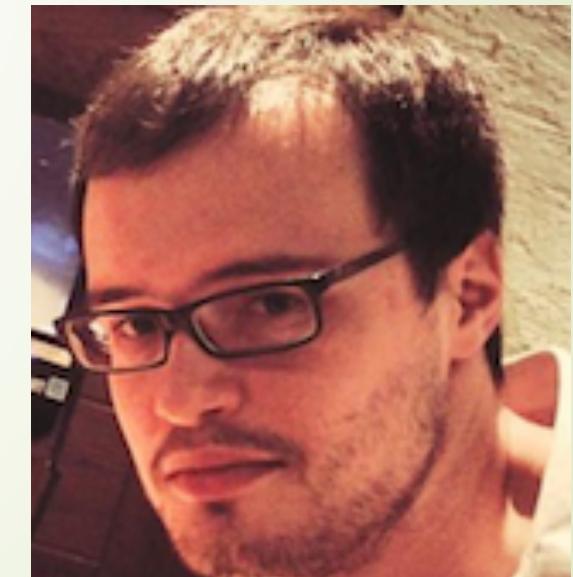


1

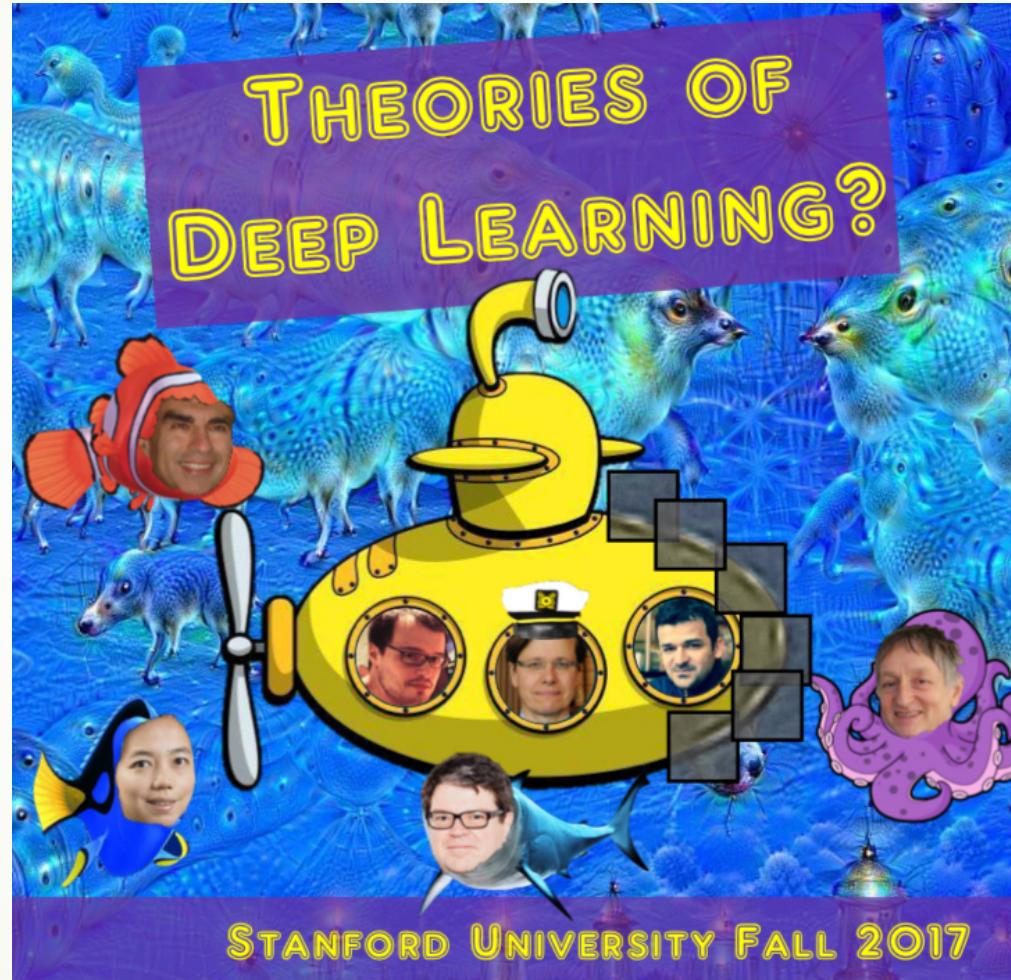
Sparsity in Convolutional Neural Networks

Yuan YAO
HKUST

Based on Qingyun Sun, Vardan Papyan talks etc.



Acknowledgement



A following-up course at HKUST: <https://deeplearning-math.github.io/>

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:

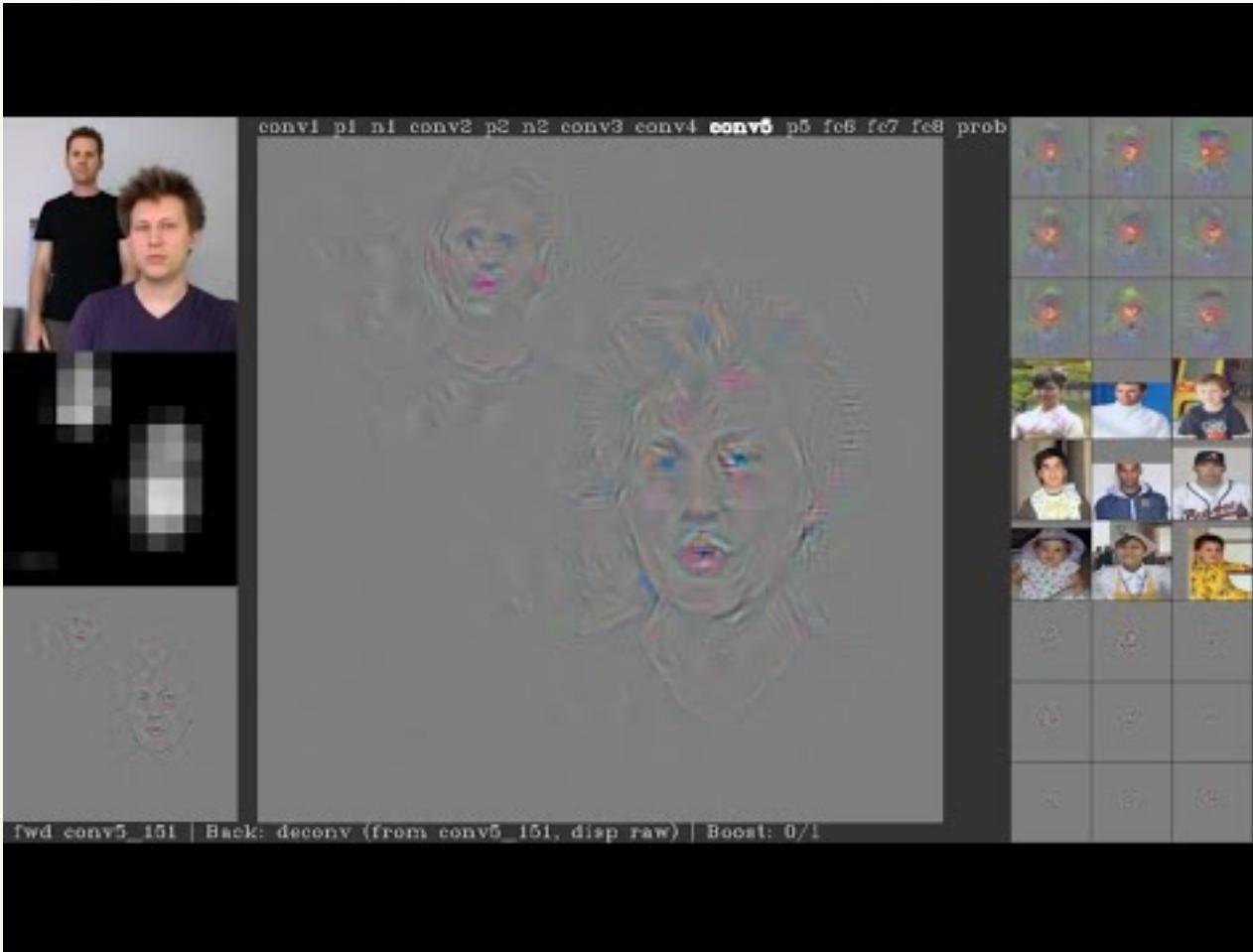
$$\mathbf{Y} = \mathbf{D}\boldsymbol{\Gamma} + \mathbf{E}$$

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



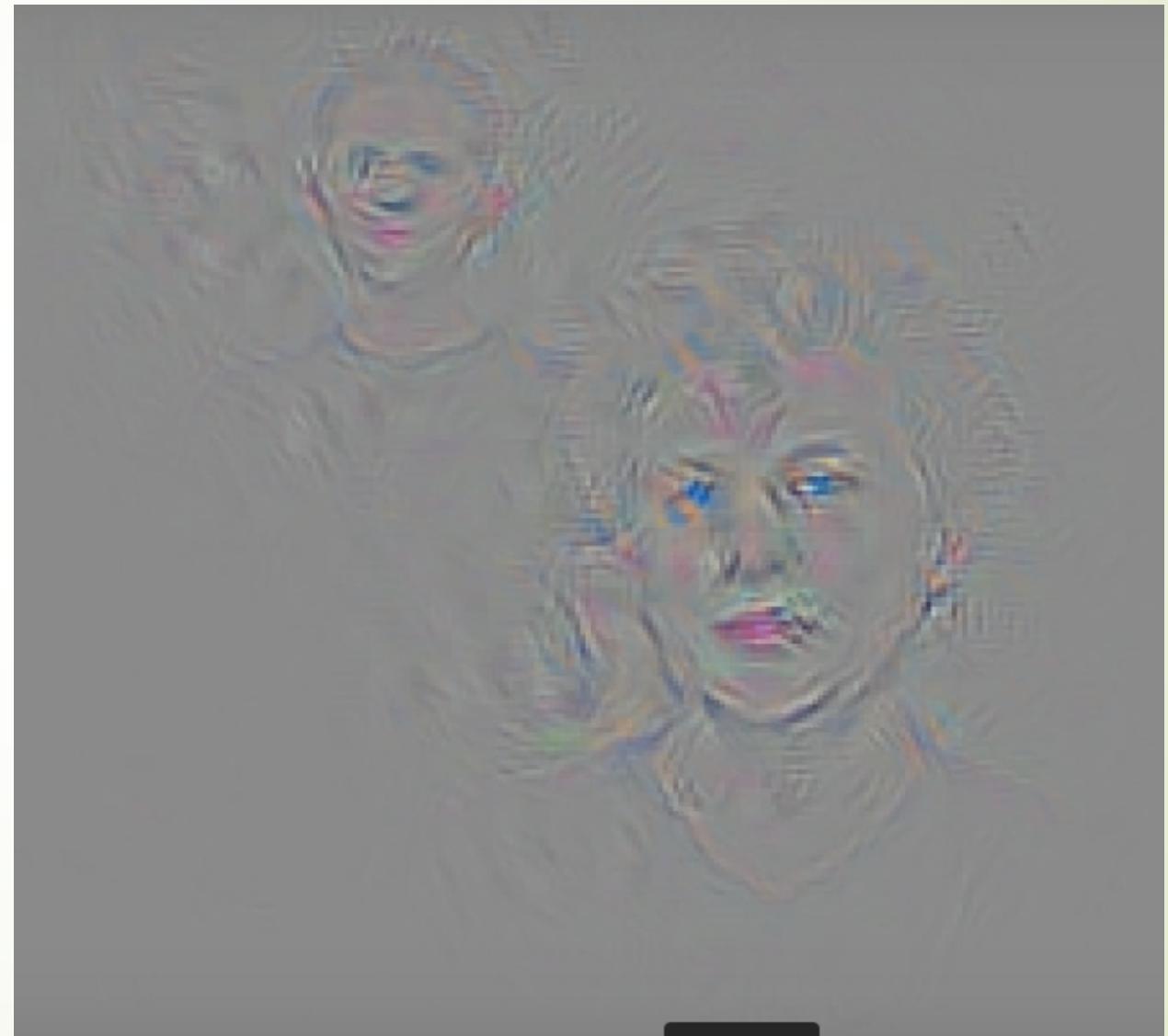
After Deep Revolution, is
sparsity still important?

