

# Sparse Modeling and Deep Learning

Michael Elad

Computer Science Department  
The Technion - Israel Institute of Technology  
Haifa 32000, Israel



The research leading to these results has been received funding  
from the European union's Seventh Framework Program  
(FP/2007-2013) ERC grant Agreement ERC-SPARSE- 320649

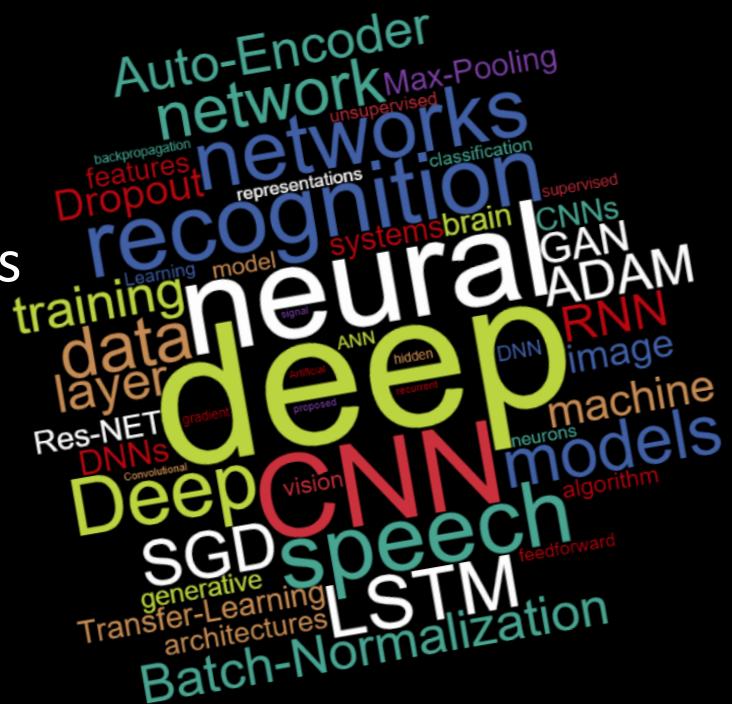


# This Lecture is About ...

# A Proposed Theory for Deep-Learning (DL)

## Explanation:

- DL has been extremely successful in solving a variety of learning problems
  - DL is an empirical field, with numerous tricks and know-how, but almost no theoretical foundations
  - A theory for DL has become the holy-grail of current research in Machine-Learning and related fields



# Who Needs Theory ?

We All Do !!

... because ... A theory

- ... could bring the next rounds of ideas to this field, breaking existing barriers and opening new opportunities
- ... could map clearly the limitations of existing DL solutions, and point to key features that control their performance
- ... could remove the feeling with many of us that DL is a “dark magic”, turning it into a solid scientific discipline

Ali Rahimi:  
NIPS 2017  
Test-of-Time  
Award



“Machine learning has become alchemy”



Yan LeCun



Understanding is a good thing ... but another goal is inventing methods. In the history of science and technology, engineering preceded theoretical understanding:

- Lens & telescope → Optics
- Steam engine → Thermodynamics
- Airplane → Aerodynamics
- Radio & Comm. → Info. Theory
- Computer → Computer Science



# A Theory for DL ?

Stephane Mallat (ENS) & Joan Bruna (NYU): Proposed the scattering transform (wavelet-based) and emphasized the treatment of invariances in the input data

Richard Baraniuk & Ankit Patel (RICE): Offered a generative probabilistic model for the data, showing how classic architectures and learning algorithms relate to it



Raja Giryes (TAU): Studied the architecture of DNN in the context of their ability to give distance-preserving embedding of signals

Gitta Kutyniok (TU) & Helmut Bolcskei (ETH): Studied the ability of DNN architectures to approximate families of functions

Data

Architecture

Algorithms

Rene Vidal (JHU): Explained the ability to optimize the typical non-convex objective and yet get to a global minima

Naftali Tishby (HUJI): Introduced the Information Bottleneck (IB) concept and demonstrated its relevance to deep learning

Stefano Soatto's team (UCLA): Analyzed the Stochastic Gradient Descent (SGD) algorithm, connecting it to the IB objective



# So, is there a Theory for DL ?



The answer is tricky:

There are already  
various such attempts,  
and some of them are  
truly impressive

... but ...

none of them is  
complete



# Interesting Observations

- Theory origins: Signal Proc., Control Theory, Info. Theory, Harmonic Analysis, Sparse Represen., Quantum Physics, PDE, Machine learning ...



Ron Kimmel: "*DL is a dark monster covered with mirrors. Everyone sees his reflection in it ...*"



David Donoho: "... these mirrors are taken from *Cinderella's story*, telling each that he is the **most beautiful**"

- Today's talk is on our proposed theory:



Yaniv Romano



Vardan Papyan



Jeremias Sulam



Aviad Aberdam



Data



Architecture

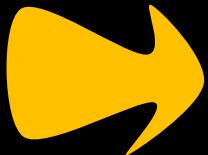
Algorithms

... and our theory is the best 😊



# Our Story: More Specifically

**Sparseland**  
Sparse  
Representation  
Theory



**CSC**  
Convolutional  
Sparse  
Coding



**ML-CSC**  
Multi-Layered  
Convolutional  
Sparse Coding

Sparsity-Inspired Models



Deep-Learning

- In this talk we shall start with a brief overview of the first two models, and then step directly to the ML-CSC model and its connection to deep-learning
- If you feel that you are missing key information, you can complement this by viewing my YouTube IPAM talk from February 2018



# Brief Background on Sparse Modeling



# Our Data is Structured

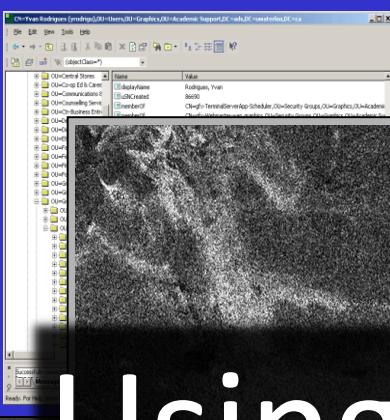
Stock Market



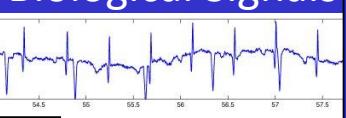
Videos



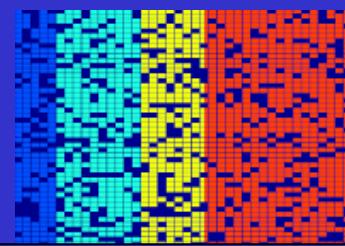
Text Documents



Biological Signals



Matrix Data



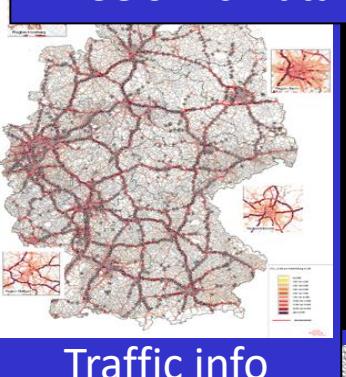
Still Images



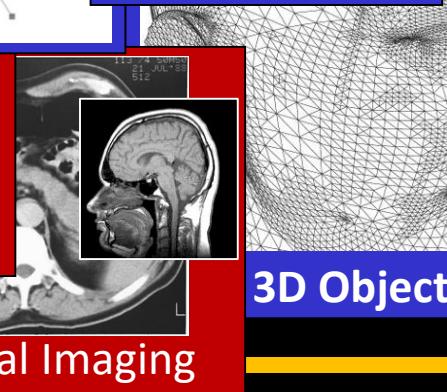
Social Networks



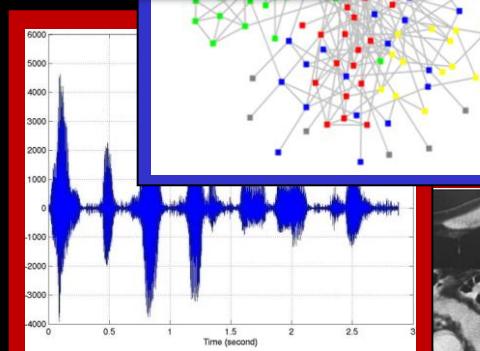
Seismic Data



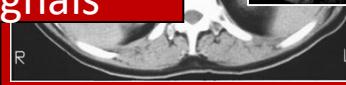
Traffic info



3D Objects



Voice Signals



Medical Imaging

# Using models

- We are surrounded by various diverse sources of massive information
- Each of these sources have an internal structure, which can be exploited
- This structure, when identified, is the engine behind the ability to process data
- How to identify structure?



