# A Tale of Signal Modeling Evolution SparseLand $\rightarrow$ CSC $\rightarrow$ CNN

#### Michael Elad

The Computer Science Department Technion – Israel Institute of Technology

Workshop on Frame Theory and Sparse Representation for Complex Data Institute for Mathematical Sciences May 29th – June 2nd



Joint work with



Yaniv Romano



Vardan Papyan Jeremias Sulam



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



#### In This Talk



CNN<sup>\*</sup> Convolutional Neural Networks

The Underlying Idea

# Modeling

data sources enables a theoretical analysis of algorithms' performance

\* Only CNN? What about other architectures ?



# Part I Motivation and Background



#### Our Starting Point: Image Denoising

Original Image

White Gaussian Noise

Noisy Image





Many (thousands) image denoising algorithms have been proposed over the years, some of which are extremely effective





### Leading Image Denoising Methods...

are built upon powerful patch-based local models:









Popular local models: GMM

Sparse-Representation Example-based Low-rank Field-of-Experts & Neural networks





## Patch-Based Image Denoising

- K-SVD: sparse representation modeling of image patches
   [Elad & Aharon, '06]
- BM3D: combines sparsity and self-similarity
   [Dabov, Foi, Katkovnik & Egiazarian '07]
- EPLL: uses GMM of the image patches
   [Zoran & Weiss '11]
- MLP: multi-layer perceptron
   [Burger, Schuler & Harmeling '12]
- NCSR: non-local sparsity with centralized coefficients
   [Dong, Zhang, Shi & Li '13]
- WNNM: weighted nuclear norm of image patches
   [Gu, Zhang, Zuo & Feng '14]
- SSC–GSM: nonlocal sparsity with a GSM coefficient model [Dong, Shi, Ma & Li '15]







#### The SparseLand Model for Patches

N

- Assumes that every patch is a linear combination of a few atoms, from a dictionary
- The operator  $\mathbf{R}_i$  extracts the i-th *n*-dimensional patch from  $\mathbf{X} \in \mathbb{R}^N$



echnion

KA

$$orall {\mathbf{i}}$$
 ,  $\mathbf{R}_{\mathbf{i}} \mathbf{X} = \mathbf{\Omega} oldsymbol{\gamma}_{\mathbf{i}}$  where  $\|oldsymbol{\gamma}_{\mathbf{i}}\|_0 \ll n$ 

i-th location



\*  $\boldsymbol{R}_i$  for 1D signals

#### Patch Denoising

Given a noisy patch  $\mathbf{R}_i \mathbf{Y}$ , solve  $(\mathbf{P}_0^{\epsilon})$ :  $\hat{\mathbf{y}}_i = \operatorname{argmin} \|\mathbf{y}_i\|_0$ 



Clean patch:  $\mathbf{\Omega} \widehat{\mathbf{\gamma}}_{i}$ 

s.t.  $\|\mathbf{R}_{i}\mathbf{Y} - \mathbf{\Omega}\mathbf{\gamma}_{i}\|_{2} \leq \epsilon$ 

 $(\mathbf{P}_0^{\varepsilon})$  is hard to solve



Greedy methods such as Orthogonal Matching Pursuit (OMP) or Thresholding (F Convex relaxations such as Basis Pursuit (BP)

 $\mathbf{g} \quad (\mathbf{P}_1^{\epsilon}): \min_{\mathbf{y}_i} \|\mathbf{y}_i\|_1 + \xi \|\mathbf{R}_i\mathbf{Y} - \mathbf{\Omega}\mathbf{y}_i\|_2^2$ 





#### Recall K-SVD Denoising [Elad & Aharon, '06]



- Despite its simplicity, this is a very well-performing algorithm
   Its origins can be traced back to Guleryuz's local DCT recovery
- A small modification of this method leads to state-of-the-art results [Mairal, Bach, Ponce, Spairo, Zisserman, `09]





## What is Missing?

 Over the years, many kept revisiting this algorithm and its line of thinking, with a clear feeling that key features are still lacking



- What is missing? Here is what **WE** thought of...
  - A multi-scale treatment [Ophir, Lustig & Elad '11] [Sulam, Ophir & Elad '14] [Papyan & Elad '15]
  - Exploiting self-similarities [Ram & Elad '13] [Romano, Protter & Elad '14]
  - Pushing to better agreement on the overlaps [Romano & Elad '13] [Romano & Elad '15]
  - Enforcing the local model on the final patches (EPLL) [Sulam & Elad '15]
- $\,\circ\,$  Eventually, we realized that the key part that is missing is

#### **A Theoretical Backbone**



## Missing Theoretical Backbone?

 $\circ$  The core global-local model assumption on  $\mathbf{X} \in \mathbb{R}^N$ :

 $\forall i \quad \mathbf{R}_i \mathbf{X} = \mathbf{\Omega} \mathbf{\gamma}_i \quad \text{where} \quad \|\mathbf{\gamma}_i\|_0 \leq k$ 



Every patch in the unknown signal is expected to have a sparse representation w.r.t. the same dictionary  $\boldsymbol{\Omega}$ 

Questions to consider:

- Who are the signals belonging to this model? Do they exist?
- How should we project a signal on this model (pursuit)?
- Could we offer theoretical guarantees for this model/algorithms?
- Could we offer a global pursuit algorithm that operates locally?
- How should we learn  $\mathbf{\Omega}$  if this is indeed the model?