Sparsity observed in CNN



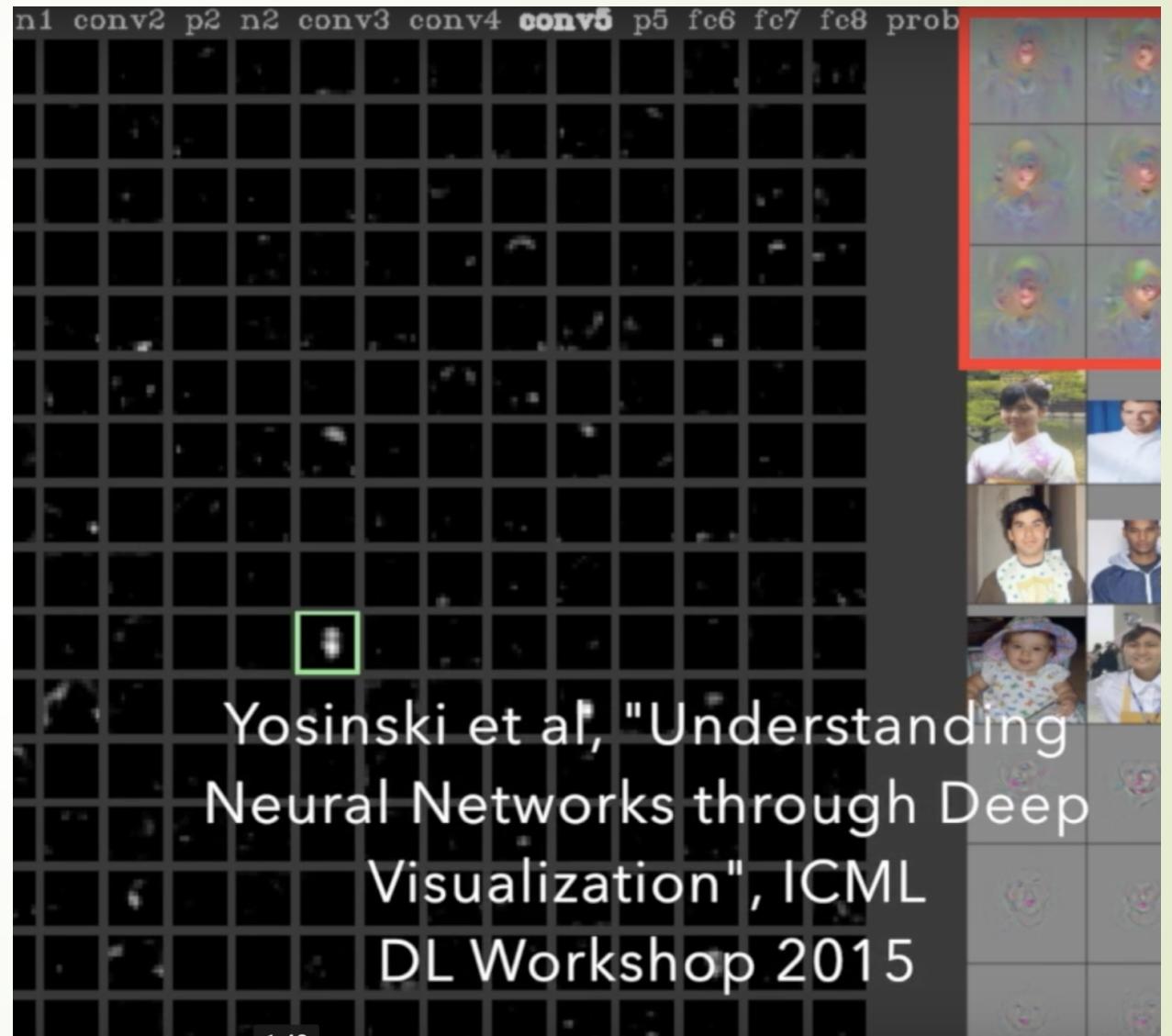
Sparsity in Practice

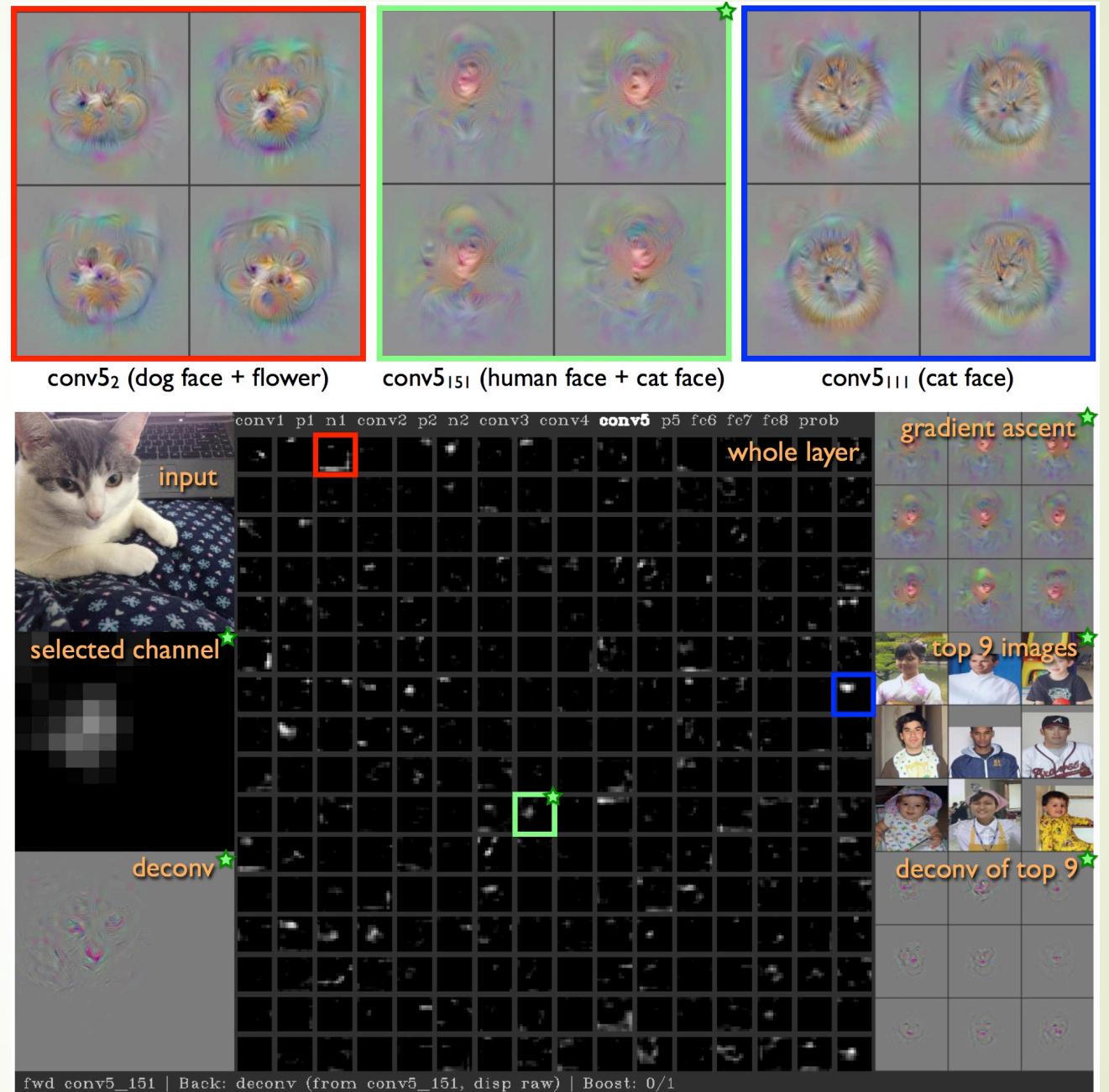
The activation of RELU layer is sparse.



Sparsity in Practice

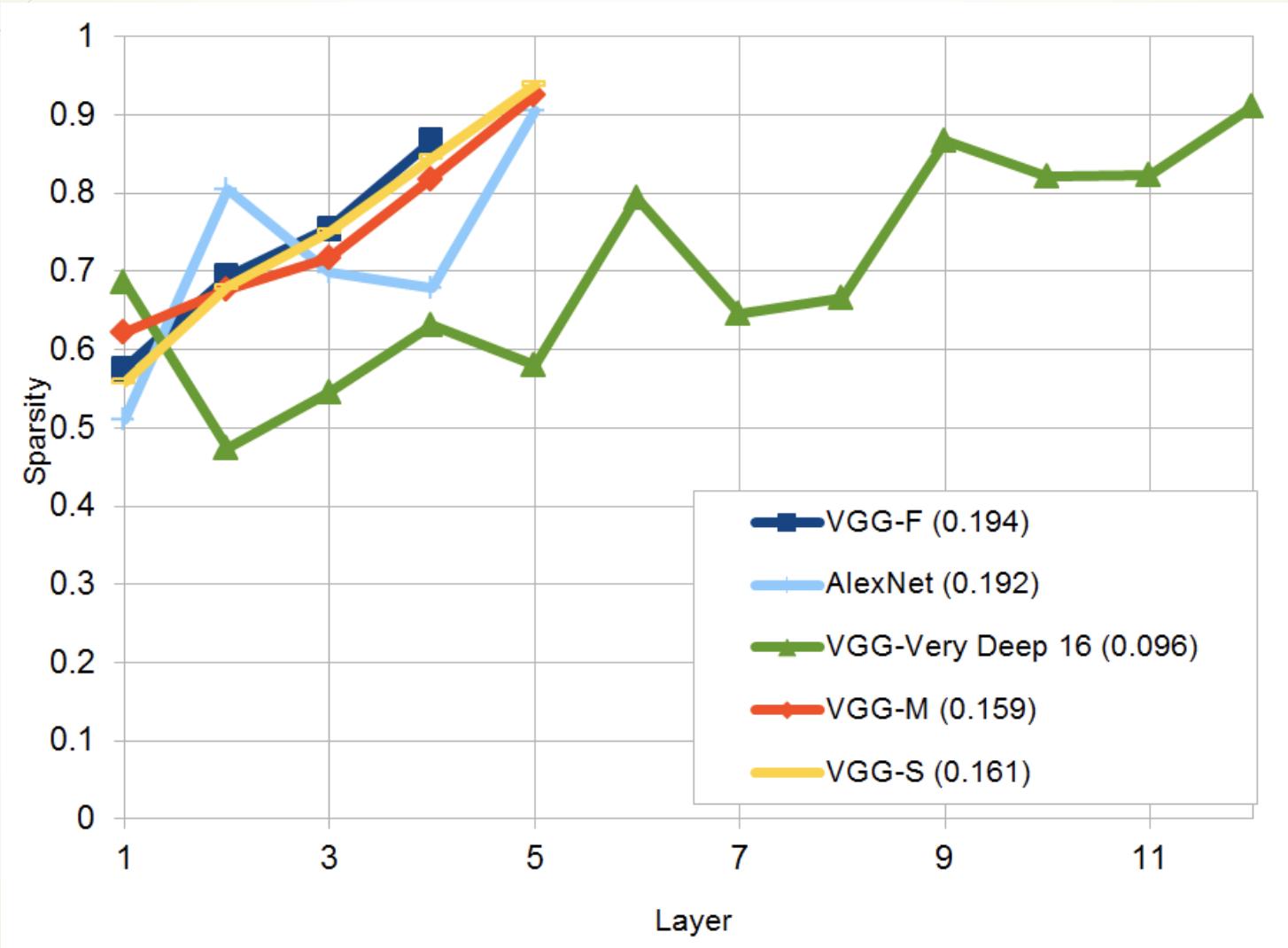
The activation of RELU layer is sparse.





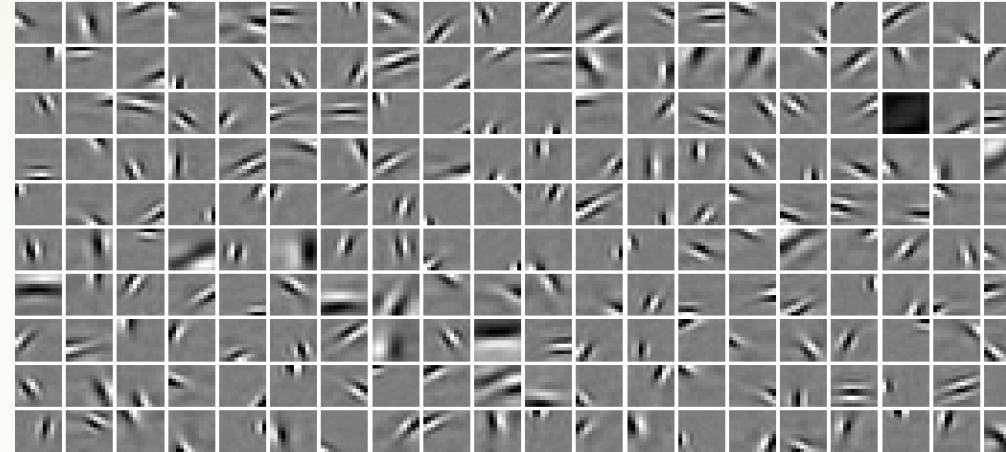
Yosinski et al. Understanding
Neural Networks Through Deep
Visualization. ICML workshop
2015.

Sparsity in Practice



Olshausen & Field and AlexNet

Olshausen & Field



explicit sparsity

AlexNet



implicit sparsity

Credit to: Vardan, Stats385@Stanford

[Fu, Li, Zhang, Sun, Yao, 2018]

- ▶ Pruning in Training, submitted to ICLR, 2018.

Ratio	100%	25%	12.5%	6.25%	3.13%	1.57%
Plain	99.12±0.02	80.46±5.46	62.61±9.05	45.49±0.97	32.34±1.53	21.30±4.83
Rand	99.19±0.01	62.23±10.12	37.71±4.34	23.58±6.96	18.58±4.70	14.36±3.27
R-P	99.16±0.06	75.47±8.11	60.31±5.19	37.97±2.99	26.11±2.13	18.11±1.05
L-P	98.95±0.04	98.95±0.04	90.29±1.30	60.37±5.45	32.91±3.35	20.31±1.10
S-P	98.97±0.07	98.96±0.08	98.67±0.15	68.27±11.22	42.10±7.83	24.95±8.92

(a) Pruning the conv.c5 layer

Ratio	100%	25%	12.5%	6.25%	3.13%	1.57%
Plain	99.10±0.02	96.73±1.05	95.65±1.76	89.60±3.49	78.40±5.48	64.17±6.93
Rand	99.09±0.01	91.56±3.89	71.05±6.08	51.92±8.87	33.65±6.17	29.91±8.49
R-P	99.13±0.05	96.39±0.48	95.31±0.79	91.29±3.46	82.75±5.59	68.35±6.29
L-P	99.10±0.03	98.89±0.05	98.89±0.05	98.89±0.05	98.89±0.05	98.89±0.05
S-P	99.13±0.03	98.73±0.10	98.61±0.15	98.23±0.40	96.75±0.88	92.53±3.17

(a) Pruning the fc.c6 layer

(%)	100	50	25	12.5	6.25	3.13	1.57
Plain	83.92	13.53	8.12	5.32	5.29	5.92	6.31
Rand	82.36	17.90	6.52	6.38	7.90	9.58	8.67
R-P	82.79	13.72	7.10	6.38	6.29	6.52	5.76
L-P / GL-P	81.09	81.09	76.43	75.06	68.42	55.25	33.49
S-P / GS-P	78.95	77.81	73.92	70.65	68.67	67.58	65.17
Com-Rat(%)	100	63.00	44.51	35.26	30.64	28.33	27.17

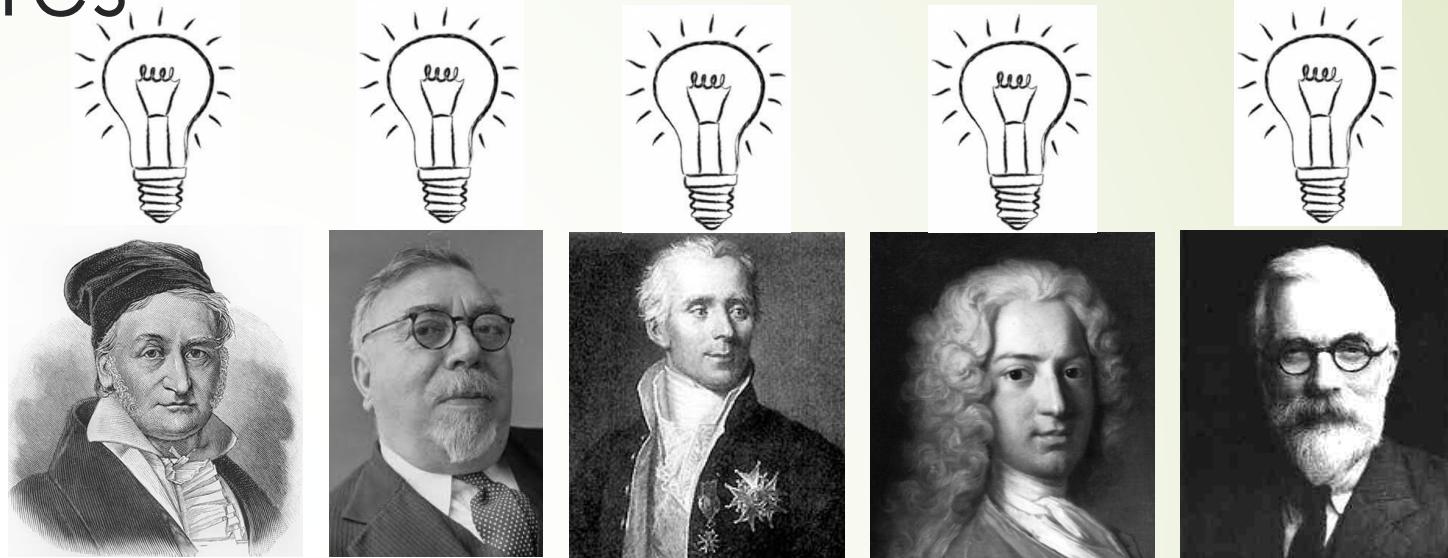
Table 5: Top 5 accuracy on *miniImagenet* by pruning ResNet-18, the fully connected layer, Block#4.0 and #4.1 layers.



Theory of Sparsity in CNN?

Breiman's “Two Cultures”

Generative modeling



Gauss

Wien

Laplac

Bernoul

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.