# What this Talk is all About?

## Data Models and Their Use

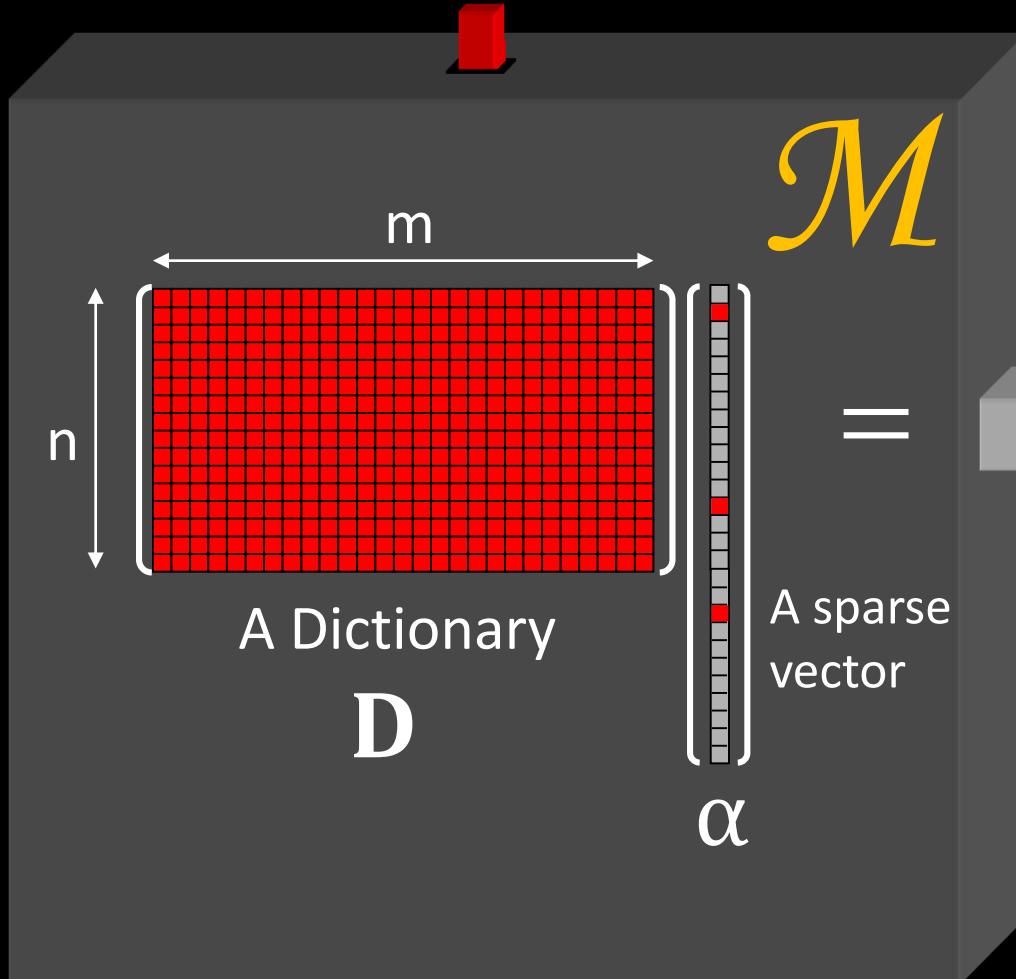
- Almost any task in data processing requires a model – true for denoising, deblurring, super-resolution, inpainting, compression, anomaly-detection, sampling, recognition, separation, and more
- Sparse and Redundant Representations offer a new and highly effective model – we call it

*Sparseland*

- We shall describe this and descendant versions of it that lead all the way to ... **deep-learning**



# *Sparseland*: A Formal Description



- Every column in  $D$  (**dictionary**) is a prototype signal (atom)
- The vector  $\underline{\alpha}$  is generated with few non-zeros at arbitrary locations and values
- This is a generative model that describes how (**we believe**) signals are created

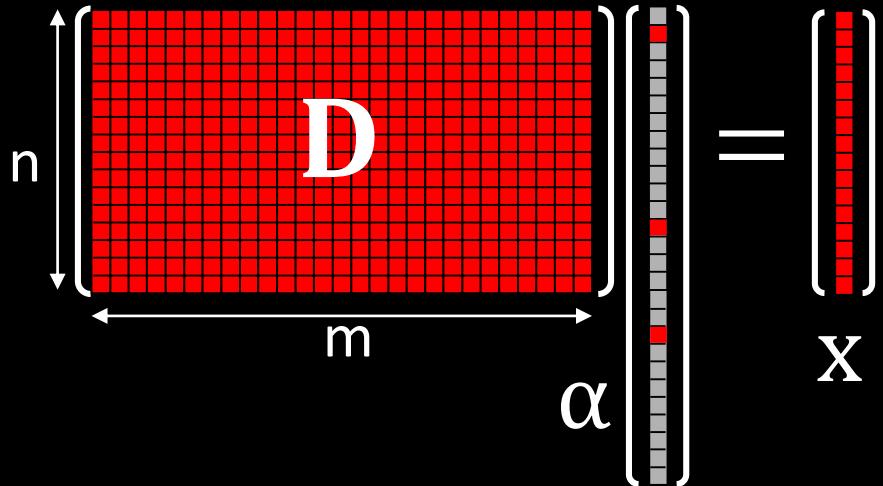


# Atom Decomposition

$$\min_{\alpha} \|\alpha\|_0 \text{ s.t. } x = D\alpha$$



$$\min_{\alpha} \|\alpha\|_0 \text{ s.t. } \|D\alpha - y\|_2 \leq \varepsilon$$



Approximation Algorithms



Relaxation methods

Basis-Pursuit



Greedy methods

Thresholding/OMP

- $\ell_0$  – counting number of non-zeros in the vector
- This is a projection onto the *Sparseland* model
- These problems are known to be NP-Hard problem



# Pursuit Algorithms

$$\min_{\alpha} \|\alpha\|_0 \text{ s. t. } \|\mathbf{D}\alpha - \mathbf{y}\|_2 \leq \varepsilon$$

Approximation Algorithms

Basis Pursuit

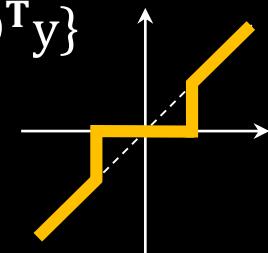
Thresholding

Change the  $L_0$  into  $L_1$  and then the problem becomes convex and manageable

Multiply  $\mathbf{y}$  by  $\mathbf{D}^T$  and apply shrinkage:

$$\min_{\alpha} \|\alpha\|_1 \text{ s. t. } \|\mathbf{D}\alpha - \mathbf{y}\|_2 \leq \varepsilon$$

$$\hat{\alpha} = \mathcal{P}_{\beta}\{\mathbf{D}^T \mathbf{y}\}$$



# The Mutual Coherence

- Compute

$$\begin{bmatrix} \mathbf{D}^T \\ \mathbf{D} \end{bmatrix} \begin{bmatrix} & & \\ & \mathbf{D} & \\ & & \end{bmatrix} = \begin{bmatrix} & & \\ & \mathbf{D}^T \mathbf{D} & \\ & & \end{bmatrix}$$

Assume  
normalized  
columns

- The **Mutual Coherence**  $\mu(\mathbf{D})$  is the largest off-diagonal entry in absolute value
- We will pose the theoretical results in this talk using this property, due to its simplicity



# Basis-Pursuit Success

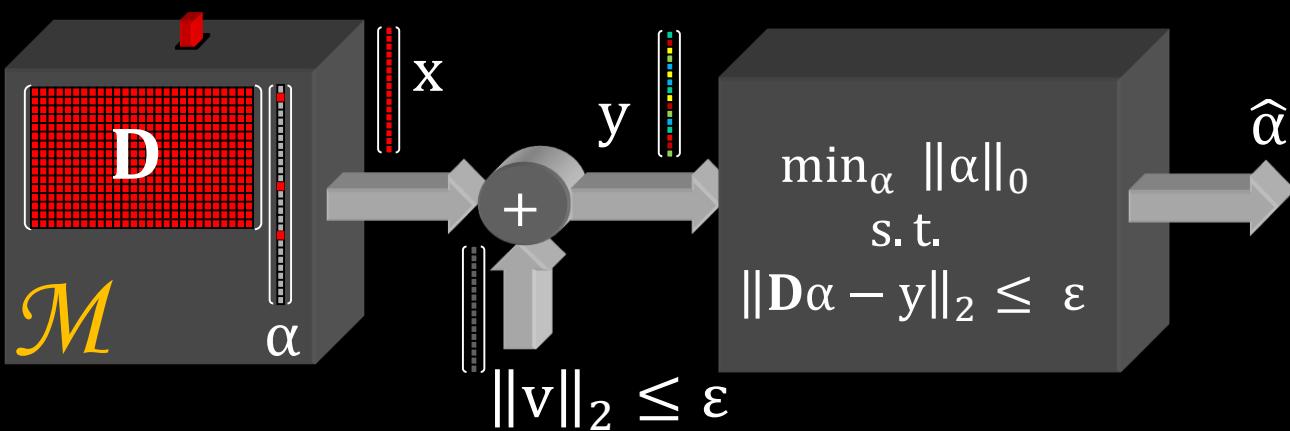


**Theorem:** Given a noisy signal  $y = D\alpha + v$  where  $\|v\|_2 \leq \varepsilon$  and  $\alpha$  is sufficiently sparse,  $\|\alpha\|_0 < \frac{1}{4} \left(1 + \frac{1}{\mu}\right)$

then **Basis-Pursuit**:  $\min_{\alpha} \|\alpha\|_1$  s. t.  $\|D\alpha - y\|_2 \leq \varepsilon$

leads to a stable result:  $\|\hat{\alpha} - \alpha\|_2^2 \leq \frac{4\varepsilon^2}{1-\mu(4\|\alpha\|_0-1)}$

Donoho, Elad & Temlyakov ('06)