 As we will see, all these questions are very relevant to recent developments in signal processing and machine learning





#### Coming Up . ||<sub>0</sub> << K Limitations of **Convolutional Sparse** Coding (CSC) model patch averaging Multi-Layer Convolutional Theoretical Sparse Coding (ML-CSC) study of CSC Convolutional neural Fresh view of CNN through networks (CNN) the eyes of sparsity



# Part II Convolutional Sparse Coding

Working Locally Thinking Globally: Theoretical Guarantees for Convolutional Sparse Coding Vardan Papyan, Jeremias Sulam and Michael Elad

Convolutional Dictionary Learning via Local Processing Vardan Papyan, Yaniv Romano, Jeremias Sulam, and Michael Elad





## Convolutional Sparse Coding (CSC)



i-th feature-map: An image of the same size as **X** holding the sparse representation related to the i-filter





#### Intuitively ...



#### CSC in Matrix Form

• Here is an alternative global sparsity-based model formulation

$$\mathbf{X} = \sum_{i=1}^{m} \mathbf{C}^{i} \mathbf{\Gamma}^{i} = \mathbf{D} \mathbf{\Gamma}$$

**C**i

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







#### Two Interpretations



#### Why CSC?





**Technion** 

18

### CSC Relation to Our Story

 $\circ$  A clear global model: every patch has a sparse representation w.r.t. to the same local dictionary  $\Omega$ , just as we have assumed

 $\odot\,\text{No}$  notion of disagreement on the patch overlaps

 $\circ$  Related to the current common practice of patch averaging ( $\mathbf{R}_i^T$  - put the patch  $\mathbf{\Omega} \mathbf{\gamma}_i$  back in the i-th location of the global vector)

$$\mathbf{X} = \mathbf{D}\boldsymbol{\Gamma} = \frac{1}{n} \sum_{i} \mathbf{R}_{i}^{\mathrm{T}} \boldsymbol{\Omega} \boldsymbol{\gamma}_{i}$$

• What about the Pursuit?

- "Patch averaging": independent sparse coding for each patch
- CSC: should seek all the representations together

 $\odot$  Is there a bridge between the two? We'll come back to this later ...





 This model has been used in the past [Lewicki & Sejnowski '99] [Hashimoto & Kurata, '00]

 Most works have focused on solving *efficiently* its associated pursuit, called **convolutional sparse coding**, using the BP algorithm

 $(\mathbf{P}_{1}^{\epsilon}): \min_{\mathbf{\Gamma}} \|\mathbf{\Gamma}\|_{1} + \lambda \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2}^{2} \qquad \text{Convolutional} \\ \checkmark \qquad \text{dictionary}$ 

Several applications were demonstrated:

- Pattern detection in images and the analysis of instruments in music signals [Mørup, Schmidt & Hansen '08]
- Inpainting [Heide, Heidrich & Wetzstein '15]
- Super-resolution [Gu, Zuo, Xie, Meng, Feng & Zhang '15]

However, little is known regrading its theoretical aspects. Why?
 Perhaps because the regular SparsLand theory is sufficient?



# Classical Sparse Theory (Noiseless) ( $\mathbf{P}_0$ ): min $\|\mathbf{\Gamma}\|_0$ s.t. $\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$ **Definition**: Mutual Coherence: $\mu(\mathbf{D}) = \max_{i \neq i} |d_i^T d_j|$ [Donoho & Elad '03] **Theorem:** For a signal $\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$ , if $\|\mathbf{\Gamma}\|_0 < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D})}\right)$ then this solution is necessarily the sparsest [Donoho & Elad '03] **Theorem**: The OMP and BP are guaranteed to recover the true sparse code assuming that $\|\Gamma\|_0 < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right)$ [Tropp '04], [Donoho & Elad '03]

21

## The Need for a Theoretical Study

 $\circ$  Assuming that m=2 and n=64 we have that [Welch, '74]  $\mu(\mathbf{D}) \geq 0.063$ 

 As a result, uniqueness and success of pursuits is guaranteed as long as

$$\|\mathbf{\Gamma}\|_{0} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right) \le \frac{1}{2} \left(1 + \frac{1}{0.063}\right) \approx 8$$

Less than 8 non-zeros GLOBALLY are allowed!!!
 This is a very pessimistic result!

- Repeating the above for the noisy case leads to even worse performance predictions
- Bottom line: Classic SparseLand Theory cannot provide good explanations for the CSC model







#### Moving to Local Sparsity: Stripes

 $\ell_{0,\infty}$  Norm:  $\|\mathbf{\Gamma}\|_{0,\infty}^{s} = \max_{i} \|\mathbf{\gamma}_{i}\|_{0}$ 

 $(\mathbf{P}_{0,\infty})$ : min  $\|\mathbf{\Gamma}\|_{0,\infty}^{s}$  s.t.  $\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$ 

 $\|\Gamma\|_{0,\infty}^s$  is low  $\to$  all  $\gamma_i$  are sparse  $\to$  every patch has a sparse representation over  $\Omega$ 

#### The Main Questions we Aim to Address:

- I. Is the solution to this problem unique ?
- II. Can we recover the solution via a global OMP/BP ?