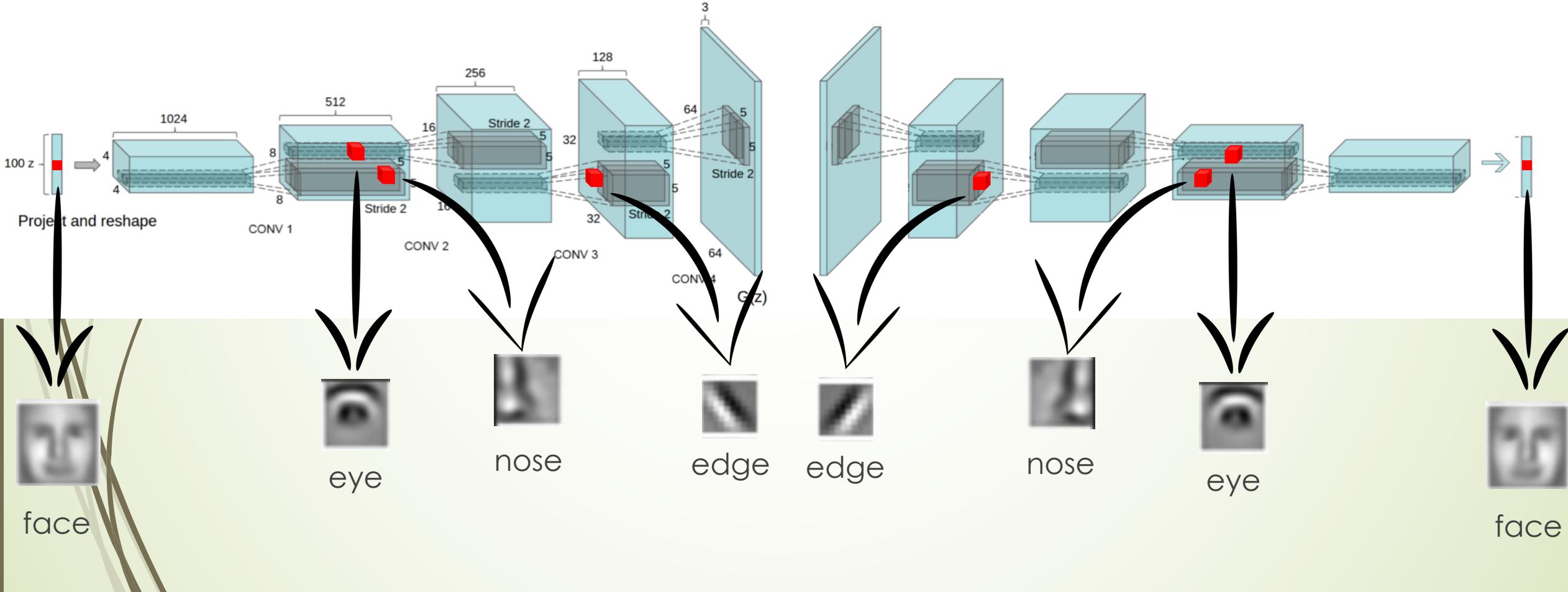
Predictive modeling

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

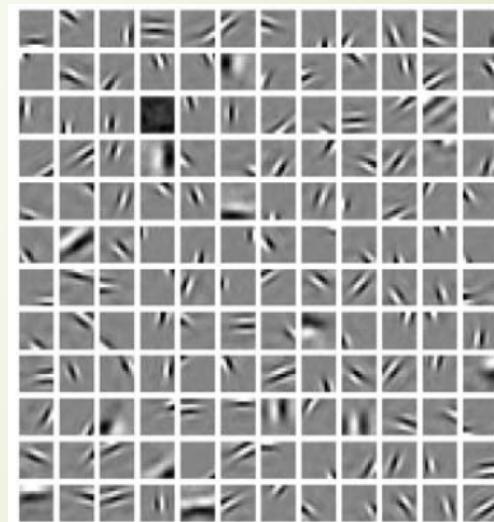


generative model



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

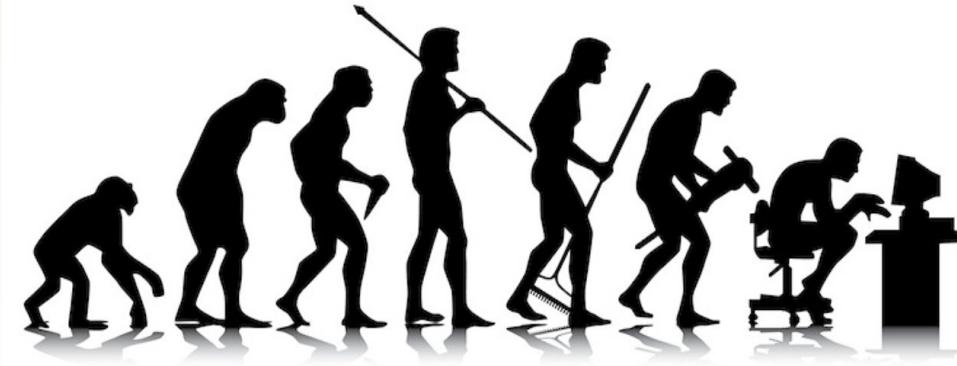
Multi-Layered
Convolutional
Sparse
Representation

First Layer of a
Convolutional
Neural Network

Convolutional
sparse
representation

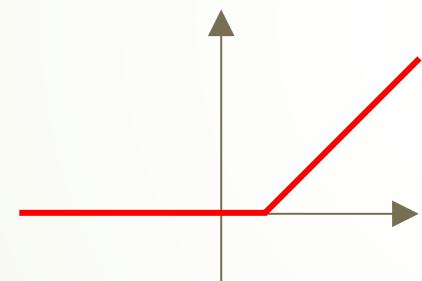
First Layer of a
Neural Network

Sparse
representations



First Layer of a Neural Network

$\hat{\Gamma}$



D^T

X

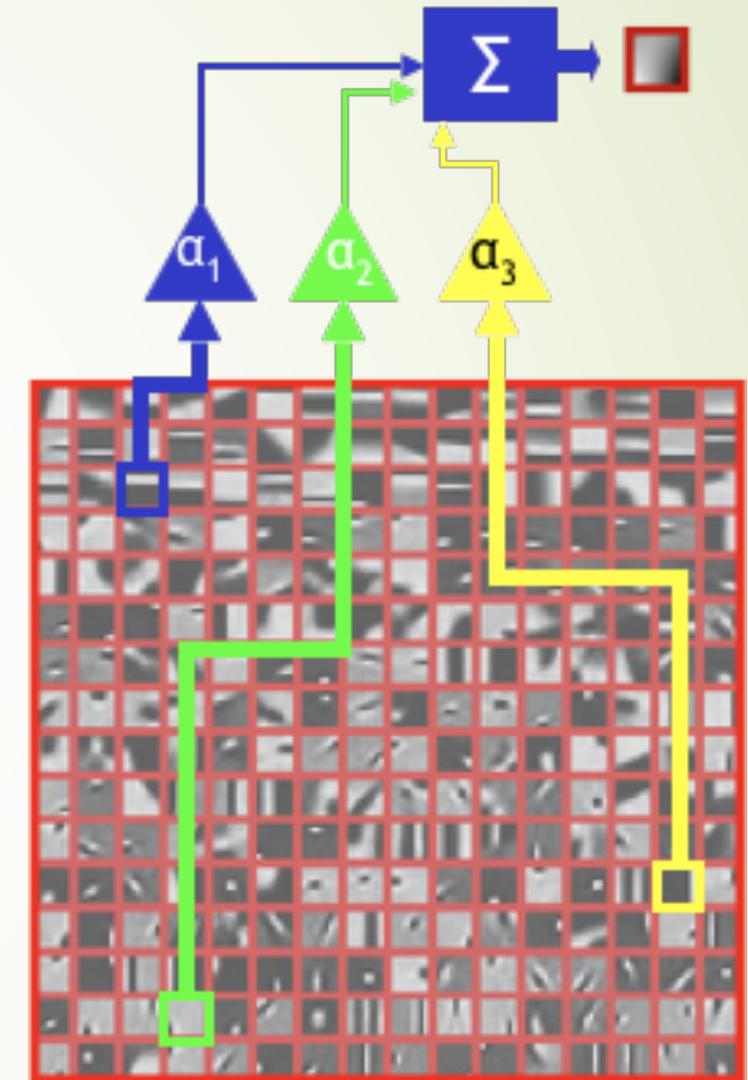
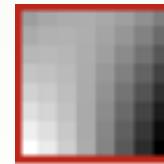
Sparse Modeling, e.g.

Task: model image patches of size 8x8 pixels

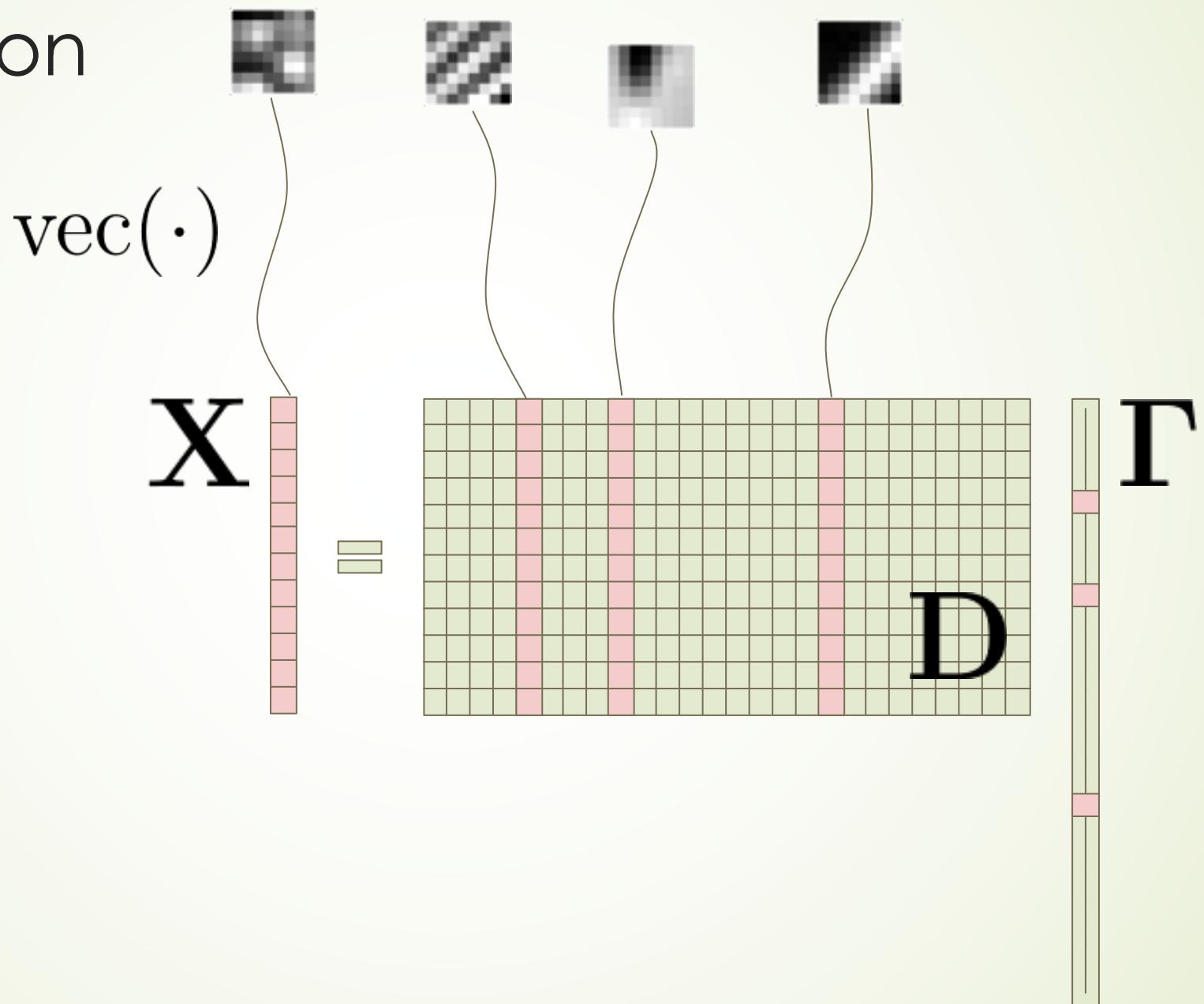
We assume a dictionary of such image patches is given, containing 256 atoms
(overcomplete)

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

Key properties: **sparsity** and **redundancy**



Matrix Notation



Sparse Coding

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

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

Convexify


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

Sparse Coding

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

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

Convexify

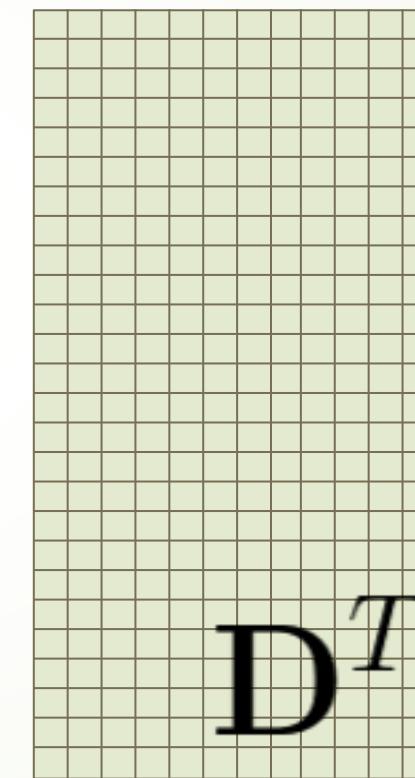
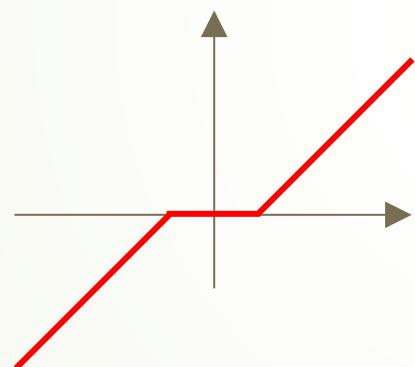
$$\min_{\Gamma} \|\Gamma\|_1 \text{ s.t. } \mathbf{X} = \mathbf{D}\Gamma$$

Crude
approximation

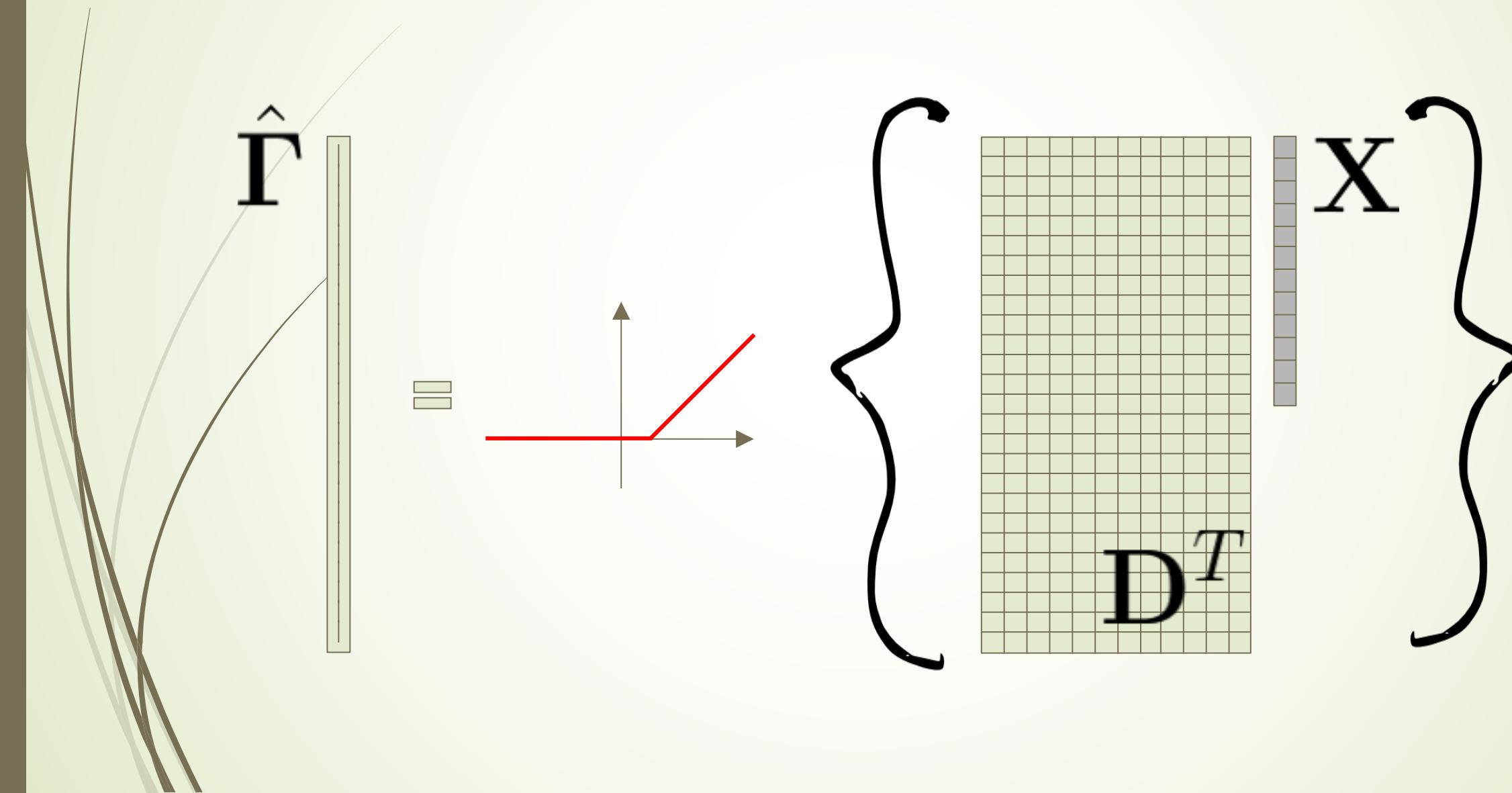
$$S_{\beta}\{\mathbf{D}^T \mathbf{X}\}$$