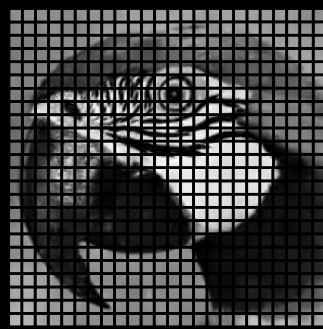


Comments:

- If  $\varepsilon=0 \rightarrow \hat{\alpha} = \alpha$
- This is a worst-case analysis – better bounds exist
- Similar theorems exist for many other pursuit algorithms



# Convolutional Sparse Coding (CSC)



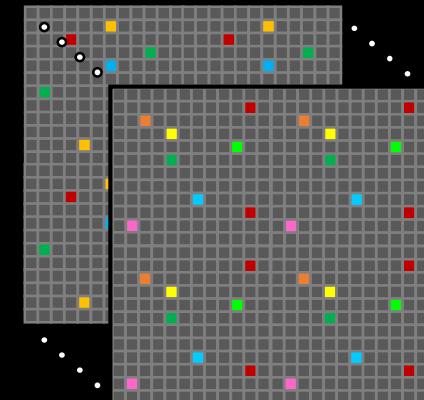
An image with  $N$  pixels

$m$  filters convolved with their sparse representations

$$[\mathbf{X}] = \sum_{i=1}^m d_i * [\Gamma_i]$$

The  $i$ -th filter of small size  $n$

$i$ -th feature-map: An image of the same size as  $\mathbf{X}$  holding the sparse representation related to the  $i$ -filter



This model emerged in 2005-2010, developed and advocated by Yan LeCun and others. It serves as the foundation of Convolutional Neural Networks

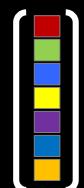


# CSC in Matrix Form

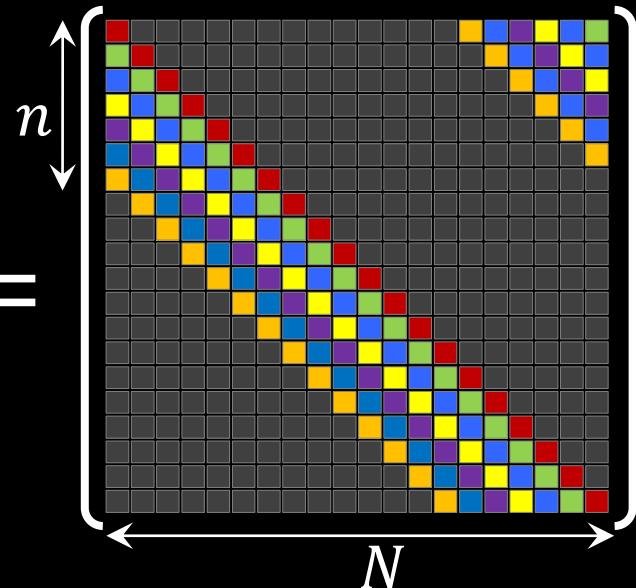
- Here is an alternative global sparsity-based model formulation

$$\mathbf{X} = \sum_{i=1}^m \mathbf{C}^i \boldsymbol{\Gamma}^i = [\mathbf{C}^1 \dots \mathbf{C}^m] \begin{bmatrix} \boldsymbol{\Gamma}^1 \\ \vdots \\ \boldsymbol{\Gamma}^m \end{bmatrix} = \mathbf{D}\boldsymbol{\Gamma}$$

- $\mathbf{C}^i \in \mathbb{R}^{N \times N}$  is a banded and Circulant matrix containing a single atom with all of its shifts



$\mathbf{C}^i =$



- $\boldsymbol{\Gamma}^i \in \mathbb{R}^N$  are the corresponding coefficients ordered as column vectors



# The CSC Dictionary

$$[\mathbf{C}^1 \ \mathbf{C}^2 \ \mathbf{C}^3] = \left[ \begin{array}{c|c|c} \text{Colorful sparse matrix} & \text{Colorful sparse matrix} & \text{Colorful sparse matrix} \\ \text{Colorful sparse matrix} & \text{Colorful sparse matrix} & \text{Colorful sparse matrix} \\ \text{Colorful sparse matrix} & \text{Colorful sparse matrix} & \text{Colorful sparse matrix} \\ \text{Colorful sparse matrix} & \text{Colorful sparse matrix} & \text{Colorful sparse matrix} \\ \text{Colorful sparse matrix} & \text{Colorful sparse matrix} & \text{Colorful sparse matrix} \end{array} \right]$$
  
$$\mathbf{D}_L \leftarrow \mathbf{D} = \left[ \begin{array}{c|c|c} \text{Red sparse matrix} & \text{Red sparse matrix} & \text{Red sparse matrix} \\ \text{Red sparse matrix} & \text{Red sparse matrix} & \text{Red sparse matrix} \\ \text{Red sparse matrix} & \text{Red sparse matrix} & \text{Red sparse matrix} \\ \text{Red sparse matrix} & \text{Red sparse matrix} & \text{Red sparse matrix} \\ \text{Red sparse matrix} & \text{Red sparse matrix} & \text{Red sparse matrix} \end{array} \right]$$

Diagram illustrating the CSC (Compressed Sparse Column) dictionary representation. The top part shows three sparse matrices  $\mathbf{C}^1, \mathbf{C}^2, \mathbf{C}^3$  as colored blocks. The bottom part shows a full matrix  $\mathbf{D}$  with a red pattern, where  $\mathbf{D}_L$  is a submatrix extracted from it. A red box highlights a  $m \times n$  block in the matrix  $\mathbf{D}$ .



# Multi-Layered Convolutional Sparse Modeling



Yaniv Romano



Vardan Papyan



Jeremias Sulam

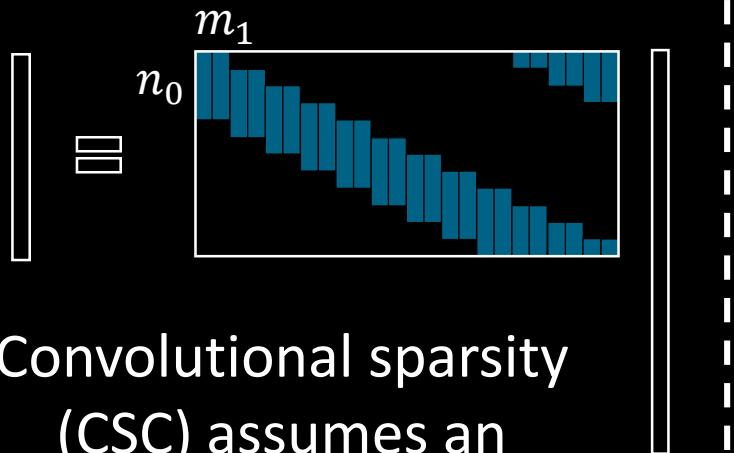


Aviad Aberdam



# From CSC to Multi-Layered CSC

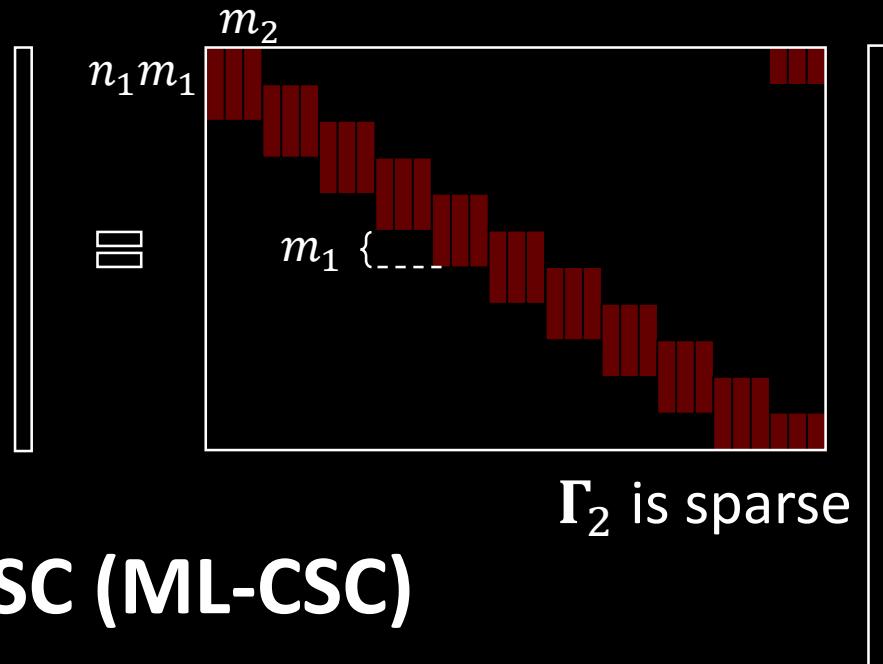
$$\mathbf{X} \in \mathbb{R}^N \quad \mathbf{D}_1 \in \mathbb{R}^{N \times Nm_1} \quad \boldsymbol{\Gamma}_1 \in \mathbb{R}^{Nm_1}$$



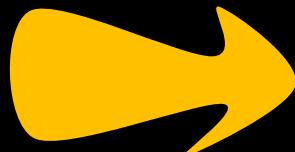
Convolutional sparsity  
(CSC) assumes an  
inherent structure is  
present in natural  
signals -  $\boldsymbol{\Gamma}_1$  is sparse