 $m = 2 \prec$ 





#### Uniqueness via Mutual Coherence

 $(\mathbf{P}_{0,\infty})$ : min  $\|\mathbf{\Gamma}\|_{0,\infty}^{s}$  s.t.  $\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$ 

**Theorem**: If a solution  $\Gamma$  is found for  $(\mathbf{P}_{0,\infty})$  such that:

$$\|\boldsymbol{\Gamma}\|_{0,\infty}^{\mathrm{s}} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right)$$

then this is necessarily the unique optimal solution to this problem

This result is exciting: This and later results pose a local constraint for a global guarantee, and as such, they are far more optimistic compared to the global guarantees For k non-zeros per stripe, and filters of length n, we get  $\|\mathbf{\Gamma}\|_0 \cong \frac{k}{2n-1} \cdot N$ non-zeros globally







#### Phase Transition Experiment

• We construct a dictionary with a low mutual coherence: m = 2, n = 64, N = 640

- We generate random sparse vectors in which the non-zero entries are drawn as random i.i.d Gaussians
- Given a sparse vector, we compute its global signal and attempt to recover it using the global OMP and BP



# From Ideal to Noisy Signals

 $\odot$  So far, we have assumed an ideal signal  $X=D\Gamma$ 

- $\odot$  However, in practice we usually have  $Y=D\Gamma+E$  where E is due to noise or model deviations
- $\odot\,\mbox{To}$  handle this, we redefine our problem as:

 $(\mathbf{P}_{0,\infty}^{\epsilon}): \min_{\mathbf{\Gamma}} \|\mathbf{\Gamma}\|_{0,\infty}^{s} \text{ s.t. } \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2} \leq \epsilon$ 

#### $\odot$ The Main Questions We Aim to Address:

- I. Stability of the solution to this problem ?
- II. Stability of the solution obtained via global OMP/BP?
- III. Could the same recovery be done via local (patch) operations ?





#### Stability of via Stripe-RIP

$$(\mathbf{P}_{0,\infty}^{\epsilon}): \min_{\mathbf{\Gamma}} \|\mathbf{\Gamma}\|_{0,\infty}^{s} \text{ s.t. } \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2} \leq \epsilon \longrightarrow \widehat{\mathbf{\Gamma}}$$

**Definition: D** is said to satisfy Stripe-RIP with constant  $\delta_k$  if  $(1 - \delta_k) \|\Delta\|_2^2 \le \|\mathbf{D}\Delta\|_2^2 \le (1 + \delta_k) \|\Delta\|_2^2$ for any vector  $\Delta$  with  $\|\Delta\|_{0,\infty}^s = k$ 



[Candes & Tao '05]

#### Local Noise Assumption

- $\circ$  Thus far, our analysis relied on the local sparsity of the underlying solution **Γ**, which was enforced through the  $\ell_{0,\infty}$  norm
- $\odot$  In what follows, we present stability guarantees for both OMP and BP that will also depend on the local energy in the noise vector E
- $\circ$  This will be enforced via the  $\ell_{2,\infty}$  norm, defined as:

E

₽ 2,∞	= 1	nax i	<b>R</b> <sub>i</sub> E

2



#### Stability of OMP

Theorem: If  $\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$  where  $\|\mathbf{\Gamma}\|_{0,\infty}^{s} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right) - \frac{1}{\mu(\mathbf{D})} \cdot \frac{\|\mathbf{E}\|_{2,\infty}^{p}}{|\Gamma_{\min}|}$ then OMP run for  $\|\mathbf{\Gamma}\|_{0}$  iterations will 1. Find the correct support 2.  $\|\mathbf{\Gamma}_{OMP} - \mathbf{\Gamma}\|_{2}^{2} \le \frac{\|\mathbf{E}\|_{2}^{2}}{1 - (\|\mathbf{\Gamma}\|_{0,\infty}^{s} - 1)\mu(\mathbf{D})}$ 





# Stability of Lagrangian BP

$$(\mathbf{P}_{1}^{\epsilon}): \quad \mathbf{\Gamma}_{\mathrm{BP}} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2}^{2} + \lambda \|\mathbf{\Gamma}\|_{1}$$

**Theorem:** For  $\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$ , if  $\lambda = 4||\mathbf{E}||^p$ 

$$\|\boldsymbol{\Gamma}\|_{0,\infty}^{\mathrm{s}} < \frac{1}{3} \left(1 + \right)$$

Then we are guaranteed that

- 1. The support of  $m{\Gamma}_{
  m BP}$  is contain
- 2.  $\|\mathbf{\Gamma}_{\mathrm{BP}} \mathbf{\Gamma}\|_{\infty} \le 7.5 \|\mathbf{E}\|_{2,\infty}^{\mathrm{p}}$
- 3. Every entry greater than 7.5
- 4.  $\Gamma_{\rm BP}$  is unique

IIEIP and Theoretical foundation for recent works tackling the convolutional sparse coding problem via BP [Bristow, Eriksson & Lucey '13] [Wohlberg '14] [Kong & Fowlkes '14] [Bristow & Lucey '14] [Heide, Heidrich & Wetzstein '15] [Šorel & Šroubek '16]





#### Phase Transition - Noisy

 $\odot$  We use the same dictionary as in the noiseless case

- We generate random sparse vectors in which the non-zero entries are drawn randomly in the range [-a, a] for different a values
- Given a sparse vector, we compute its global signal and attempt to recover it using the global OMP and BP



#### **Global** Pursuit via Local Processing

$$(\mathbf{P}_{1}^{\epsilon}): \quad \mathbf{\Gamma}_{\mathrm{BP}} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2}^{2} + \xi \|\mathbf{\Gamma}\|_{1}$$

- While CSC is a global model, its theoretical guarantees rely on local properties
- We aim to show that this global-local
  - relation can also be exploited for solving the global BP problem using only local operations





 $\alpha_i$ 

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

## Global Pursuit via Local Processing (1)

$$(\mathbf{P}_{1}^{\epsilon}): \quad \mathbf{\Gamma}_{\mathrm{BP}} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2}^{2} + \xi \|\mathbf{\Gamma}\|_{1}$$

