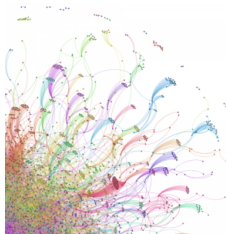


# Neural Networks as Sparsity Enforcing Algorithms

Jeremias Sulam



Deep Geometric Learning of Big Data and Applications  
IPAM – 2019

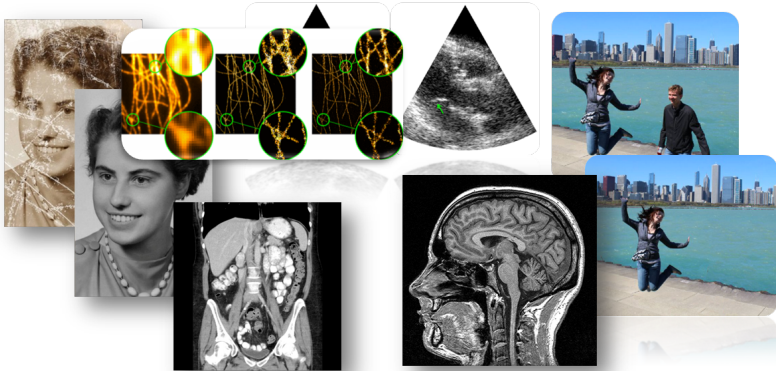


JOHNS HOPKINS  
WHITING SCHOOL  
of ENGINEERING



JOHNS HOPKINS  
MATHEMATICAL INSTITUTE  
for DATA SCIENCE

## Inverse Problems in Image Processing



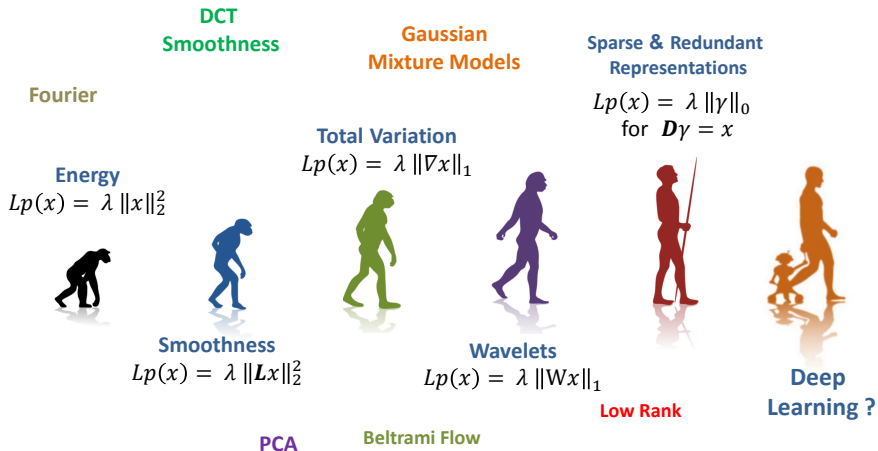
- All data has inherent **structure** than can be exploited
- This structure enables different **processing** tasks to be carried out

} **Signal Models**

find image that : is *related to the input* + looks like a "good" image

$$\min_{\mathbf{x}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 \quad - \quad \log P(\mathbf{x})$$

## Image Models



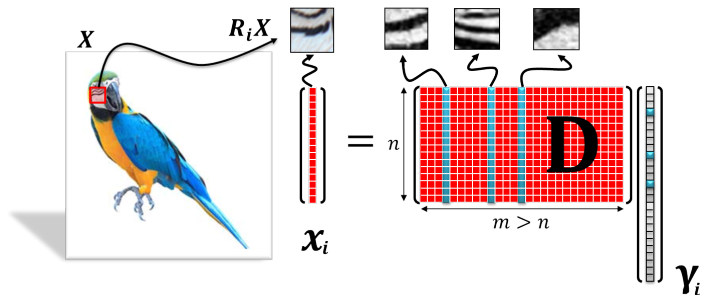
## Sparsity $\in$ (artificial) Neural Networks