We propose to impose the  
same structure on the  
representations **themselves**

$$\boldsymbol{\Gamma}_1 \in \mathbb{R}^{Nm_1} \quad \mathbf{D}_2 \in \mathbb{R}^{Nm_1 \times Nm_2} \quad \boldsymbol{\Gamma}_2 \in \mathbb{R}^{Nm_2}$$



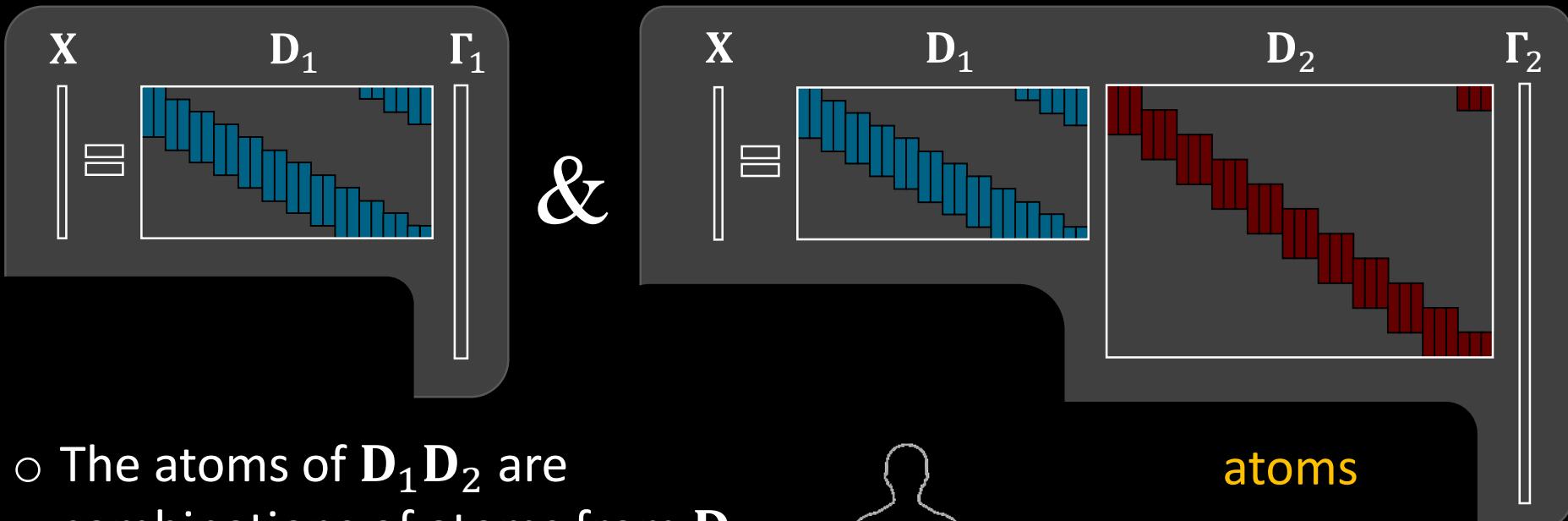
$\boldsymbol{\Gamma}_2$  is sparse



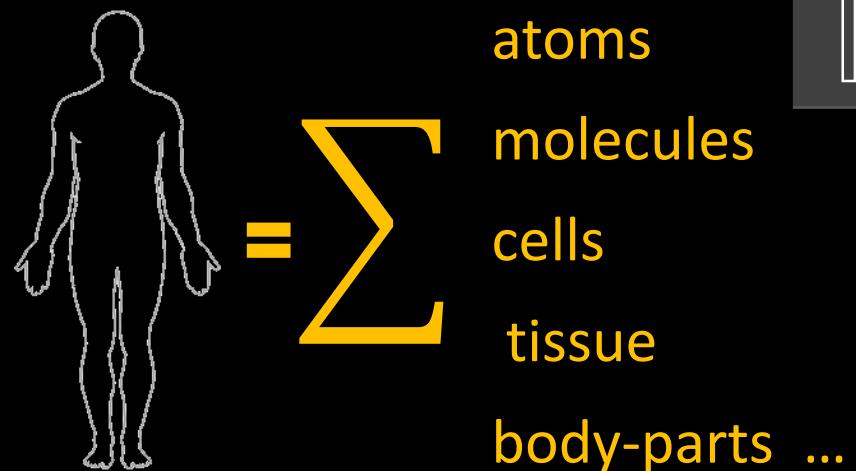
**Multi-Layer CSC (ML-CSC)**



# Intuition: From Atoms to Molecules



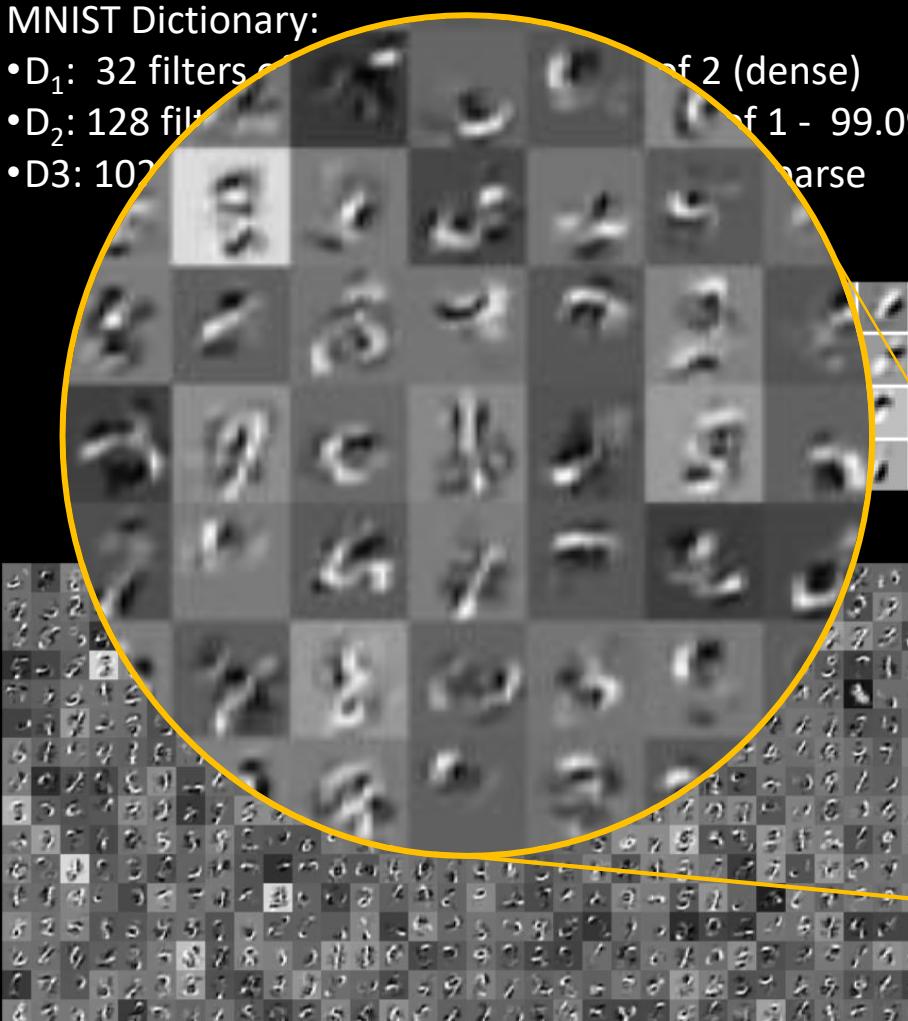
- The atoms of  $D_1 D_2$  are combinations of atoms from  $D_1$  - these are now **molecules**
- Thus, this model offers different **levels of abstraction** in describing  $X$



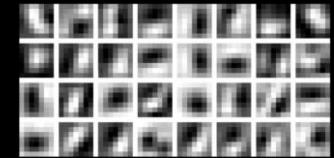
# A Small Taste: Model Training (MNIST)

MNIST Dictionary:

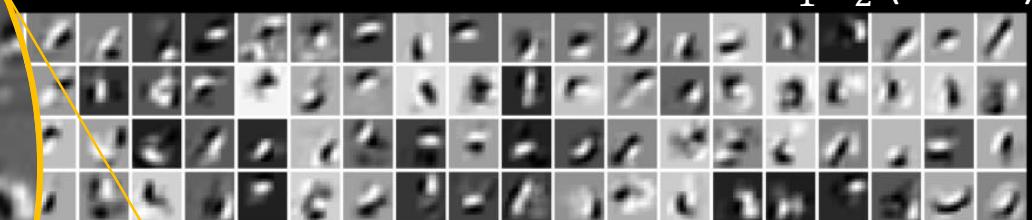
- $D_1$ : 32 filters of size 7x7 - 99.99 % sparse
- $D_2$ : 128 filters of size 3x3 - 99.09 % sparse
- $D_3$ : 1024 filters of size 2x2 (dense)