 Recall: Iterative Soft Thresholding is an appealing method for handling the above minimization task

Projection 
$$\Gamma^{t} = S_{\xi/c} \left( \Gamma^{t-1} + \frac{1}{c} D^{T} (Y - D\Gamma^{t-1}) \right)$$
  
onto  $L_{1}$  ball  $\mathcal{I}$  Gradient step

This algorithm is guaranteed to solve the above problem
 [Daubechies, Defrise, De-Mol, 2004] [Blumensath & Davies '08]

 $\odot$  Proposal: We shall manipulate this algorithm to an equivalent form that operates locally \* c > 0.5  $\lambda_{max}$ (**D**<sup>T</sup>**D**)



#### Global Pursuit via Local Processing (1)




#### Simulation

#### Details:

• Signal length: N = 300• Patch size: n = 25• Unique atoms: p = 5• Local sparsity (k) is 11 • Global sparsity: k = 40• Number of iterations: 400 • Lagrangian:  $\xi = 4 ||\mathbf{E}||_{2,\infty}^{p}$ • Noise level: PSNR= 0.03



True Sparse Code
 Iterative Soft Thresholding



#### Global Pursuit via Local Processing (2)

$$(\mathbf{P}_{1}^{\epsilon}): \quad \mathbf{\Gamma}_{\mathrm{BP}} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2}^{2} + \xi \|\mathbf{\Gamma}\|_{1}$$

 Here is an alternative approach, based on a different interpretation of this linear system

 s<sub>i</sub> are slices – local patches that overlap to form the full image

$$\mathbf{X} = \mathbf{D}\mathbf{\Gamma} = \sum_{i} \mathbf{R}_{i}^{\mathrm{T}} \mathbf{D}_{\mathrm{L}} \boldsymbol{\alpha}_{i} = \sum_{i} \mathbf{R}_{i}^{\mathrm{T}} \mathbf{s}_{i} \quad \boldsymbol{\mathbf{A}}_{i}$$



F

 $X = D\Gamma$ 

## Global Pursuit via Local Processing (2)

$$(\mathbf{P}_{1}^{\epsilon}): \quad \mathbf{\Gamma}_{\mathrm{BP}} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2}^{2} + \lambda \|\mathbf{\Gamma}\|_{1}$$

Turning to the local form and using the Augmented Lagrangian

$$\min_{\boldsymbol{\alpha}_{i},\boldsymbol{s}_{i}} \frac{1}{2} \left\| \mathbf{Y} - \sum_{i} \mathbf{R}_{i}^{\mathrm{T}} \mathbf{s}_{i} \right\|_{2}^{2} + \sum_{i} \left( \lambda \|\boldsymbol{\alpha}_{i}\|_{1} + \frac{\rho}{2} \|\mathbf{s}_{i} - \mathbf{D}_{\mathrm{L}} \boldsymbol{\alpha}_{i} + \mathbf{u}_{i}\|_{2}^{2} \right)$$

 $\circ$  These two problems are equivalent, and convex w.r.t their variables

 The new formulation targets the local slices, and their sparse representations

 $\odot$  The vectors  $u_i$  are the Lagrange multipliers for the constraints  $s_i = D_L \alpha_i$ 



#### Global Pursuit via Local Processing (2)

$$\min_{\boldsymbol{\alpha}_{i},\boldsymbol{s}_{i}} \frac{1}{2} \left\| \mathbf{Y} - \sum_{i} \mathbf{R}_{i}^{\mathrm{T}} \mathbf{s}_{i} \right\|_{2}^{2} + \sum_{i} \left( \lambda \|\boldsymbol{\alpha}_{i}\|_{1} + \frac{\rho}{2} \|\mathbf{s}_{i} - \mathbf{D}_{\mathrm{L}} \boldsymbol{\alpha}_{i} + \mathbf{u}_{i}\|_{2}^{2} \right)$$

$$ADMM$$

• Slice-update: 
$$\min_{s_i} \frac{1}{2} \| \mathbf{Y} - \sum_{i} \mathbf{R}_{i}^{T} \mathbf{s}_{i} \|^{2} + \sum_{i} \rho_{i} |\mathbf{s}_{i} - \mathbf{D}_{L} \boldsymbol{\alpha}_{i} + \mathbf{u}_{i} \|_{2}^{2}$$
  
• Sparse-Upda  
• Sparse-Upda  
• Sparse-Upda  
• Sparse-Upda  
• Sparse-Upda  
• Sparse-Upda

Comment: One iteration of this procedure amounts to ... the very same patch-averaging algorithm we started with



#### Two Comments About this Scheme

#### 

Patches extracted from natural images, and their corresponding slices. Observe how the slices are far simpler, and contained by their corresponding patches



The Proposed Scheme can be used for Dictionary (D $_L$ ) Learning

Slice-based DL algorithm using standard patch-based tools, leading to a faster and simpler method, compared to existing methods

[Wohlberg, 2016]

Ours

## Partial Summary of CSC

- O What we have seen so far is a new way to analyze the global CSC model using local sparsity constraints. We proved:
  - <u>k</u>j
    - **Uniqueness** of the solution for the noiseless problem
    - Stability of the solution for the noisy problem
  - NA F
- Guarantee of success and stability of both OMP and BP
- ¥)
- We obtained guarantees and algorithms that operate locally while claiming global optimality



We mentioned briefly the mater of learning the model (i.e. dictionary learning for CSC), and presented our competitive approach which is based on simple local steps





# Part III Going Deeper

Convolutional Neural Networks Analyzed via Convolutional Sparse Coding

Vardan Papyan, Yaniv Romano and Michael Elad





#### CSC and CNN

• There is an analogy between CSC and CNN:

- Convolutional structure
- Data driven models
- ReLU is a sparsifying operator

 $\odot$  We propose a principled way to analyze CNN

○ But first, a short review of CNN...