Soft-Thresholding Algorithm

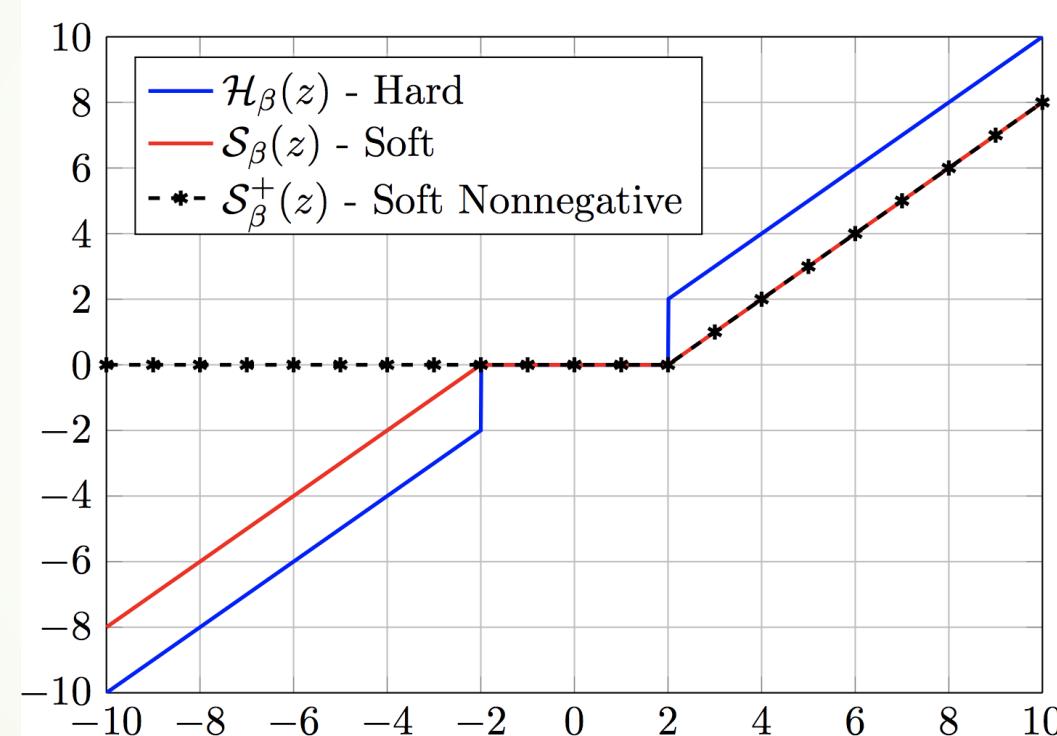
$\hat{\Gamma}$



First Layer of a Neural Network

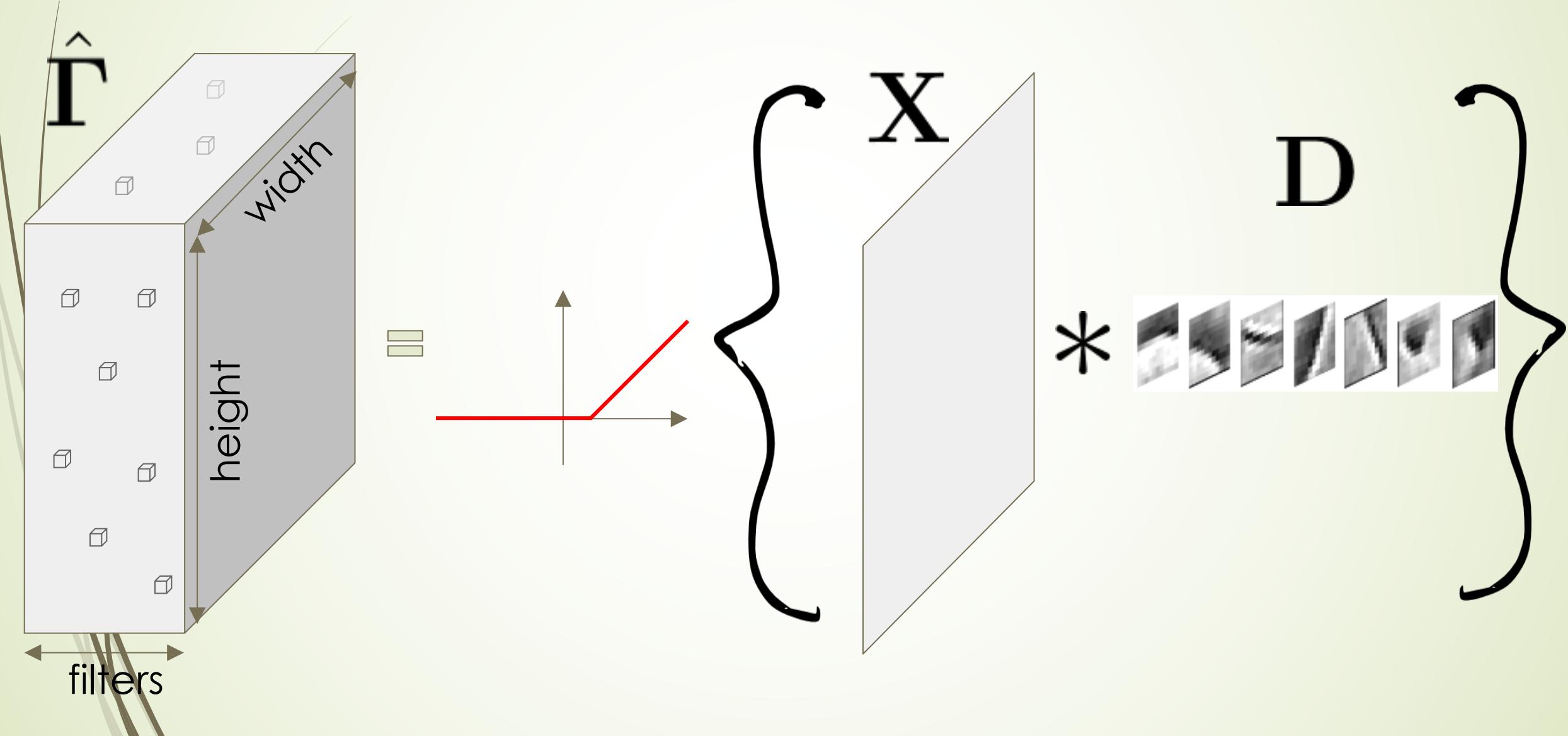


ReLU = Soft Nonnegative Thresholding



ReLU is equivalent to soft nonnegative thresholding

First layer of a **Convolutional** Neural Network

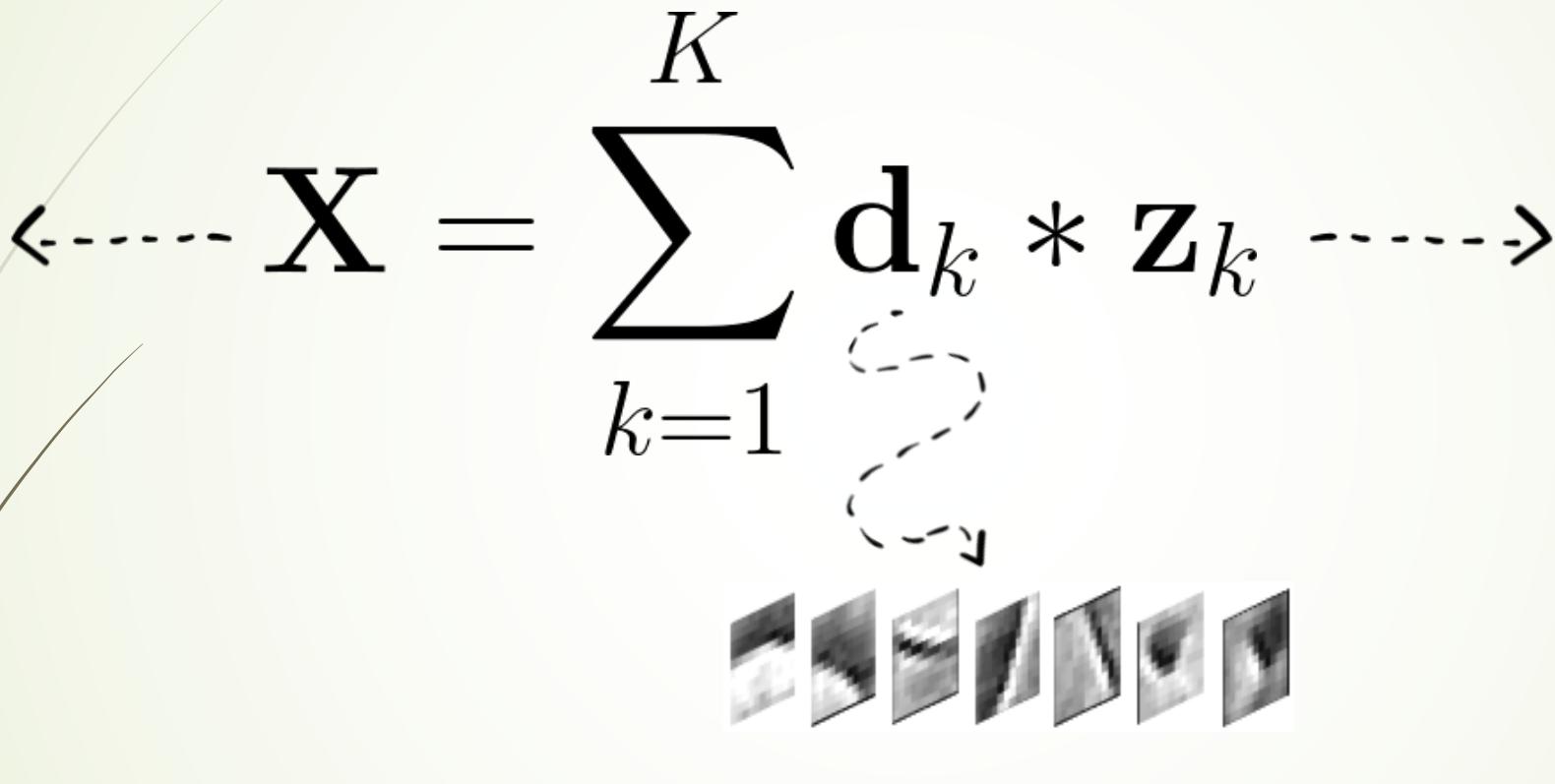


Convolutional Sparse Modeling

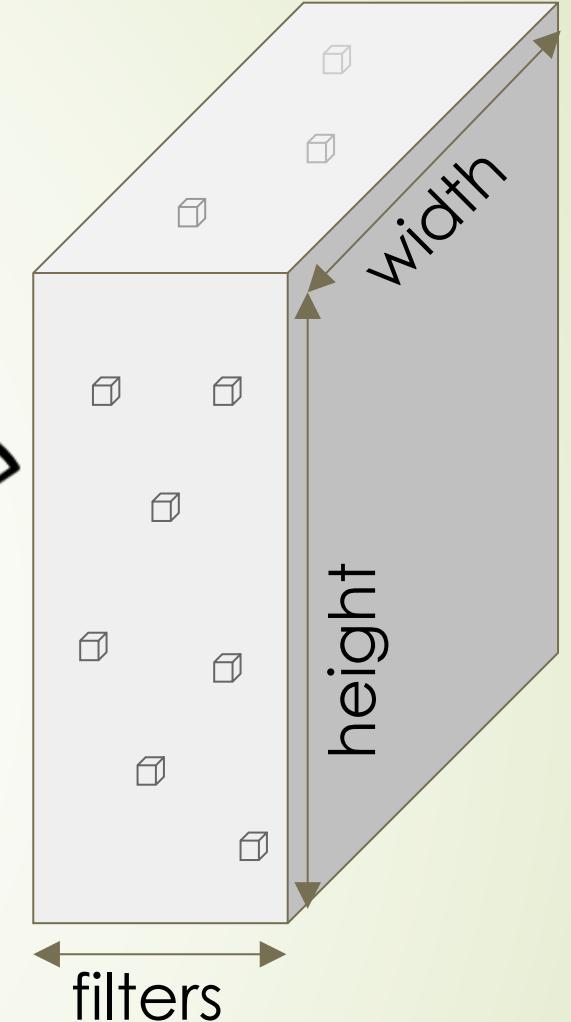
$$\mathbf{X} = \mathbf{D} \mathbf{T}$$

The diagram illustrates Convolutional Sparse Modeling. On the left, a vertical bar labeled \mathbf{X} contains several curved lines of varying colors (brown, tan, grey, black) representing the input signal. To the right of an equals sign is a large rectangular grid representing the convolutional dictionary \mathbf{D} . The grid has a repeating pattern of colored blocks (purple, blue, green) that form a diagonal step-like structure, representing the convolutional basis functions. To the right of the grid is a vertical bar labeled \mathbf{T} containing a series of vertical dots, representing the sparse coefficients or codebook entries.