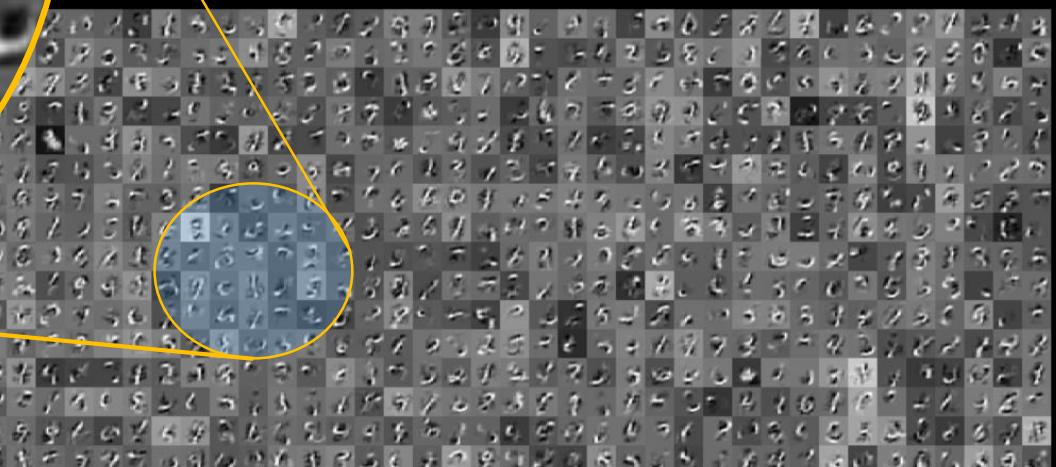
$D_1$  (7x7)



$D_1 D_2$  (15x15)



$D_1 D_2 D_3$  (28x28)



# ML-CSC: Pursuit

- Deep–Coding Problem ( $\mathbf{DCP}_\lambda$ ) (dictionaries are known):

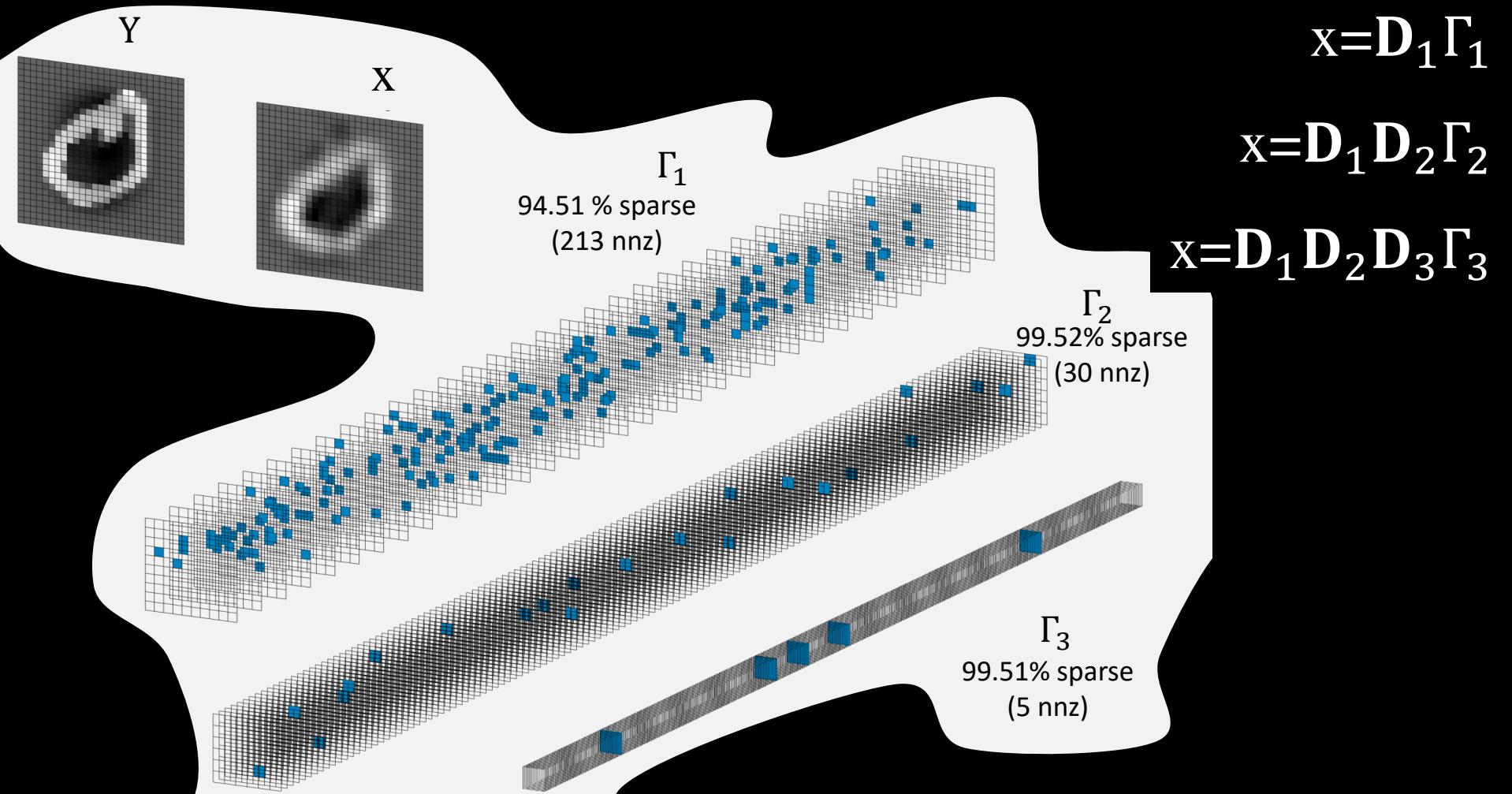
$$\left\{ \begin{array}{ll} \mathbf{X} = \mathbf{D}_1 \boldsymbol{\Gamma}_1 & \|\boldsymbol{\Gamma}_1\|_0 \leq \lambda_1 \\ \boldsymbol{\Gamma}_1 = \mathbf{D}_2 \boldsymbol{\Gamma}_2 & \|\boldsymbol{\Gamma}_2\|_0 \leq \lambda_2 \\ \vdots & \vdots \\ \boldsymbol{\Gamma}_{K-1} = \mathbf{D}_K \boldsymbol{\Gamma}_K & \|\boldsymbol{\Gamma}_K\|_0 \leq \lambda_K \end{array} \right\}$$

- Or, more realistically for noisy signals,

$$\text{Find } \{\boldsymbol{\Gamma}_j\}_{j=1}^K \text{ s.t. } \left\{ \begin{array}{ll} \|\mathbf{Y} - \mathbf{D}_1 \boldsymbol{\Gamma}_1\|_2 \leq \varepsilon & \|\boldsymbol{\Gamma}_1\|_0 \leq \lambda_1 \\ \boldsymbol{\Gamma}_1 = \mathbf{D}_2 \boldsymbol{\Gamma}_2 & \|\boldsymbol{\Gamma}_2\|_0 \leq \lambda_2 \\ \vdots & \vdots \\ \boldsymbol{\Gamma}_{K-1} = \mathbf{D}_K \boldsymbol{\Gamma}_K & \|\boldsymbol{\Gamma}_K\|_0 \leq \lambda_K \end{array} \right\}$$



# A Small Taste: Pursuit



# Consider this for Solving the DCP

- Layered Thresholding (LT):  
Estimate  $\Gamma_1$  via the THR algorithm

$$\widehat{\Gamma}_2 = \underbrace{\mathcal{P}_{\beta_2} \left( \mathbf{D}_2^T \mathcal{P}_{\beta_1} \left( \mathbf{D}_1^T \mathbf{Y} \right) \right)}$$

Estimate  $\Gamma_2$  via the THR algorithm

$$(\mathbf{DCP}_\lambda^\varepsilon): \text{Find } \{\Gamma_j\}_{j=1}^K \text{ s.t.}$$
$$\begin{cases} \|\mathbf{Y} - \mathbf{D}_1 \Gamma_1\|_2 \leq \varepsilon & \|\Gamma_1\|_0 \leq \lambda_1 \\ \Gamma_1 = \mathbf{D}_2 \Gamma_2 & \|\Gamma_2\|_0 \leq \lambda_2 \\ \vdots & \vdots \\ \Gamma_{K-1} = \mathbf{D}_K \Gamma_K & \|\Gamma_K\|_0 \leq \lambda_K \end{cases}$$