#### CNN



[LeCun, Bottou, Bengio and Haffner '98][Krizhevsky, Sutskever & Hinton '12][Simonyan & Zisserman '14][He, Zhang, Ren & Sun '15]





ReLU(z) = max(Thr, z)

#### CNN



Notice that we do not include a pooling stage:

- Can be replaced by a convolutional layer with increased stride without loss in performance [Springenberg, Dosovitskiy, Brox & Riedmiller '14]
- The current state-of-the-art in image recognition does not use it [He, Zhang, Ren & Sun '15]



#### Mathematically...

 $f(\mathbf{Y}, {\mathbf{W}_{i}}, {\mathbf{b}_{i}}) = \text{ReLU}(\mathbf{b}_{2} + \mathbf{W}_{2}^{T} \text{ReLU}(\mathbf{b}_{1} + \mathbf{W}_{1}^{T}\mathbf{X}))$ 





## Training Stage of CNN

 $\odot$  Consider the task of classification

 $\circ$  Given a set of signals  $\{Y_j\}_j$  and their corresponding labels  $\{h(Y_j)\}_i$ , the CNN learns an end-to-end mapping





#### Back to CSC

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



Convolutional sparsity (CSC) assumes an inherent structure is present in natural signals



#### Technion Israel Institute of Technology

#### Intuition: From Atoms to Molecules



#### A Small Taste: Model Training (MNIST)

MNIST Dictionary:

•  $D_1$ : 32 filters of size 7×7, with stride of 2 (dense)

- D<sub>2</sub>: 128 filters of size 5×5×32 with stride of 1 99.09 % sparse
- •D3: 1024 filters of size 7×7×128 99.89 % sparse

# $D_1D_2$ (15×15)

#### $D_1 D_2 D_3$ (28×28)

 $D_1$  (7×7)



#### A Small Taste: Pursuit





#### A Small Taste: Pursuit



## A Small Taste: Model Training (CFAR)

 $\mathbf{D}_{1}\mathbf{D}_{2}$  (13×13)

CIFAR Dictionary:

- D<sub>1</sub>: 64 filters of size 5x5x3, stride of 2 dense
- D<sub>2</sub>: 256 filters of size 5x5x64, stride of 2 82.99 % sparse
- D<sub>3</sub>: 1024 filters of size 5x5x256
   90.66 % sparse

**D**<sub>1</sub> (5×5×3)

 $D_1 D_2 D_3$  (32×32)



#### ML-CSC: Pursuit

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

Find 
$$\{\boldsymbol{\Gamma}_{j}\}_{j=1}^{K}$$
 s.t. 
$$\begin{cases} \mathbf{X} = \mathbf{D}_{1}\boldsymbol{\Gamma}_{1} & \|\boldsymbol{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \boldsymbol{\Gamma}_{1} = \mathbf{D}_{2}\boldsymbol{\Gamma}_{2} & \|\boldsymbol{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \boldsymbol{\Gamma}_{K-1} = \mathbf{D}_{K}\boldsymbol{\Gamma}_{K} & \|\boldsymbol{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{cases}$$

• Or, more realistically for noisy signals,

Find 
$$\{\mathbf{\Gamma}_{j}\}_{j=1}^{K}$$
 s.t. 
$$\begin{cases} \|\mathbf{Y} - \mathbf{D}_{1}\mathbf{\Gamma}_{1}\|_{2} \leq \mathcal{E} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} = \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} = \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{cases}$$



#### **ML-CSC:** Dictionary Learning

o Deep-Learning Problem (**DLP** $_{\lambda}$ ):

$$\text{and } \{\mathbf{D}_{i}\}_{i=1}^{K} \quad s.t. \begin{cases} \left\| \mathbf{Y}_{j} - \mathbf{D}_{1}\mathbf{\Gamma}_{1}^{j} \right\|_{2}^{2} \leq \mathcal{E} \quad \left\| \mathbf{\Gamma}_{1}^{j} \right\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{2}^{j} = \mathbf{D}_{2}\mathbf{\Gamma}_{1}^{2} \quad \left\| \mathbf{\Gamma}_{2}^{j} \right\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K}^{j} = \mathbf{D}_{K}\mathbf{\Gamma}_{K}^{j} \quad \left\| \mathbf{\Gamma}_{K}^{j} \right\|_{0,\infty}^{s} \leq \lambda_{K} \end{pmatrix}_{j=1}^{j}$$

 While the above is an unsupervised DL, a supervised version can be envisioned

Fi

 $\min_{\{\mathbf{D}_i\}_{i=1}^{K}, \mathbf{U}} \sum_{j} \ell\left(h(\mathbf{Y}_j), \mathbf{U}, \mathbf{DCP}^{\star}(\mathbf{Y}_j, \{\mathbf{D}_i\})\right)$ 

The deepest representation  $\Gamma_{\!K}$  obtained by solving the DCP





ML-CSC: The Simplest Pursuit

Keep it simple! • The simplest pursuit algorithm (single-layer case) is the THR algorithm, which operates on a given input signal Y by:

$$\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$$
 and  $\mathbf{\Gamma}$  is sparse  $\sum \widehat{\mathbf{\Gamma}} = \mathcal{P}_{\beta}(\mathbf{D}^{\mathrm{T}}\mathbf{Y})$ 