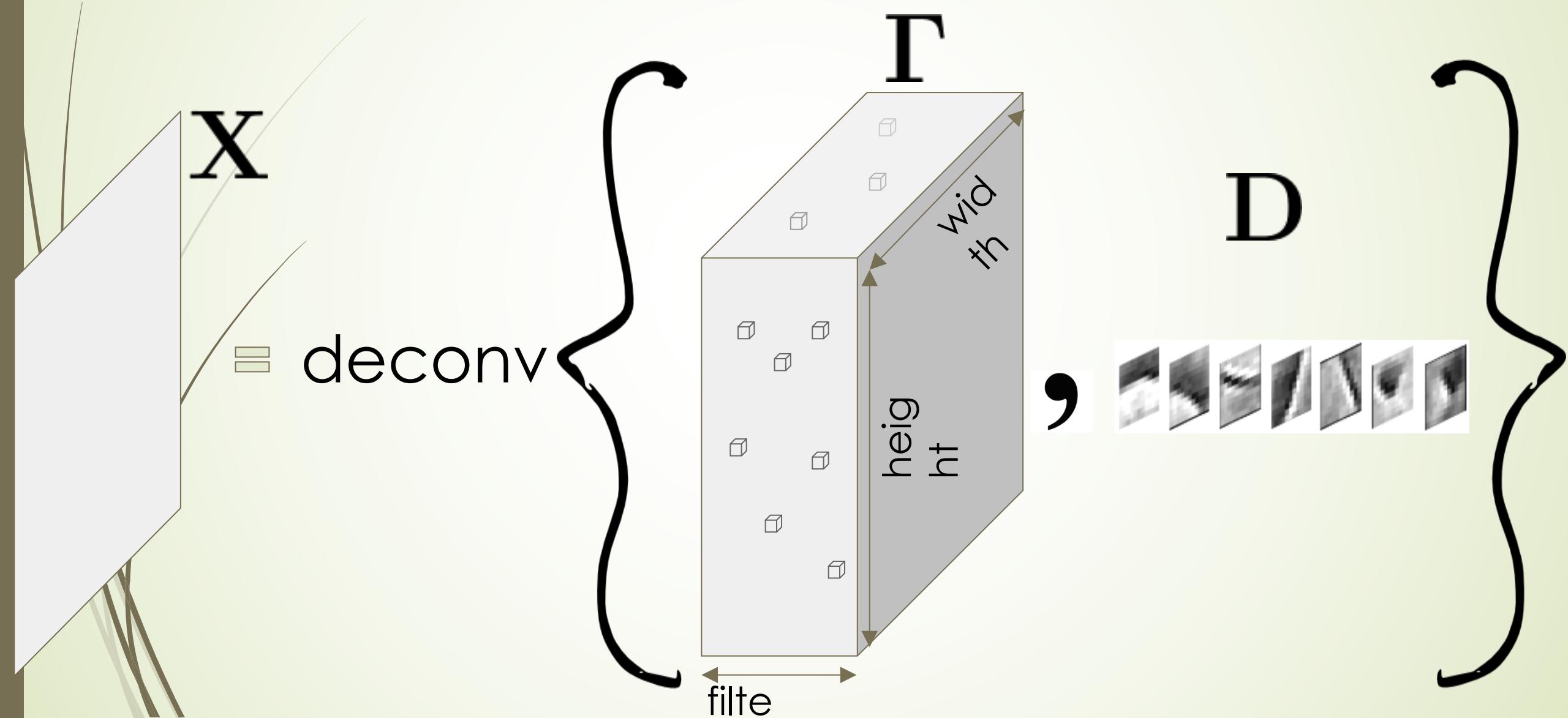
Convolutional Sparse Modeling

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


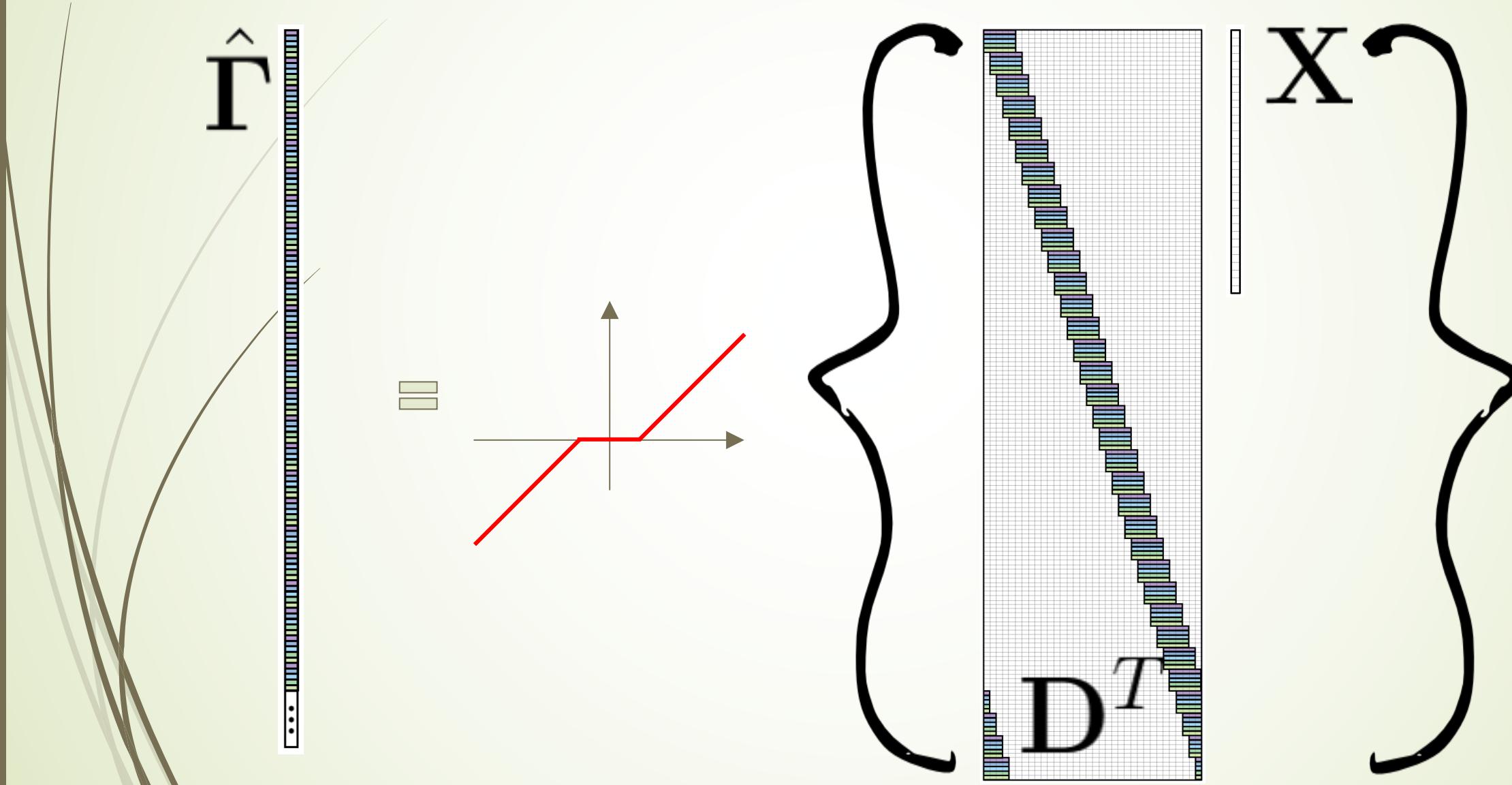
The diagram illustrates the convolutional sparse modeling equation. On the left, a sparse input vector X is shown as a stack of filters \mathbf{d}_k , each multiplied by a sparse coefficient \mathbf{z}_k . The filters are represented as small grayscale images, and the coefficients \mathbf{z}_k are represented as small white cubes. A dashed arrow points from the equation to the filters, and another dashed arrow points from the filters to the right side of the equation.



Convolutional Sparse Modeling

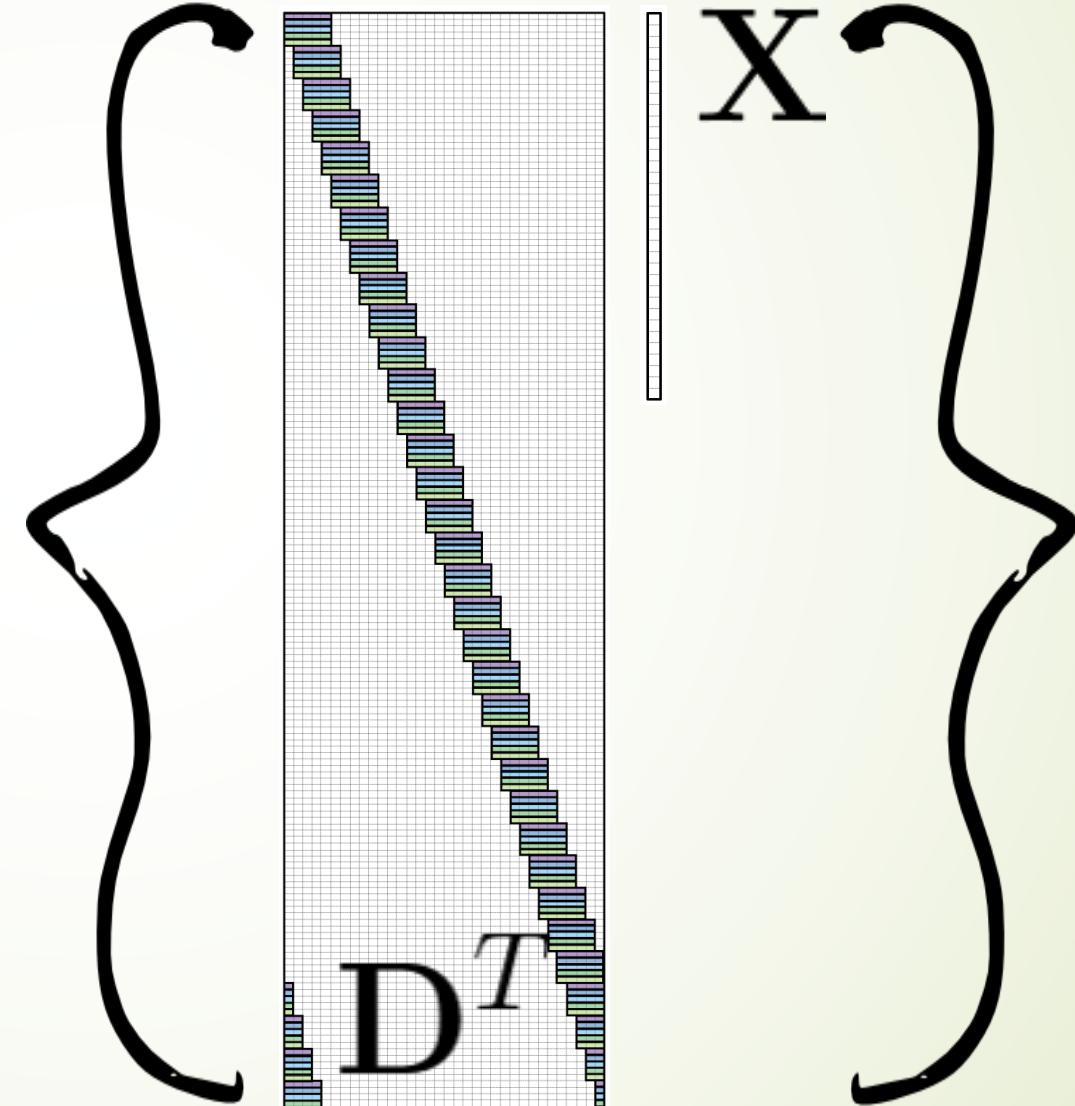
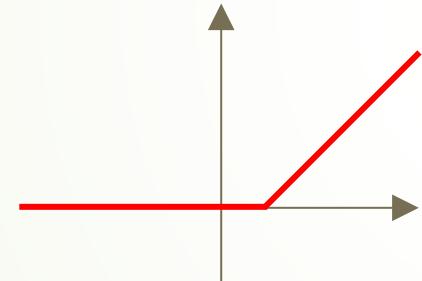