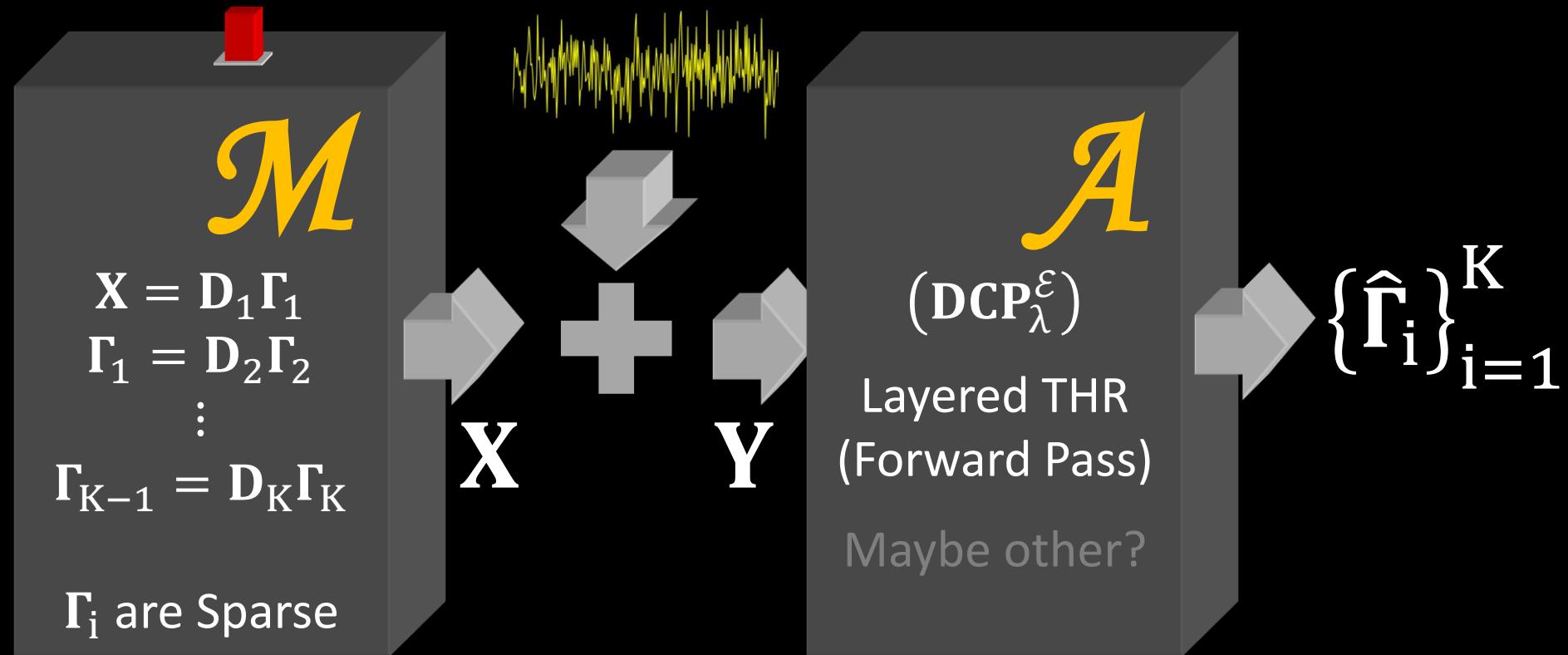
- Now let's take a look at how Conv. Neural Network operates:

$$f(\mathbf{Y}) = \text{ReLU} \left( \mathbf{b}_2 + \mathbf{W}_2^T \text{ReLU} \left( \mathbf{b}_1 + \mathbf{W}_1^T \mathbf{Y} \right) \right)$$

The layered (soft nonnegative)  
thresholding and the CNN forward pass  
algorithm are the very same thing !!!



# Theoretical Path



Armed with this view of a generative source model, we may ask new and daring theoretical questions



# Success of the Layered-THR

**Theorem:** If  $\|\Gamma_i\|_0 < \frac{1}{2} \left( 1 + \frac{1}{\mu(D_i)} \cdot \frac{|\Gamma_i^{\min}|}{|\Gamma_i^{\max}|} \right) - \frac{1}{\mu(D_i)} \cdot \frac{\varepsilon_L^{i-1}}{|\Gamma_i^{\max}|}$

then the **Layered Hard THR** (with the proper thresholds) **finds the correct supports** and  $\|\Gamma_i^{\text{LT}} - \Gamma_i\|_{2,\infty}^p \leq \varepsilon_L^i$ , where we have defined  $\varepsilon_L^0 = \|\mathbf{E}\|_2$  and

$$\varepsilon_L^i = \sqrt{\|\Gamma_i\|_0} \cdot (\varepsilon_L^{i-1} + \mu(D_i)(\|\Gamma_i\|_0 - 1)|\Gamma_i^{\max}|)$$

Papyan, Romano & Elad ('17)

The stability of the forward pass is guaranteed if the underlying representations are sparse and the noise is bounded

- Problems:**
1. Contrast
  2. Error growth
  3. Error even if no noise



# Layered Basis Pursuit (BP)

- We chose the Thresholding algorithm due to its simplicity, but we do know that there are better pursuit methods – how about using them?
- Lets use the Basis Pursuit instead ...

$$\text{(DCP}_{\lambda}^{\varepsilon}\text{): Find } \{\Gamma_j\}_{j=1}^K \text{ s.t.}$$
$$\left\{ \begin{array}{l} \|\mathbf{Y} - \mathbf{D}_1 \Gamma_1\|_2 \leq \varepsilon \quad \|\Gamma_1\|_0 \leq \lambda_1 \\ \Gamma_1 = \mathbf{D}_2 \Gamma_2 \\ \vdots \\ \Gamma_{K-1} = \mathbf{D}_K \Gamma_K \quad \|\Gamma_K\|_0 \leq \lambda_K \end{array} \right.$$

$$\Gamma_1^{\text{LBP}} = \min_{\Gamma_1} \frac{1}{2} \|\mathbf{Y} - \mathbf{D}_1 \Gamma_1\|_2^2 + \lambda_1 \|\Gamma_1\|_1$$



$$\Gamma_2^{\text{LBP}} = \min_{\Gamma_2} \frac{1}{2} \|\Gamma_1^{\text{LBP}} - \mathbf{D}_2 \Gamma_2\|_2^2 + \lambda_2 \|\Gamma_2\|_1$$



⋮

Does this algorithm work ?  
Is it better than the Layered-THR ?  
Can we provide theoretical guarantees for it?



# Success of the Layered BP

**Theorem:** Assuming that  $\|\Gamma_i\|_0 < \frac{1}{3} \left( 1 + \frac{1}{\mu(D_i)} \right)$

then the Layered Basis Pursuit performs very well:

1. The support of  $\Gamma_i^{\text{LBP}}$  is contained in that of  $\Gamma_i$
2. The error is bounded:  $\|\Gamma_i^{\text{LBP}} - \Gamma_i\|_2 \leq \varepsilon_L^i$ , where

$$\varepsilon_L^i = 7.5^i \|\mathbf{E}\|_2 \prod_{j=1}^i \sqrt{\|\Gamma_j\|_0}$$

3. Every entry in  $\Gamma_i$  greater than

$$\varepsilon_L^i / \sqrt{\|\Gamma_i\|_0} \text{ will be found}$$

Papyan, Romano & Elad ('17)

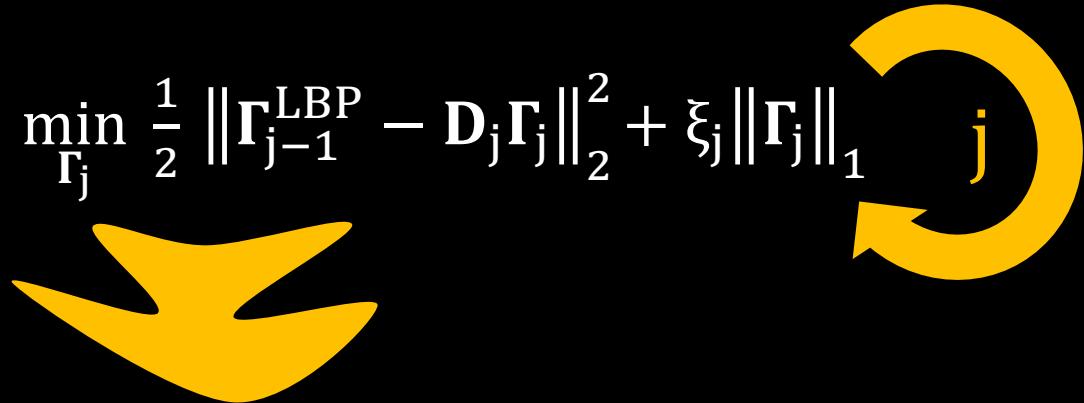
## Problems:

1. ~~Contrast~~
2. Error growth
3. ~~Error even if no noise~~



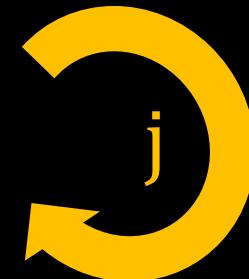
# Layered Iterative Thresholding

$$\text{Layered BP: } \Gamma_j^{\text{LBP}} = \min_{\Gamma_j} \frac{1}{2} \left\| \Gamma_{j-1}^{\text{LBP}} - D_j \Gamma_j \right\|_2^2 + \xi_j \left\| \Gamma_j \right\|_1$$



Layered Iterative Soft-Thresholding:

$$\Gamma_j^t = \mathcal{S}_{\xi_j/c_j} \left( \Gamma_j^{t-1} + D_j^T (\widehat{\Gamma}_{j-1} - D_j \Gamma_j^{t-1}) \right)$$



Note that our suggestion  
implies that groups of layers  
share the same dictionaries