• Restricting the coefficients to be nonnegative does not restrict the expressiveness of the model





#### Consider this for Solving the DCP

 $\circ$  Layered thresholding (LT): Estimate  $\Gamma_1$  via the THR algorithm

$$\widehat{\boldsymbol{\Gamma}}_{2} = \mathcal{P}_{\beta_{2}}\left(\boldsymbol{D}_{2}^{\mathrm{T}}\mathcal{P}_{\beta_{1}}\left(\boldsymbol{D}_{1}^{\mathrm{T}}\boldsymbol{Y}\right)\right)$$

Estimate  $\Gamma_{\!2}$  via the THR algorithm

• Forward pass of CNN:

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

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





$\left(\mathbf{DCP}_{\lambda}^{\mathcal{E}}\right)$ : Find	$\left\{\mathbf{\Gamma}_{\mathbf{j}}\right\}_{\mathbf{j}=1}^{\mathbf{K}}$ s.t.
$\ \mathbf{Y} - \mathbf{D}_1 \mathbf{\Gamma}_1\ _2 \le \mathcal{E}$	$\ \mathbf{\Gamma}_1\ _{0,\infty}^{\mathrm{s}} \leq \lambda_1$
$\Gamma_1 = \mathbf{D}_2 \Gamma_2$	$\ \mathbf{\Gamma}_2\ _{0,\infty}^{\mathrm{s}} \leq \lambda_2$
6 6 6	• • •
$\mathbf{\Gamma}_{\mathrm{K-1}} = \mathbf{D}_{\mathrm{K}}\mathbf{\Gamma}_{\mathrm{K}}$	$\ \mathbf{\Gamma}_{\mathbf{K}}\ _{0,\infty}^{s} \leq \lambda_{\mathbf{K}}$

## Consider this for Solving the DLP

• DLP (supervised<sup>\*</sup>):

$$\min_{\{\mathbf{D}_i\}_{i=1}^{K}, \mathbf{U}} \sum_{j} \ell\left(h(\mathbf{Y}_j), \mathbf{U}, \mathbf{D}\mathbf{C}\mathbf{P}^{\star}(\mathbf{Y}_j, \{\mathbf{D}_i\})\right)$$

The thresholds for the DCP should also learned

Estimate via the layered THR algorithm

• CNN training:

$$\min_{\{\mathbf{W}_i\},\{\mathbf{b}_i\},\mathbf{U}}\sum_{j}\ell\left(h(\mathbf{Y}_j),\mathbf{U},f(\mathbf{Y},\{\mathbf{W}_i\},\{\mathbf{b}_i\})\right)$$

The problem solved by the training stage of CNN and the DLP are equivalent as well, assuming that the DCP is approximated via the layered thresholding algorithm

 Recall that for the ML-CSC, there exists an unsupervised avenue for training the dictionaries that has no simple parallel in CNN





#### Theoretical Path



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





#### Theoretical Path: Possible Questions

 Having established the importance of the ML-CSC model and its associated pursuit, the DCP problem, we now turn to its analysis

• The main questions we aim to address:

- I. Uniqueness of the solution (set of representations) to the (**DCP**<sub> $\lambda$ </sub>)?
- II. Stability of the solution to the  $(\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}})$  problem ?
- III. Stability of the solution obtained via the hard and soft layered THR algorithms (forward pass) ?
- IV. Limitations of this (very simple) algorithm and alternative pursuit?

V. Algorithms for training the dictionaries  $\{\mathbf{D}_i\}_{i=1}^K$  vs. CNN ? VI. New insights on how to operate on signals via CNN ?



61

## Uniqueness of $(DCP_{\lambda})$

 $(\mathbf{DCP}_{\lambda})$ : Find a set of representations satisfying  $\mathbf{X} = \mathbf{D}_1 \mathbf{\Gamma}_1 \qquad \|\mathbf{\Gamma}_1\|_{0,\infty}^{\mathsf{s}} \leq \lambda_1$ Is this set  $\Gamma_1 = \mathbf{D}_2 \Gamma_2 \qquad \|\Gamma_2\|_{0,\infty}^{s} \le \lambda_2$ unique?  $\Gamma_{K-1} = \mathbf{D}_{K}\Gamma_{K} \quad \|\Gamma_{K}\|_{0,\infty}^{s} \leq \lambda_{K}$ **Theorem:** If a set of solutions  $\{\Gamma_i\}_{i=1}^K$  is found for  $(\mathbf{DCP}_{\lambda})$  such that:  $\|\boldsymbol{\Gamma}_{\mathbf{i}}\|_{0,\infty}^{\mathbf{s}} \leq \lambda_{\mathbf{i}} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D}_{\mathbf{i}})}\right)$ then these are necessarily the unique solution to the DCP problem The feature maps CNN aims to recover are unique





# Stability of $(\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}})$

• The problem we aim to solve is this

$$\begin{split} \left( \mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}} \right) &: \text{Find a set of representations satisfying} \\ \|\mathbf{Y} - \mathbf{D}_{1}\mathbf{\Gamma}_{1}\|_{2} \leq \mathcal{E} \quad \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} &= \mathbf{D}_{2}\mathbf{\Gamma}_{2} \quad \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ &\vdots &\vdots \\ \mathbf{\Gamma}_{K-1} &= \mathbf{D}_{K}\mathbf{\Gamma}_{K} \quad \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{split}$$



• Suppose that we manage to solve the  $(\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}})$  and find a feasible set of representations satisfying all the conditions



• The question we pose is How close is  $\widehat{\Gamma}_i$  to  $\Gamma_i$ ?