Thresholding Algorithm

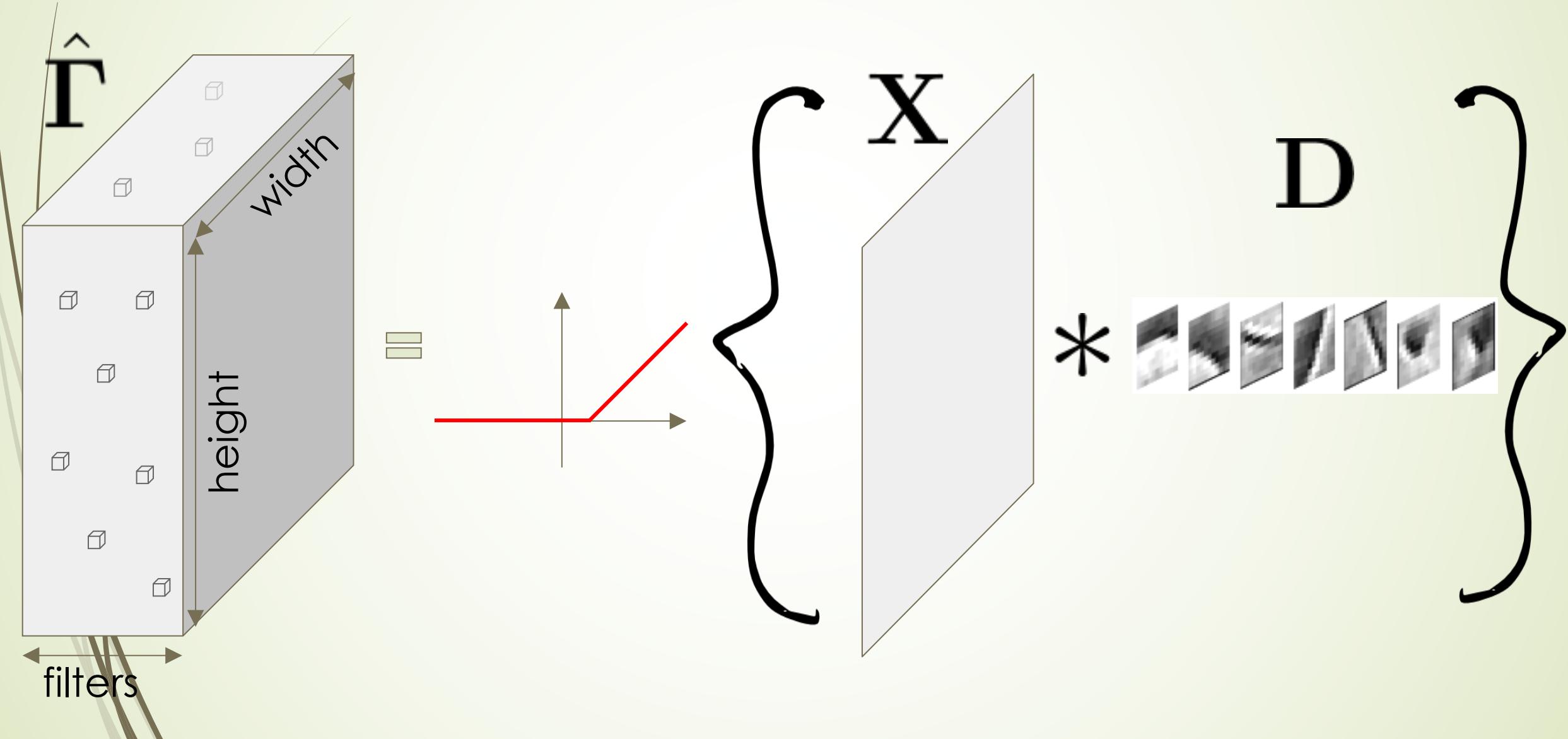


First layer of a Convolutional Neural Network

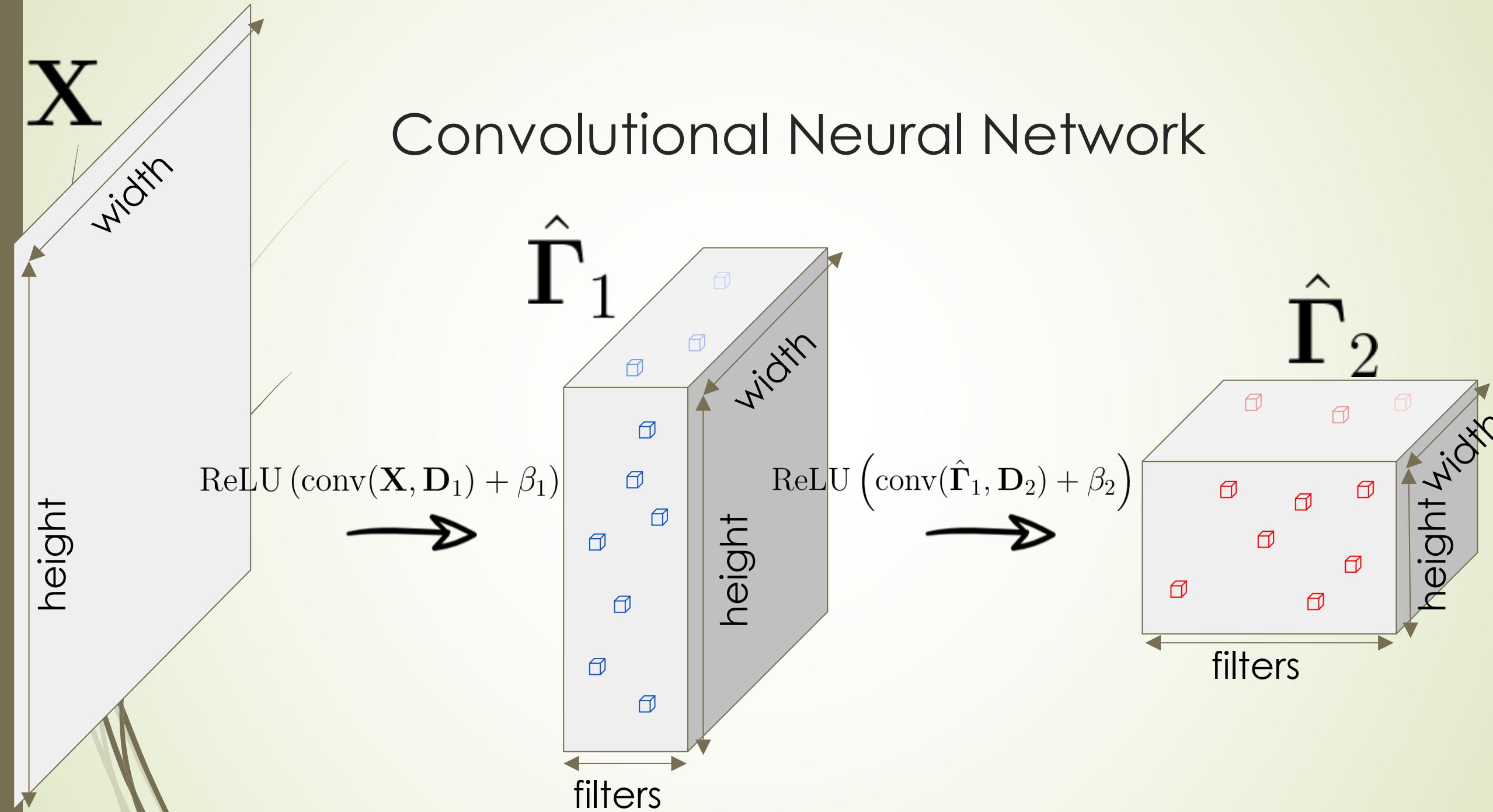
$\hat{\Gamma}$



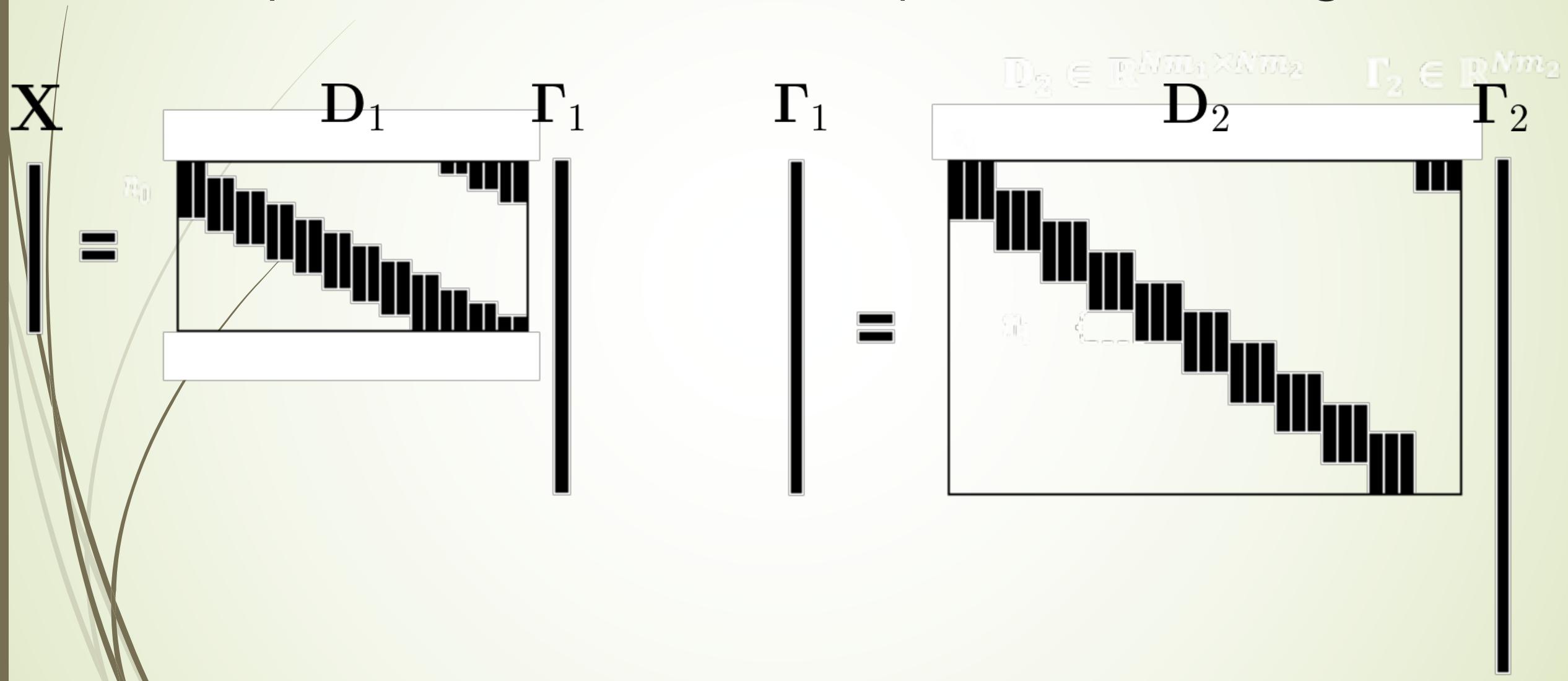
First layer of a Convolutional Neural Network



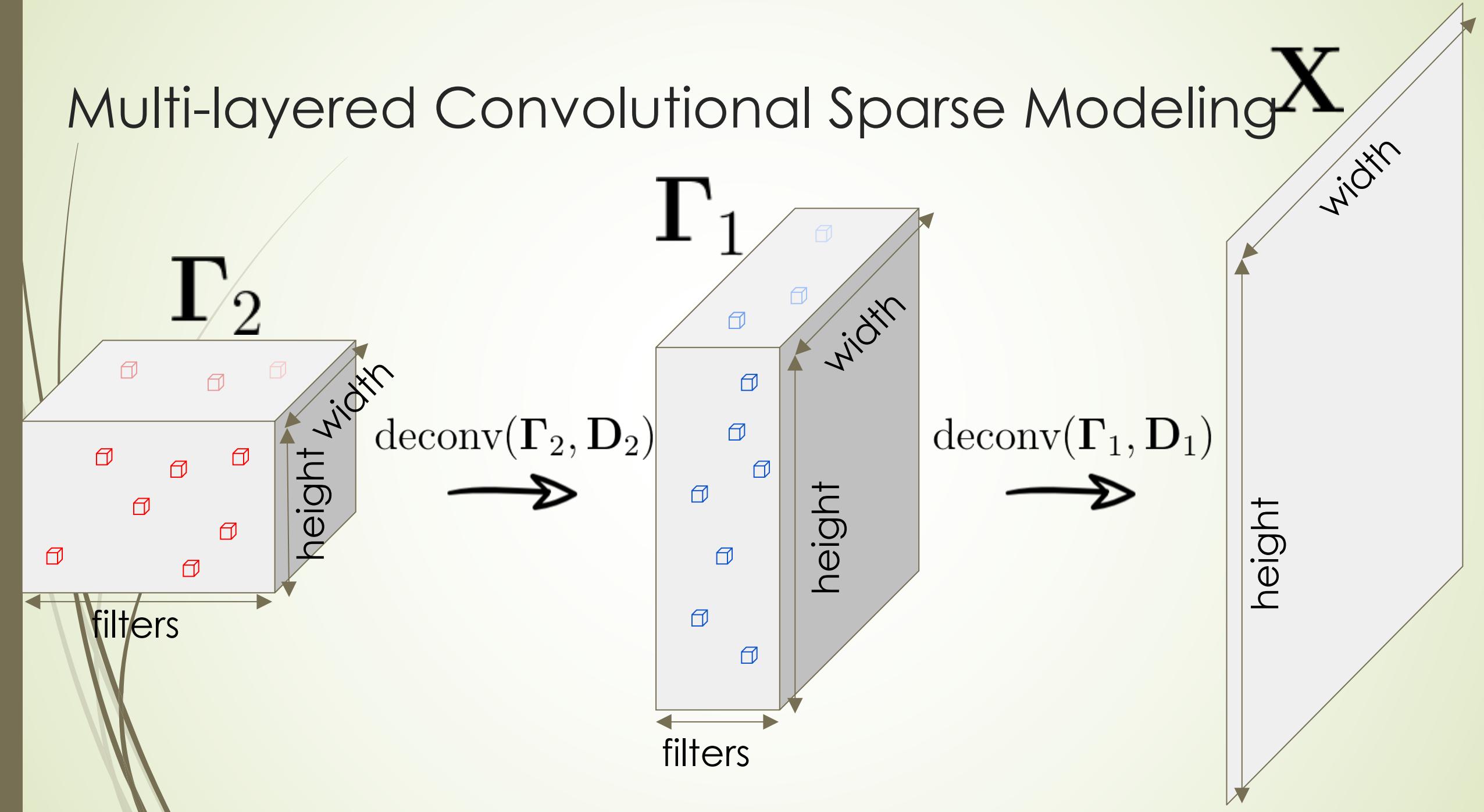
Convolutional Neural Network



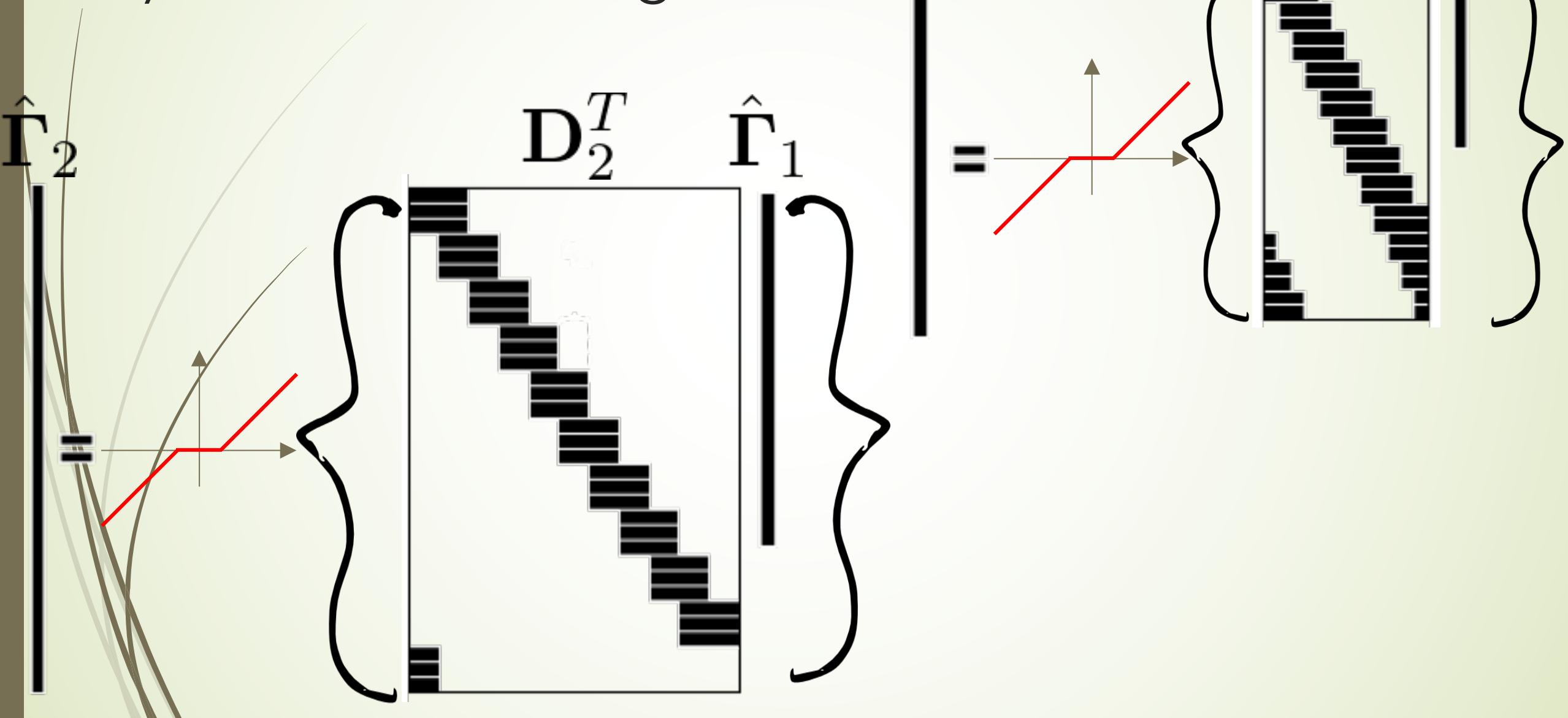
Multi-layered Convolutional Sparse Modeling



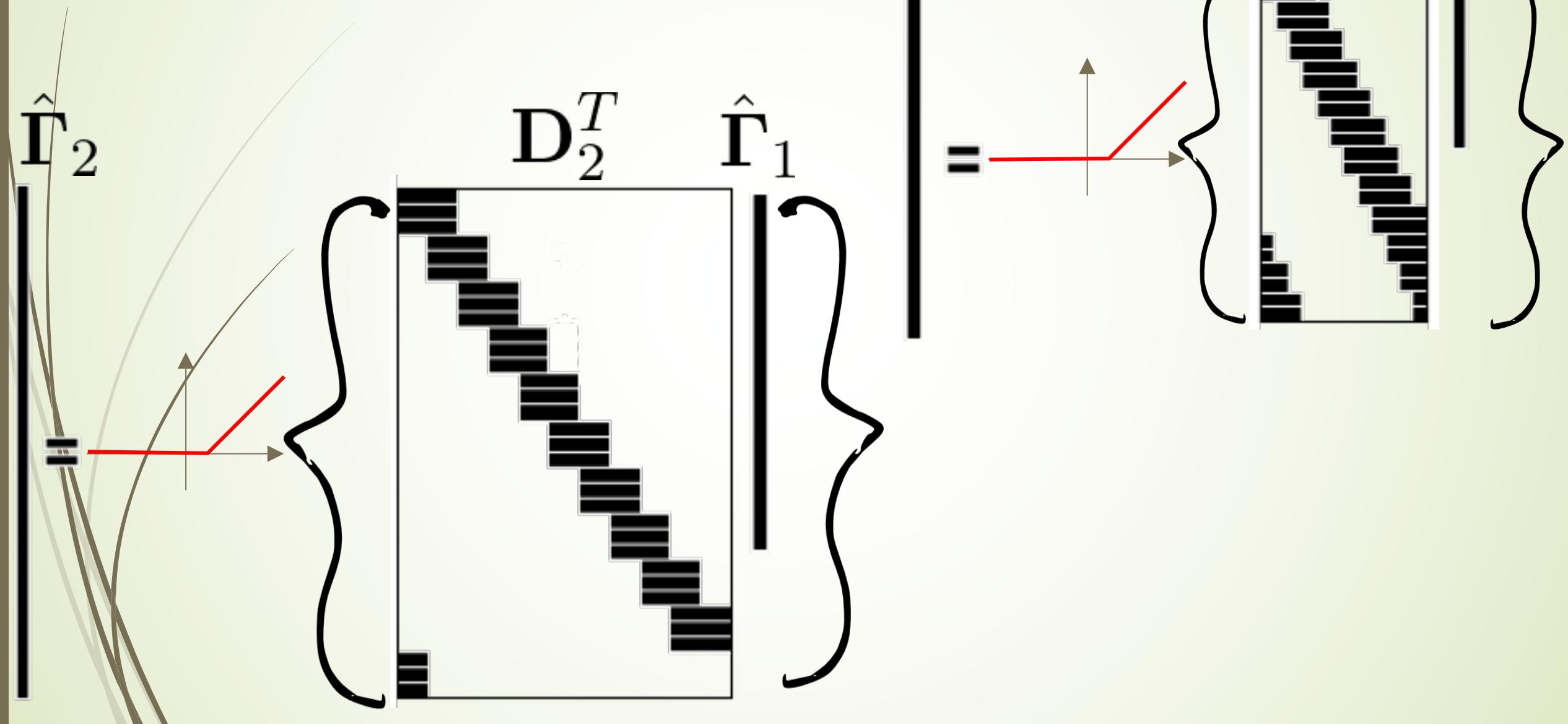
Multi-layered Convolutional Sparse Modeling