- Ranzato & LeCun et al, *Efficient Learning of Sparse Representations with an Energy-Based Model* NIPS '07
- Lee & Ng et al, "Sparse deep belief net model for visual area V2" NIPS '08
- ...
- Koray & LeCun et al. "Learning convolutional feature hierarchies for visual recognition" NIPS '10.
- He, A. Szlam et al "Unsupervised feature learning by deep sparse coding" ICDM '14.
- Zeiler & Fergus "Deconvolutional networks" CVPR. '10
- Alireza & Frey. "K-sparse autoencoders" '13, "Winner Take All Auto-encoders" NIPS '15

## Neural Networks $\in$ Sparsity

- Gregor & LeCun "Learning fast approximations of sparse coding" ICML '10
- Moreau & Bruna "Understanding trainable sparse coding via matrix factorization" ICLR '17
- Liu et al. "On the convergence of learning-based iterative methods for nonconvex inverse problems" (2018).

## Sparse Representations (synthesis)



**Pursuit Problems:** How to find  $\boldsymbol{\gamma}$ ?

$$(P_0^\epsilon) : \min_{\boldsymbol{\gamma}_i} \|\mathbf{y}_i - \mathbf{D}\boldsymbol{\gamma}_i\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\gamma}_i\|_0 \leq s$$

## Pursuit Algorithms

$$(P_0^\epsilon) : \min_{\gamma} \|\mathbf{y} - \mathbf{D}\gamma\|_2^2 \quad \text{s.t.} \quad \|\gamma\|_0 \leq s$$

- **Greedy Algorithms**

- (Orthogonal) Matching Pursuit
- Iterative Hard Thresholding (P.G.D.)

- **Relaxation Approaches**

$$(P_1) : \min_{\gamma} \|\mathbf{y} - \mathbf{D}\gamma\|_2^2 + \lambda \|\gamma\|_1$$

- Convex optimization tools
- Soft Thresholding
- **Iterative Soft Thresholding (ISTA/FISTA)**

$$\hat{\gamma}^{t+1} = \mathcal{T}_{\lambda\eta} (\hat{\gamma}^t - \eta \mathbf{D}^T (\mathbf{D}\hat{\gamma}^t - \mathbf{y}))$$

## Sparse Representations

### Characterization of the Dictionary

- **Restricted Isometry Property (R.I.P.)** smallest  $\delta_s$  so that  $\forall \gamma : \|\gamma\|_0 = s$ ,

$$(1 - \delta_s)\|\gamma\|_2^2 \leq \|\mathbf{D}\gamma\|_2^2 \leq (1 + \delta_s)\|\gamma\|_2^2$$

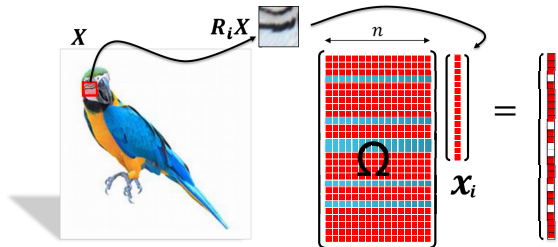
- **Mutual Coherence**  $\mu(\mathbf{D}) = \max_{i \neq j} \frac{|\mathbf{d}_i^T \mathbf{d}_j|}{\|\mathbf{d}_i\|_2 \|\mathbf{d}_j\|_2}$

### Sparse Representations are useful

Say  $\mathbf{y} = \mathbf{D}\gamma + \mathbf{w}$ ,  $\mathbf{w} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$

$$\mathbb{E} \|\gamma - \hat{\gamma}\|_2^2 \leq \frac{\sigma^2 s}{1 - \delta_{\lambda_s}} \ll \sigma^2 n$$