Can be seen as a very deep  
recurrent neural network

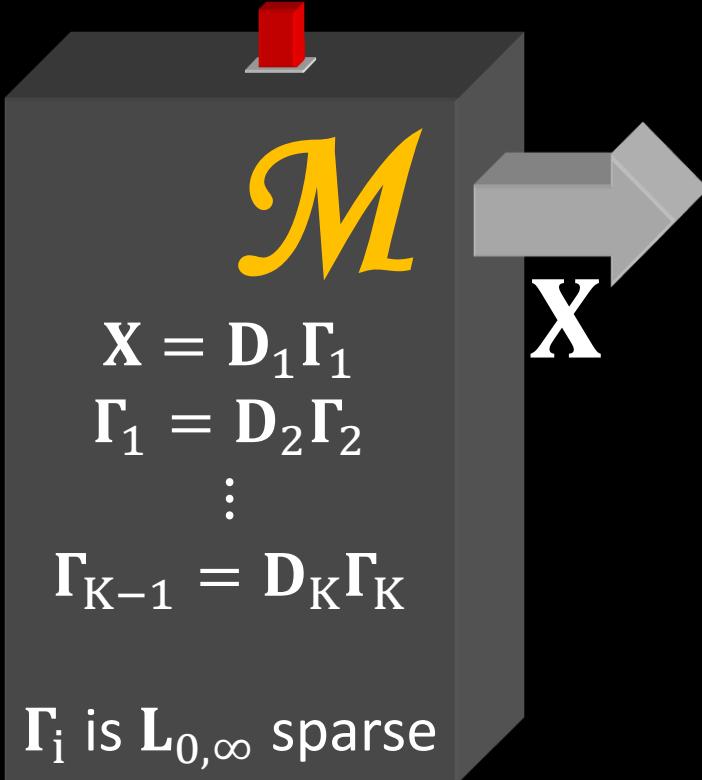
[Gregor & LeCun '10]



# Reflections and Recent Results



# Where are the Labels?



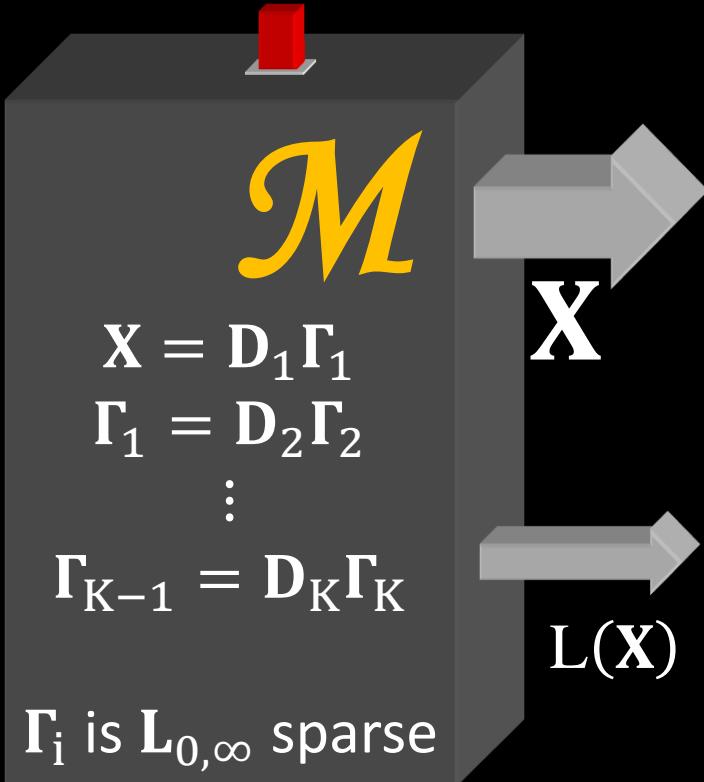
Answer 1:

- We do not need labels because everything we show refer to the **unsupervised** case, in which we operate on signals, not necessarily in the context of recognition

We presented the ML-CSC as a machine that produces signals  $X$



# Where are the Labels?



We presented the ML-CSC as a machine that produces signals  $\mathbf{X}$

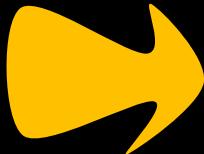
## Answer 2:

- In fact, this model could be augmented by a synthesis of the corresponding label by:  
$$L(\mathbf{X}) = sign\{c + \sum_{j=1}^K w_j^T \boldsymbol{\Gamma}_j\}$$
- This assumes that knowing the representations suffices for classification → **supervised** mode
- Thus, a successful pursuit algorithm can lead to an accurate recognition if the network is augmented by a FC classification layer
- In fact, we can analyze theoretically the classification accuracy and the sensitivity to adversarial noise – see later



# What About Learning?

**Sparseland**  
Sparse  
Representation  
Theory



**CSC**  
Convolutional  
Sparse  
Coding



**ML-CSC**  
Multi-Layered  
Convolutional  
Sparse Coding

All these models rely on proper  
**Dictionary Learning Algorithms** to fulfil their mission:

- Sparseland: We have unsupervised and supervised such algorithms, and a beginning of theory to explain how these work
- CSC: We have few and only unsupervised methods, and even these are not fully stable/clear
- ML-CSC: Two algorithms were proposed – unsupervised (to appear in IEEE-TSP) and supervised (submitted to IEEE-TPAMI)



# Fresh from the Oven (1)

## Main Focus:

- Better pursuit &
- Dictionary learning

## Contributions:

- Proposed a projection based pursuit (i.e. Verifying that the obtained signal obeys the synthesis equations), accompanied by better theoretical guarantees
- Proposes **the first dictionary learning algorithm** for the ML-CSC model for an unsupervised mode of work (as an auto-encoder, and trading representations' sparsities by dictionary sparsity)

## Multilayer Convolutional Sparse Modeling: Pursuit and Dictionary Learning

Jeremias Sulam , Member, IEEE, Vardan Petyan , Yaniv Romano , and Michael Elad , Fellow, IEEE

*Abstract*—The recently proposed multilayer convolutional sparse coding (ML-CSC) model, consisting of a cascade of convolutional sparse layers, provides a new interpretation of convolutional neural networks (CNNs). Under this framework, the forward pass in a CNN is equivalent to a pursuit algorithm aiming to estimate the nested sparse representation vectors from a given input signal. Despite having served as a pivotal connection between CNNs and sparse modeling, a deeper understanding of the ML-CSC is still lacking. In this paper, we propose a sound pursuit algorithm for the ML-CSC model by adopting a projection approach. We provide new and improved bounds on the stability of the solution of such pursuit and we analyze different practical alternatives to implement this in practice. We show that the training of the filters is essential to allow for nontrivial signals in the model, and we derive an online algorithm to learn the dictionaries from real

as atoms [1]. Backed by elegant theoretical results, this model led to a series of works dealing either with the problem of the pursuit of such decompositions, or with the design and learning of better atoms from real data [2]. The latter problem, termed dictionary learning, empowered sparse enforcing methods to achieve remarkable results in many different fields from signal and image processing [3]–[5] to machine learning [6]–[8].

Neural networks, on the other hand, were introduced around forty years ago and were shown to provide powerful classification algorithms through a series of function compositions [9], [10]. It was not until the last half-decade, however, that through a series of incremental modifications these methods

To appear in IEEE-TSP



# Fresh from the Oven (2)

## Main Focus:

- Holistic pursuit &
- Relation to the Co-Sparse analysis model

## Contributions:

- Proposed a systematic way to synthesize signals from the ML-CSC model
- Develop performance bounds for the oracle in various pursuit strategies
- Constructs the first provable **holistic pursuit** that mixes greedy-analysis and relaxation-synthesis pursuit algorithms

### MULTI LAYER SPARSE CODING: THE HOLISTIC WAY

AVIAD ABERDAM\*, JEREMIAS SULAM†, AND MICHAEL ELAD‡

**Abstract.** The recently proposed multi-layer sparse model has raised insightful connections between sparse representations and convolutional neural networks (CNN). In its original conception, this model was restricted to a cascade of *convolutional synthesis* representations. In this paper, we start by addressing a more general model, revealing interesting ties to fully connected networks. We then show that this multi-layer construction admits a brand new interpretation in a unique symbiosis between synthesis and analysis models: while the deepest layer indeed provides a synthesis representation, the mid-layers decompositions provide an analysis counterpart. This new perspective exposes the suboptimality of previously proposed pursuit approaches, as they do not fully leverage all the information comprised in the model constraints. Armed with this understanding, we address fundamental theoretical issues, revisiting previous analysis and expanding it. Motivated by the limitations of previous algorithms, we then propose an integrated – *holistic* – alternative that estimates all representations in the model simultaneously, and analyze all these different schemes under stochastic noise assumptions. Inspired by the synthesis-analysis duality, we further present a Holistic Pursuit algorithm, which alternates between synthesis and analysis sparse coding steps, eventually solving for the entire model as a whole, with provable improved performance. Finally, we present numerical results that demonstrate the practical advantages of our approach.