# Stability of $(\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}})$

Theorem: If the true representations  $\{\Gamma_i\}_{i=1}^K$  satisfy  $\|\Gamma_i\|_{0,\infty}^s \leq \lambda_i < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D}_i)}\right)$ then the set of solutions  $\{\widehat{\Gamma}_i\}_{i=1}^K$  obtained by solving this problem (somehow) must obey  $\|\widehat{\Gamma}_i - \Gamma_i\|_2^2 \leq \mathcal{E}_i^2$  for  $\mathcal{E}_0^2 = 4\mathcal{E}^2, \qquad \mathcal{E}_i^2 = \frac{\mathcal{E}_{i-1}^2}{1 - (2\lambda_i - 1)\mu(\mathbf{D}_i)}$ 

The problem CNN aims to solve is stable under certain conditions

Observe this annoying effect of error magnification as we dive into the model





#### Stability of Layered-THR

$$\begin{split} \text{Theorem: If } \|\Gamma_{i}\|_{0,\infty}^{s} &< \frac{1}{2} \left(1 + \frac{1}{\mu(D_{i})} \cdot \frac{\left|\Gamma_{i}^{min}\right|}{\left|\Gamma_{i}^{max}\right|}\right) - \frac{1}{\mu(D_{i})} \cdot \frac{\epsilon_{L}^{i-1}}{\left|\Gamma_{i}^{max}\right|} \\ \text{then the layered hard THR (with the proper thresholds) will find the correct supports* and \\ \left\|\Gamma_{i}^{LT} - \Gamma_{i}\right\|_{2,\infty}^{p} &\leq \epsilon_{L}^{i} \\ \text{where we have defined } \epsilon_{L}^{0} &= \|\mathbf{E}\|_{2,\infty}^{p} \text{ and } \\ \epsilon_{L}^{i} &= \sqrt{\|\Gamma_{i}\|_{0,\infty}^{p}} \cdot \left(\epsilon_{L}^{i-1} + \mu(\mathbf{D}_{i})\left(\|\Gamma_{i}\|_{0,\infty}^{s} - 1\right)|\Gamma_{i}^{max}|\right) \end{split}$$

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

\* Least-Squares update of the non-zeros?





#### Limitations of the Forward Pass

 The stability analysis reveals several inherent limitations of the forward pass (a.k.a. Layered THR) algorithm:

- Even in the noiseless case, the forward pass is incapable of recovering the perfect solution of the DCP problem
- Its success depends on the ratio  $|\Gamma_i^{\min}|/|\Gamma_i^{\max}|$ . This is a direct consequence of relying on a simple thresholding operator
- The distance between the true sparse vector and the estimated one increases exponentially as a function of the layer depth

 We now turn to propose a new algorithm that attempts to solve some of these problems



#### Special Case – Sparse Dictionaries

 $\odot$  Throughout the theoretical study we assumed that the representations in the different layers are  $L_{0,\infty}$ -sparse

 $\circ$  Do we know of a simple example of a set of dictionaries  $\{D_i\}_{i=1}^K$ and their corresponding signals X that will obey this property?



• Assuming the dictionaries are sparse:

$$\left\|\boldsymbol{\Gamma}_{j}\right\|_{0,\infty}^{s} \leq \|\boldsymbol{\Gamma}_{K}\|_{0,\infty}^{s} \prod_{i=j+1}^{K} \|\boldsymbol{D}_{i}\|_{0}^{\gamma} \xrightarrow{\text{Maximal number of non-zeros in an atom in } \boldsymbol{D}_{i}}$$

 $\odot$  In the context of CNN, the above happens if a sparsity promoting regularization, such as the  $L_1$ , is employed on the filters



#### Better Pursuit ?

 $\circ~(\text{DCP}_{\lambda})$  Noiseless: Find a set of representations satisfying

$$\begin{split} \mathbf{X} &= \mathbf{D}_{1}\mathbf{\Gamma}_{1} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} &= \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} &= \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{split}$$

 $\odot$  So far we proposed the Layered THR:

$$\widehat{\mathbf{\Gamma}}_{K} = \mathcal{P}_{\beta_{K}} \left( \mathbf{D}_{K}^{\mathrm{T}} \dots \mathcal{P}_{\beta_{2}} \left( \mathbf{D}_{2}^{\mathrm{T}} \mathcal{P}_{\beta_{1}} \left( \mathbf{D}_{1}^{\mathrm{T}} \mathbf{X} \right) \right) \right)$$

 $\odot$  The motivation is clear – getting close to what CNN use

 However, this is the simplest and weakest pursuit known in the field of sparsity – Can we offer something better?