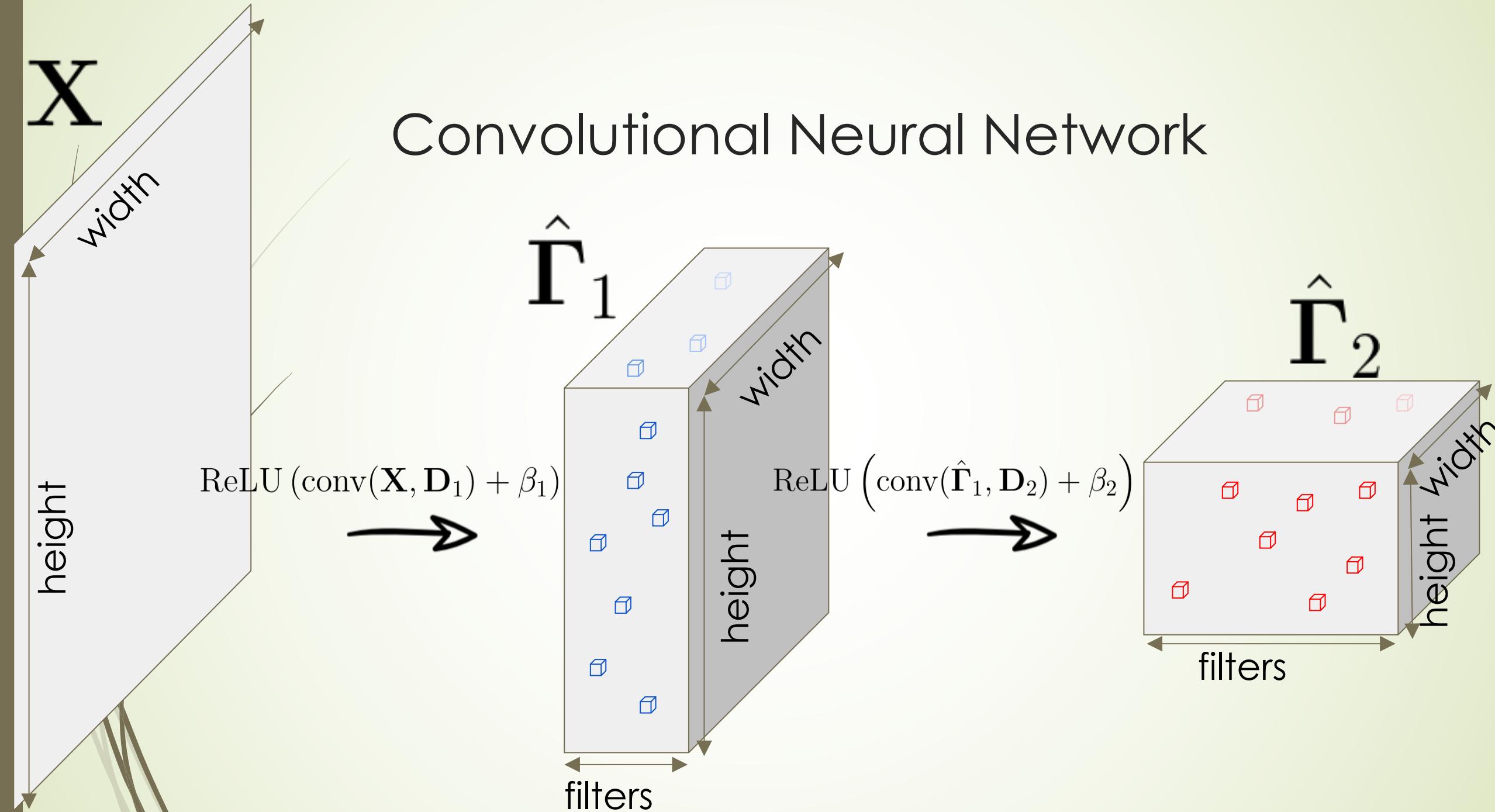
Layered Thresholding



Convolutional Neural Network



Convolutional Neural Network



Theories of Sparse Coding



Sparse Modeling

$$\mathbf{X} \quad \begin{matrix} \\ \end{matrix} = \quad \mathbf{D} \quad \Gamma$$

The diagram illustrates the sparse modeling equation $\mathbf{X} = \mathbf{D}\Gamma$. On the left, the matrix \mathbf{X} is shown as a vertical stack of gray rectangles. An equals sign ($=$) is positioned between \mathbf{X} and the product $\mathbf{D}\Gamma$. To the right of the equals sign is the matrix \mathbf{D} , which is a grid of light green squares. To the right of \mathbf{D} is the matrix Γ , represented as a vertical stack of light green and gray rectangles. The overall background features abstract, curved brown and gray lines.

Classic Sparse Theory: Basis Pursuit

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

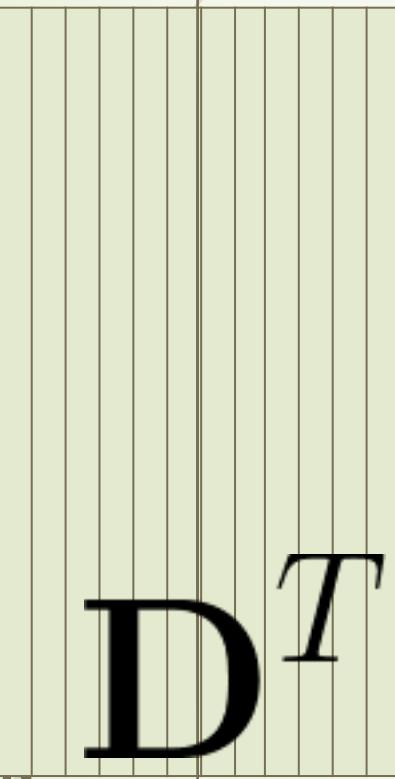
$$\hat{\boldsymbol{\Gamma}} = \arg \min_{\boldsymbol{\Gamma}} \|\boldsymbol{\Gamma}\|_1 \text{ s.t. } \mathbf{X} = \mathbf{D}\boldsymbol{\Gamma}$$

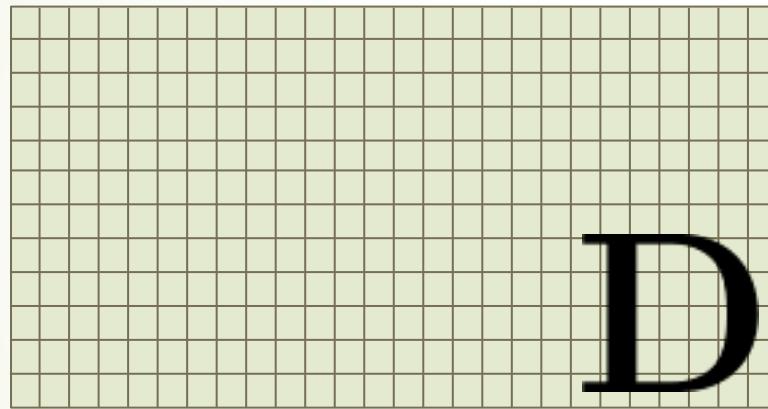
Theorem: [Donoho and Elad, 2003]

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

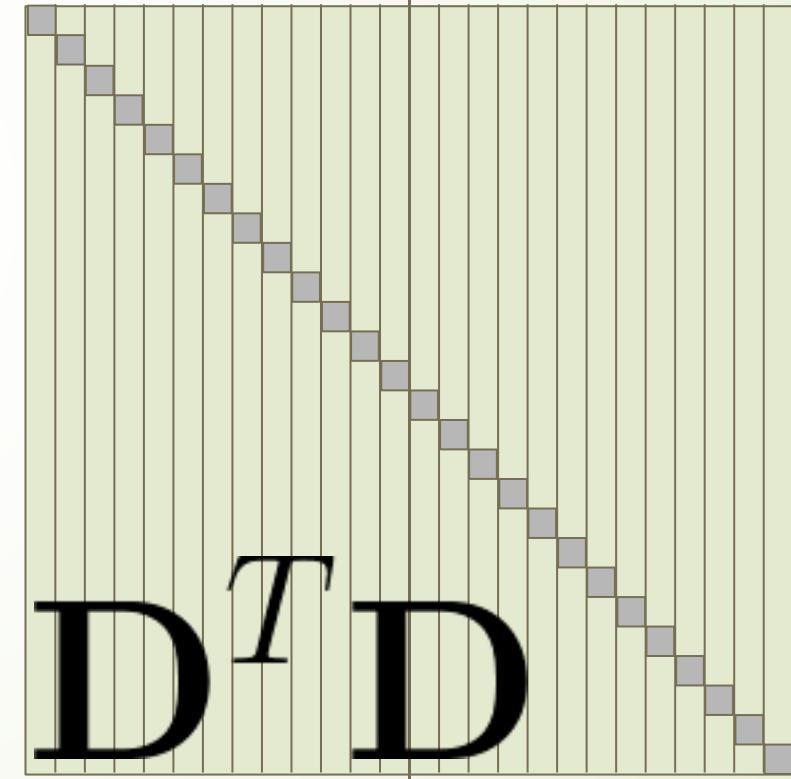
$$\|\boldsymbol{\Gamma}\|_0 < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})} \right)$$

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

$$\mathbf{D}^T$$


$$\mathbf{D}$$


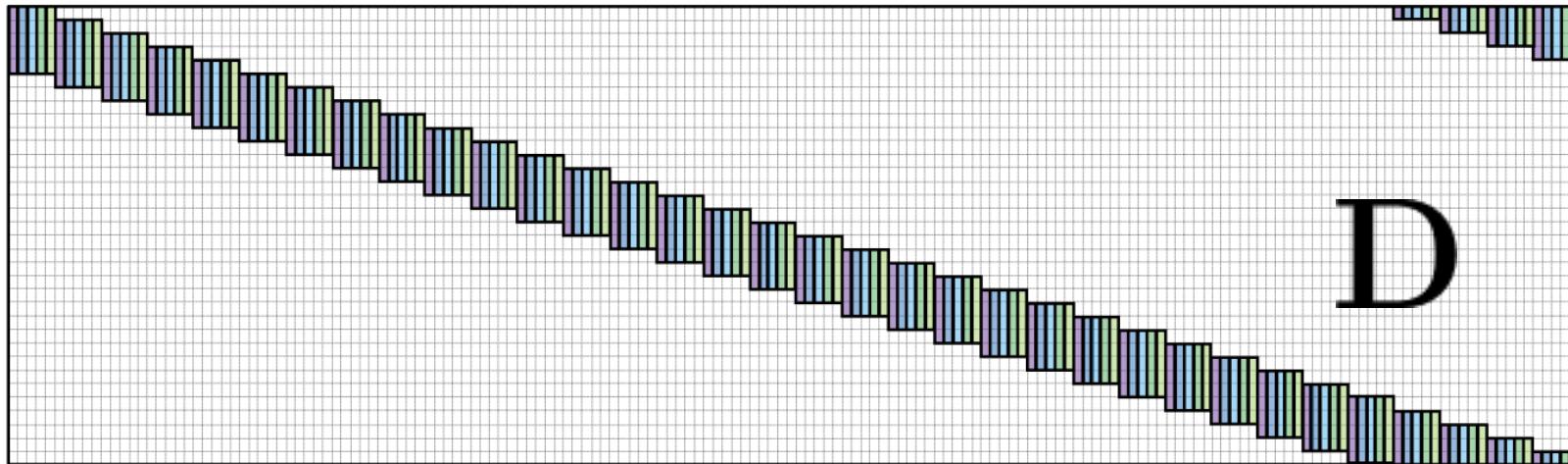
⋮

$$\mathbf{D}^T \mathbf{D}$$


Convolutional Sparse Modeling

X

=



D

Γ

Classic Sparse Theory for Convolutional Case

Theorem: [Donoho and Elad, 2003]

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

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

Assuming 2 atoms of length 64

$$\mu(\mathbf{D}) \geq 0.063 \quad [\text{Welch, 1974}]$$

Success guaranteed when

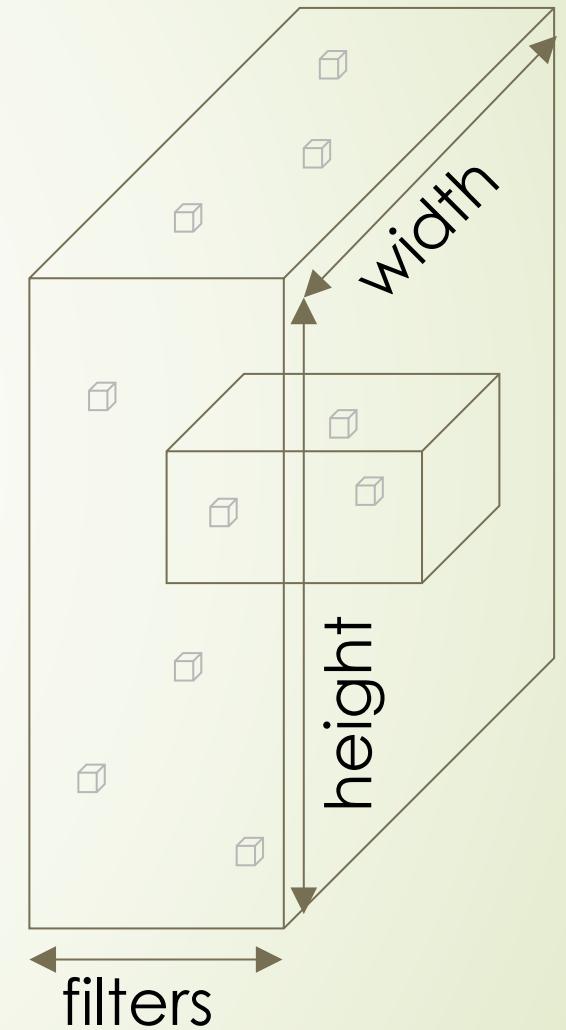
$$\|\Gamma\|_0 < 8.43$$



Local Sparsity

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

maximal number of non-zeroes
in a local neighborhood



Success of Basis Pursuit DeNoising

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

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

Theorem: [Papyan, Sulam and Elad, 2016]

