi-layered Convolutional Sparse Modeling

$\mathbf{\Gamma}_2$     $\mathbf{\Gamma}_1$     $\mathbf{X}$

width    height    filters

$\mathrm{deconv}(\mathbf{\Gamma}_2, \mathbf{D}_2)$

width    height    filters

$\mathrm{deconv}(\mathbf{\Gamma}_1, \mathbf{D}_1)$

width    height

# Summary



1 — Sparsity well established theoretically

2 — Sparsity is covertly exploited in practice: ReLU, dropout, stride, dilation, ...

3 — Sparsity is the secret sauce behind CNN

4 — Need to bring sparsity to the surface to better understand CNNs

5 — Andrej Karpathy agrees