To Appear in SIMODS



# Fresh from the Oven (3)

## Main Focus:

- Take the labels into account
- Analyze classification performance and sensitivity to adversarial noise

## Contributions:

- Develop bounds on the **maximal adversarial noise** that guarantees a proper classification
- Expose the higher sensitivity of poor pursuit methods (Layered-THR) over better ones (Layered-BP)

### Adversarial Noise Attacks of Deep Learning Architectures – Stability Analysis via Sparse Modeled Signals

Yaniv Romano · Aviad Aberdam · Jeremias Sulam · Michael Elad

Received: date / Accepted: date

**Abstract** Despite their impressive performance, deep convolutional neural networks (CNNs) have been shown to be sensitive to small adversarial perturbations. These nuisances, which one can barely notice, are powerful enough to fool sophisticated and well performing classifiers, leading to ridiculous misclassification results. In this paper we analyze the stability of state-of-the-art

## 1 Introduction

Deep learning, and in particular Convolutional Neural Networks (CNN), is one of the hottest topics in data sciences as it has led to many state-of-the-art results spanning across many domains [13,8]. Despite the evident great success of classifying images, it has been

Submitted to JMIV



# Fresh from the Oven (4)

## Main Focus:

- Better and provable ISTA-like pursuit algorithm
- Examine the effect of the number of iterations in the unfolded architecture

## Contributions:

- Develop a **novel ISTA-like algorithms** for the ML-CSC model, with proper mathematical justifications
- Demonstrate the architecture obtained when unfolding this algorithm
- Show that for the same number of parameters, more iterations lead to better classification

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL., NO., NOVEMBER 2018

1

## On Multi-Layer Basis Pursuit, Efficient Algorithms and Convolutional Neural Networks

Jeremias Sulam, *Member, IEEE*, Aviad Aberdam, Amir Beck, Michael Elad, *Fellow, IEEE*

**Abstract**—Parsimonious representations are ubiquitous in modeling and processing information. Motivated by the recent Multi-Layer Convolutional Sparse Coding (ML-CSC) model, we herein generalize the traditional Basis Pursuit problem to a multi-layer setting, introducing similar sparse enforcing penalties at different representation layers in a symbiotic relation between synthesis and analysis sparse priors. We explore different iterative methods to solve this new problem in practice, and we propose a new Multi-Layer Iterative Soft Thresholding Algorithm (ML-ISTA), as well as a fast version (ML-FISTA). We show that these nested first order algorithms converge, in the sense that the function value of near-fixed points can get arbitrarily close to the solution of the original problem. We further show how these algorithms effectively implement particular recurrent convolutional neural networks (CNNs) that generalize feed-forward ones without introducing any parameters. We present and analyze different architectures resulting unfolding the iterations of the proposed pursuit algorithms, including a new Learned ML-ISTA, providing a principled way to construct deep recurrent CNNs. Unlike other similar constructions, these architectures unfold a global pursuit holistically for the entire network. We demonstrate the emerging constructions in a supervised learning setting, consistently improving the performance of classical CNNs while maintaining the number of parameters constant.

**Index Terms**—Multi-Layer Convolutional Sparse Coding, Network Unfolding, Recurrent Neural Networks, Iterative Shrinkage Algorithms.

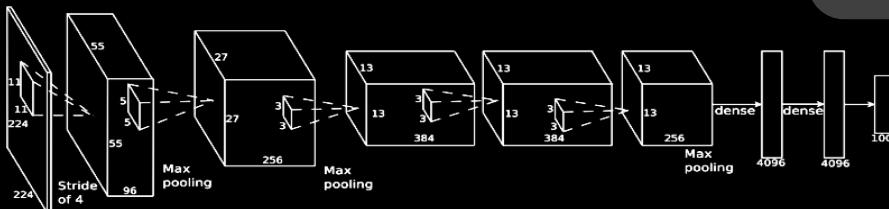
Submitted to IEEE-TPAMI



# Time to Conclude



# This Talk



A novel interpretation  
and theoretical  
understanding of CNN

What does it mean that the learned filters are sparse?  
Can we find better filters for faster, more accurate  
learning, and developing it further

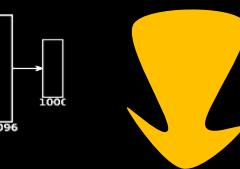
*Sparseland*



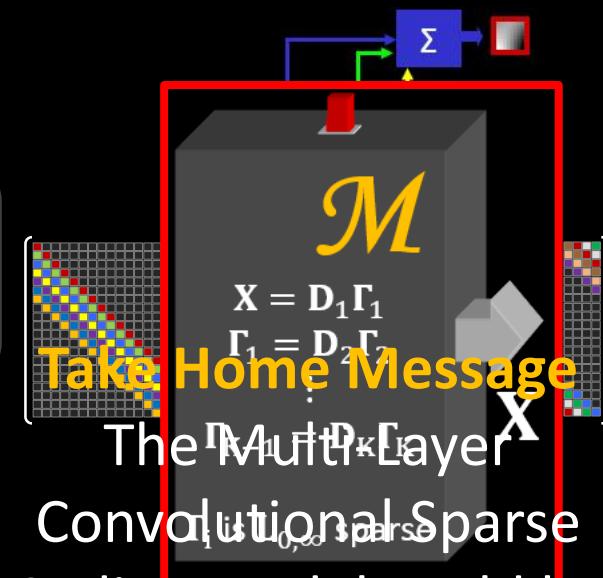
The desire to  
model data



Novel View of  
Convolutional  
Sparse Coding



Multi-Layer  
Convolutional  
Sparse Coding



Take Home Message:  
The Multi-Layer  
Convolutional Sparse  
Coding model could be  
a new platform for  
theoretically  
understanding deep  
learning, and  
developing it further



# A New Massive Open Online Course

**edX** Courses ▾ Programs ▾ Schools & Partners About ▾ Search:  Sign In Register

**Israel X**  
The National Online Education Initiative

## Sparse Representations in Signal and Image Processing

Learn the theory, tools and algorithms of sparse representations and their impact on signal and image processing.

[Start the Professional Certificate Program](#)



Courses in the Professional Certificate Program

 Sparse Representations in Signal and Image Processing: Fundamentals  
Learn about the field of sparse representations by understanding its fundamental theoretical and algorithmic foundations.  
[Learn more](#)

 Sparse Representations in Image Processing: From Theory to Practice  
Learn about the deployment of the sparse representation model to signal and image processing.  
[Learn more](#)

Starts on October 25, 2017 [Enroll Now](#)

I would like to receive email from IsraelX and learn about other offerings related to Sparse Representations in Signal and Image Processing: Fundamentals.

Starts on February 28, 2018 [Enroll Now](#)

I would like to receive email from IsraelX and learn about other offerings related to Sparse Representations in Image Processing: From Theory to Practice.

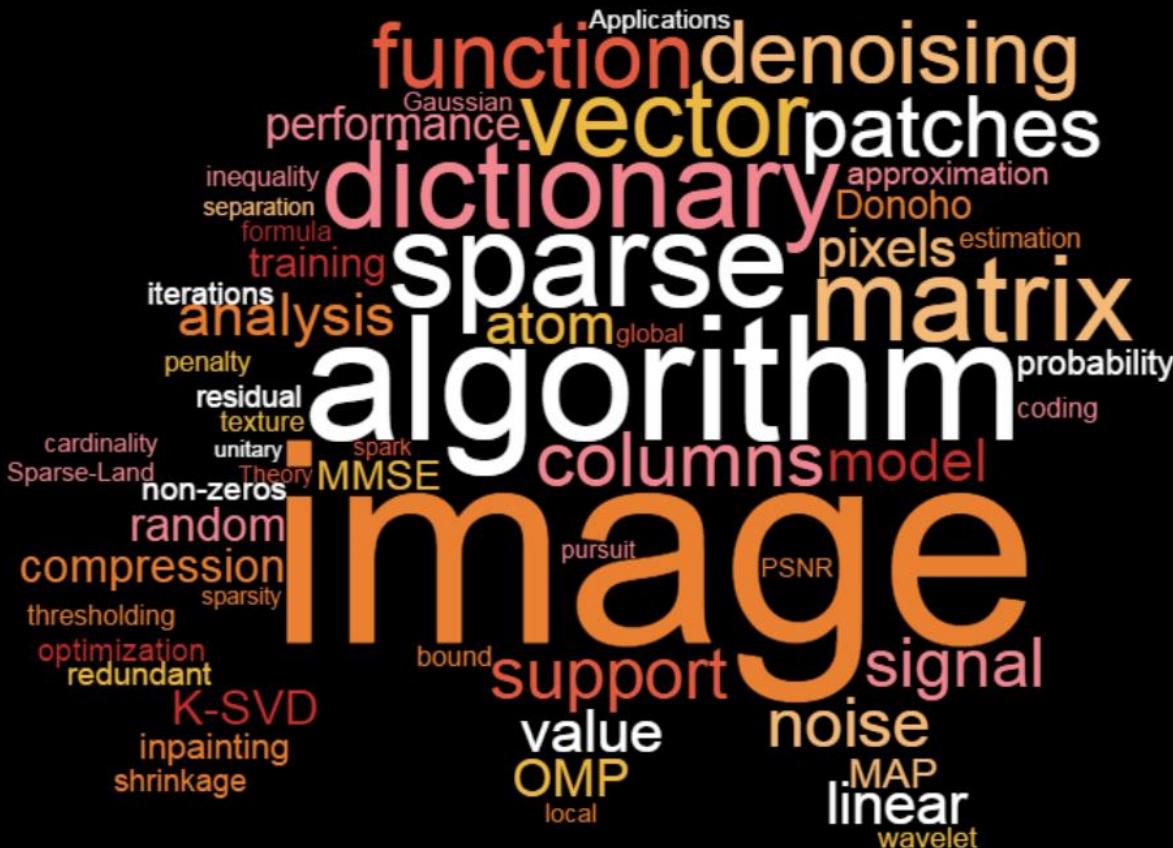
---

Instructors

 Michael Elad  
The Computer-Science Department  
The Technion

 Yaniv Romano

 Michael Elad



More on these (including these slides and the relevant papers) can be found in <http://www.cs.technion.ac.il/~elad>