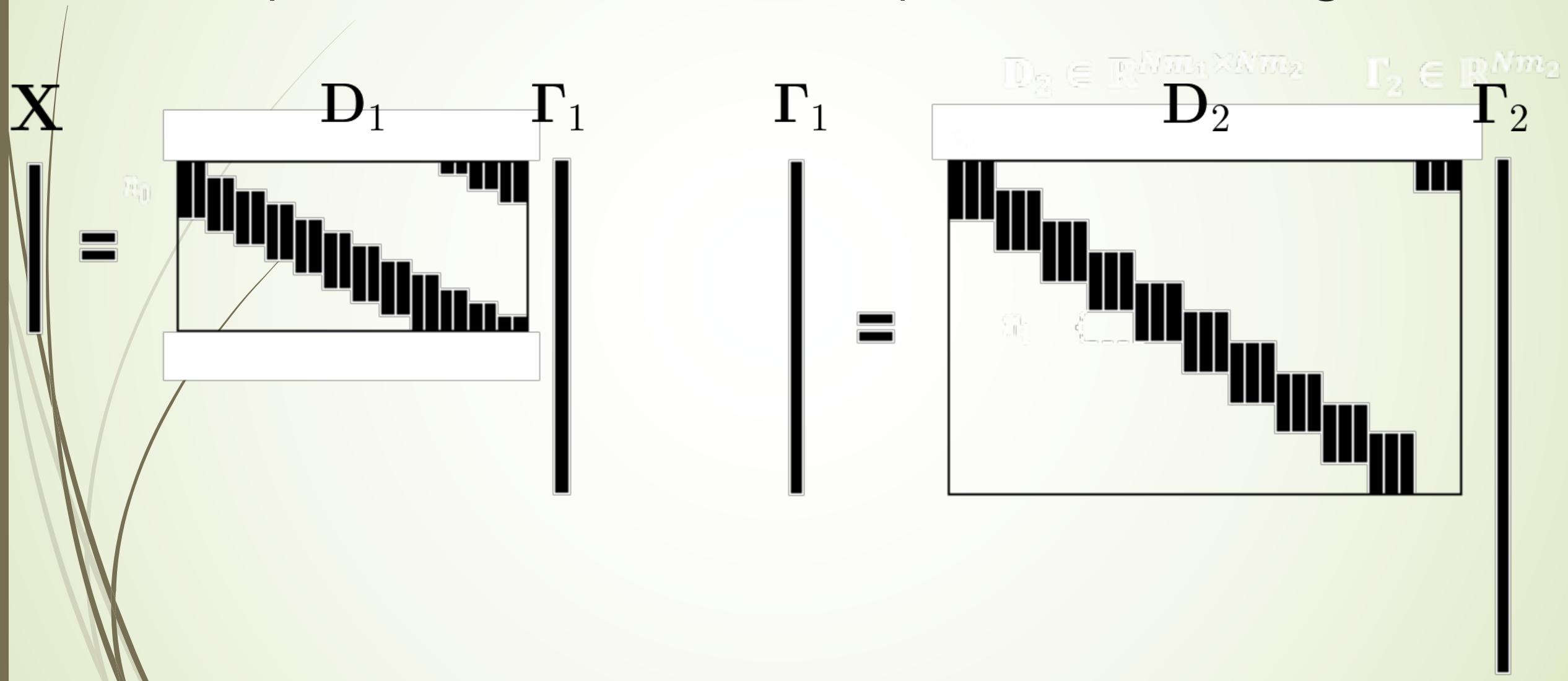
Assume: $\|\boldsymbol{\Gamma}\|_{0,\infty} < \frac{1}{3} \left(1 + \frac{1}{\mu(\mathbf{D})} \right)$

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

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:

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

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

⋮

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

Deep Coding Problem

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

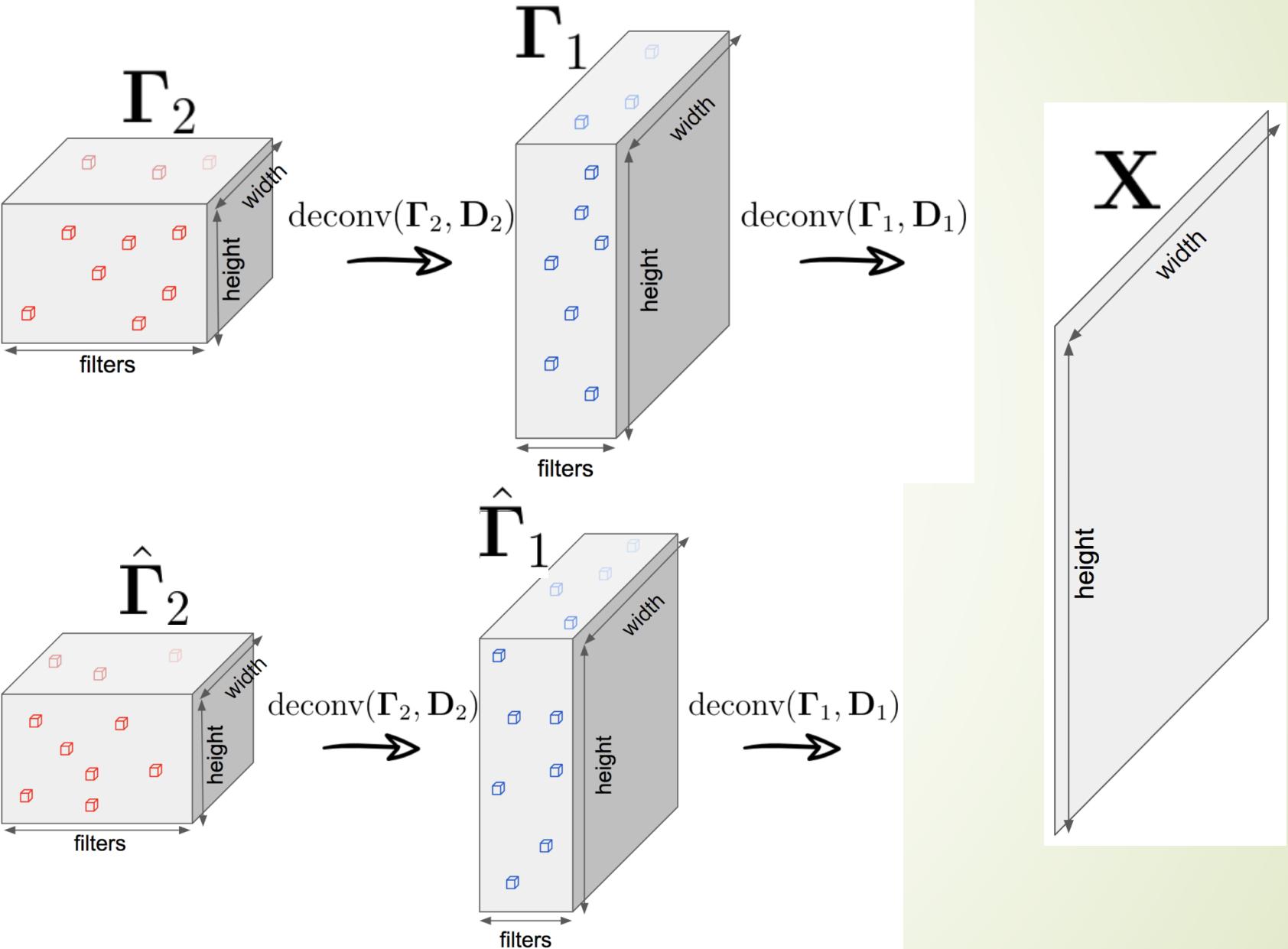
$$\|\mathbf{Y} - \mathbf{D}_1\boldsymbol{\Gamma}_1\|_2 \leq \epsilon, \quad \|\boldsymbol{\Gamma}_1\|_{0,\infty} \leq \lambda_1$$

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

⋮

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

Uniqueness



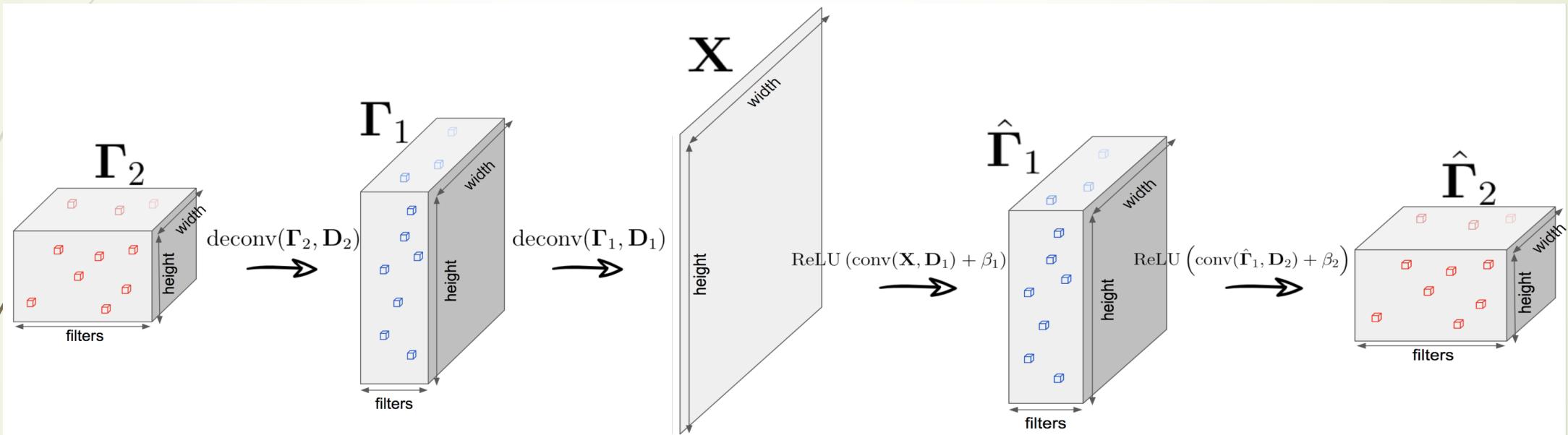
Uniqueness Theorem

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



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

Success of Forward Pass



Success of Forward Pass Theorem

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



Layered thresholding guaranteed:

1. Find correct places of nonzeros

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



Forward pass always fails at recovering representations exactly

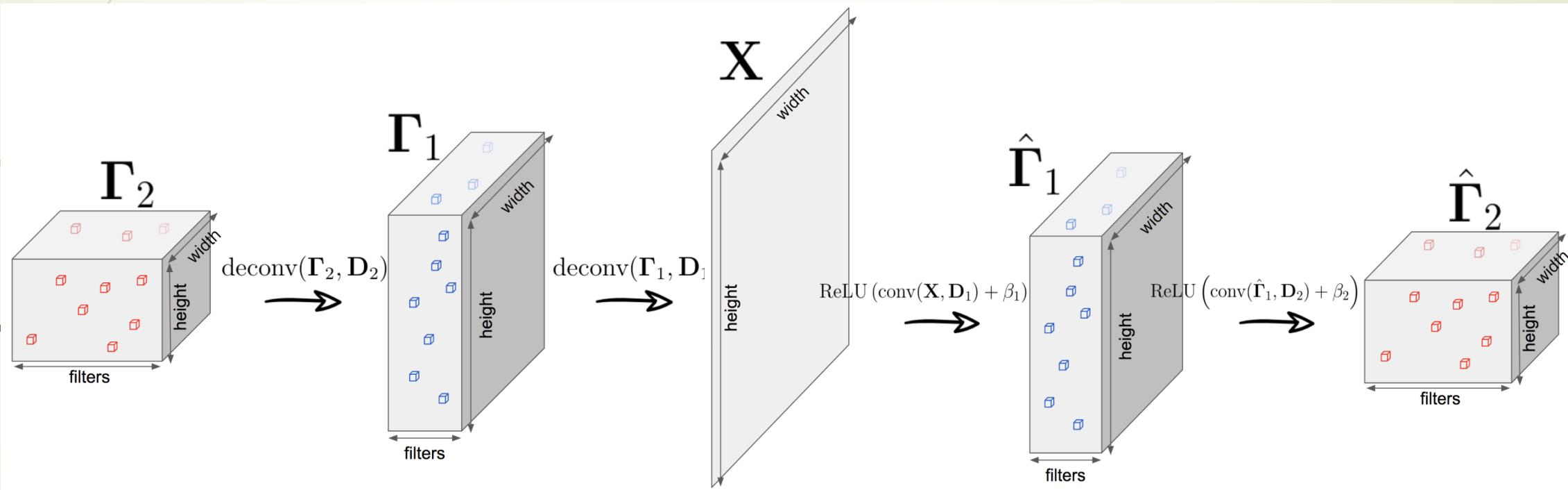


Success depends on ratio



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

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



Layered Lasso guaranteed:

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

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



Forward pass always fails at recovering representations exactly



Success depends on ratio



Distance increases with layer

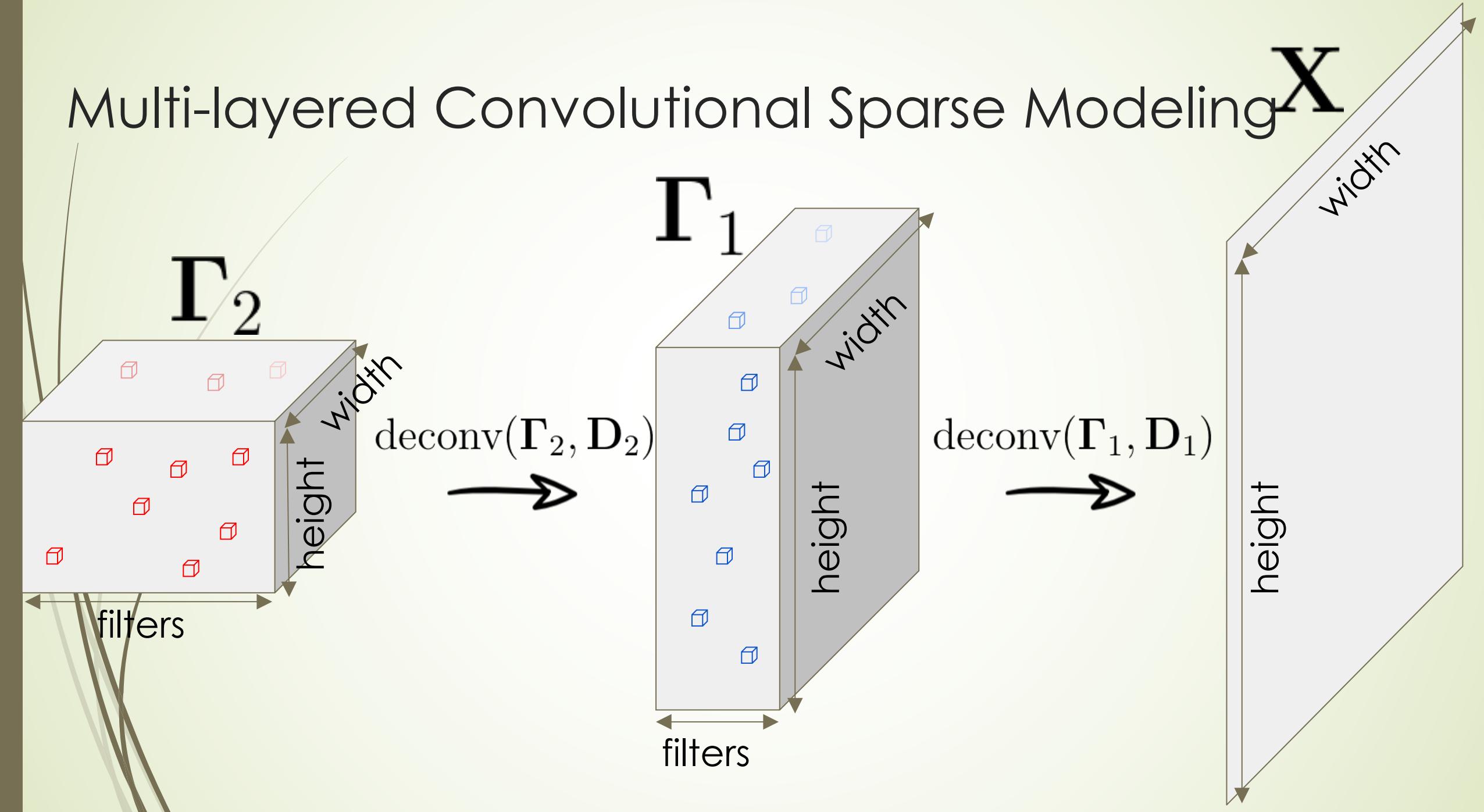
Layered Iterative Thresholding

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

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



Multi-layered Convolutional Sparse Modeling



Summary

1



2



3



4



5



Sparsity well established theoretically

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

Sparsity is the secret sauce behind CNN

Need to bring sparsity to the surface to
better understand CNNs

Andrej Karpathy agrees

Thank you!