#### Layered Basis Pursuit (Noiseless)

 $\circ$  Our Goal: (**DCP**<sub> $\lambda$ </sub>): Find a set of representations satisfying

$$\begin{split} \mathbf{X} &= \mathbf{D}_{1} \mathbf{\Gamma}_{1} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} &= \mathbf{D}_{2} \mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} &= \mathbf{D}_{K} \mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{split}$$

• We can propose a Layered Basis Pursuit Algorithm:

$$\Gamma_{1}^{\text{LBP}} = \min_{\Gamma_{1}} \|\Gamma_{1}\|_{1} \text{ s.t. } \mathbf{X} = \mathbf{D}_{1}\Gamma_{1}$$
$$\Gamma_{2}^{\text{LBP}} = \min_{\Gamma_{2}} \|\Gamma_{2}\|_{1} \text{ s.t. } \Gamma_{1}^{\text{LBP}} = \mathbf{D}_{2}\Gamma_{2}$$

Deconvolutional networks [Zeiler, Krishnan, Taylor & Fergus '10]



## Guarantee for Success of Layered BP

 As opposed to prior work in CNN, we can do far more than just proposing an algorithm – we can analyze its terms for success:



Theorem: If a set of representations  $\{\Gamma_i\}_{i=1}^K$  of the Multi-Layered CSC model satisfy

$$\|\boldsymbol{\Gamma}_{i}\|_{0,\infty}^{s} \leq \lambda_{i} < \frac{1}{2} \left(1 + \frac{1}{\mu(\boldsymbol{D}_{i})}\right)$$

then the Layered BP is guaranteed to find them

#### o Consequences:

- The layered BP can retrieve the underlying representations in the noiseless case, a task in which the forward pass fails to provide
- The Layered-BP's success does not depend on the ratio  $|\Gamma_i^{min}|/|\Gamma_i^{max}|$



#### Layered Basis Pursuit (Noisy)

$$\boldsymbol{\Gamma}_{1}^{\text{LBP}} = \min_{\boldsymbol{\Gamma}_{1}} \frac{1}{2} \| \boldsymbol{Y} - \boldsymbol{D}_{1} \boldsymbol{\Gamma}_{1} \|_{2}^{2} + \lambda_{1} \| \boldsymbol{\Gamma}_{1} \|_{1}$$
$$\boldsymbol{\Gamma}_{2}^{\text{LBP}} = \min_{\boldsymbol{\Gamma}_{2}} \frac{1}{2} \| \boldsymbol{\Gamma}_{1}^{\text{LBP}} - \boldsymbol{D}_{2} \boldsymbol{\Gamma}_{2} \|_{2}^{2} + \lambda_{2} \| \boldsymbol{\Gamma}_{2} \|_{1}$$

We can invoke a result we have seen already, referring to the BP for the CSC model:

**RECALL** For 
$$\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$$
, if  
 $\|\mathbf{\Gamma}\|_{0,\infty}^{s} < \frac{1}{3}\left(1 + \frac{1}{\mu(\mathbf{D})}\right)$   
then we are guaranteed that  
 $\|\mathbf{\Delta}\|_{2,\infty}^{p} \le 7.5 \varepsilon_{L}^{0} \sqrt{\|\mathbf{\Gamma}\|_{0,\infty}^{p}}$ 



## Stability of Layered BP

**Theorem:** Assuming that  $\|\Gamma_i\|_{0,\infty}^s < \frac{1}{3}\left(1 + \frac{1}{\mu(D_i)}\right)$ then For correctly chosen  $\{\lambda_i\}_{i=1}^K$  we are guaranteed that

. The support of  $\mathbf{\Gamma}_i^{\mathrm{LBP}}$  is contained in that of  $\mathbf{\Gamma}_i$ 

2. The error is bounded:  $\|\boldsymbol{\Gamma}_{i}^{\text{LBP}} - \boldsymbol{\Gamma}_{i}\|_{2,\infty}^{p} \leq \varepsilon_{L}^{i}$ , where

$$\varepsilon_{\mathrm{L}}^{\mathrm{i}} = 7.5^{\mathrm{i}} \|\mathbf{E}\|_{2,\infty}^{\mathrm{p}} \prod_{\mathrm{j}=1}^{\mathrm{I}} \sqrt{\|\mathbf{\Gamma}_{\mathrm{j}}\|_{0,\infty}^{\mathrm{p}}}$$

3. Every entry in  $\Gamma_i$  greater than  $\varepsilon_L^i / \sqrt{\|\Gamma_i\|_{0,\infty}^p}$  will be found


## Layered Iterative Thresholding

Layered BP: 
$$\Gamma_{j}^{LBP} = \min_{\Gamma_{j}} \frac{1}{2} \|\Gamma_{j-1}^{LBP} - D_{j}\Gamma_{j}\|_{2}^{2} + \xi_{j}\|\Gamma_{j}\|_{1}$$
  
Layered Iterative Soft-Thresholding:  
 $t = \Gamma_{j}^{t} = S_{\xi_{j}/c_{j}} \left(\Gamma_{j}^{t-1} + \frac{1}{c_{j}}D_{j}^{T}(\widehat{\Gamma}_{j-1} - D_{j}\Gamma_{j}^{t-1})\right)$   
Note that our suggestion  
implies that groups of layers  
share the same dictionaries  
 $Can be seen as a recurrent neural network [Gregor & LeCun '10]$   
 $* c_{i} > 0.5 \lambda_{max}(D_{i}^{T}D_{i})$   
 $* c_{i} > 0.5 \lambda_{max}(D_{i}^{T}D_{i})$   
 $T_{i}$ 

## Time to Conclude





## Current/Future Work

In general, we aim to leverage our theoretical insights in order to get to practical implications

More specifically, we work on:

- Developing alternative (local) pursuit methods for the CSC and ML-CSC
- Could we propose an MMSE-driven pursuit
- Training the dictionaries So far our efforts are focused on the unsupervised mode and the results are encouraging
- Explaining theoretically "known" tricks in CNN (local normalization, batchnormalization, the effect of stride, residual networks, dropout, ...
- Better understanding this model by projecting true signals on to it to see what kind of sparsities and dictionaries are obtained
- Improving the corresponding performance bounds, and
- Tying all the above to applications



76

These slides will be shared in my webpage in few days

## Questions?