# Analysis Sparsity



Model

$$y = x + w, \quad \|\Omega x\|_0 \leq m - \ell$$

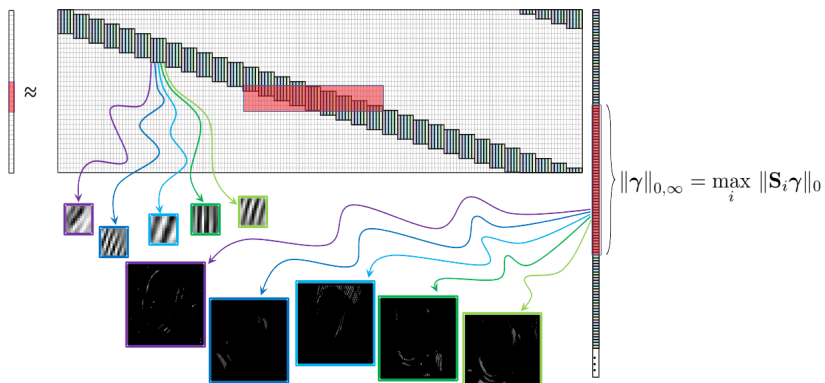
Pursuit

$$\min_x \|x - y\|_2^2 \quad \text{s.t.} \quad \|\Omega x\|_0 \leq m - \ell$$

## Convolutional Sparse Modeling of Signals

Model

$$\mathbf{y} = \mathbf{D}\boldsymbol{\gamma} + \mathbf{w}, \quad \|\boldsymbol{\gamma}\|_{0,\infty} \leq s$$



## Convolutional Sparse Modeling of Signals

Model

$$\mathbf{y} = \mathbf{D}\boldsymbol{\gamma} + \mathbf{w}, \quad \|\boldsymbol{\gamma}\|_{0,\infty} \leq s$$

Convolutional Pursuit

$$\min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{D}\boldsymbol{\gamma}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\gamma}\|_{0,\infty} \leq s$$

$\Downarrow$

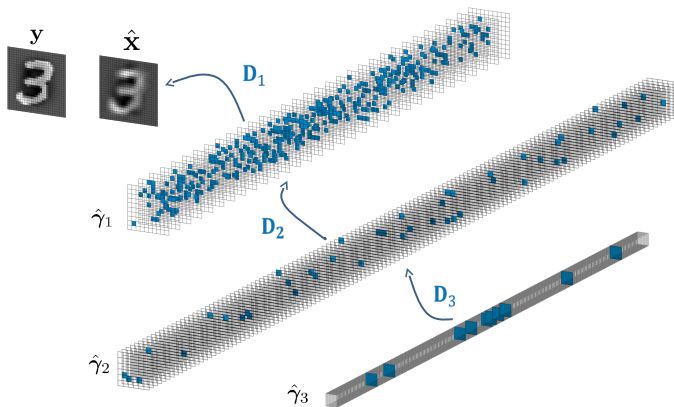
$$\min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{D}\boldsymbol{\gamma}\|_2^2 + \lambda \|\boldsymbol{\gamma}\|_1$$

[Pappyan, Sulam, Elad, IEEE TSP 2017]

## Multi-Layer Sparse Modeling of Signals

Model

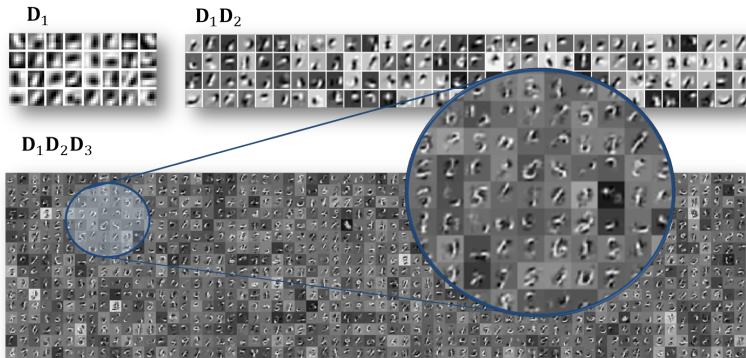
$$\mathbf{y} = \mathbf{D}_1 \boldsymbol{\gamma}_1 + \mathbf{w}, \quad \{\boldsymbol{\gamma}_{i-1} = \mathbf{D}_i \boldsymbol{\gamma}_i, \quad \|\boldsymbol{\gamma}_i\|_{0,\infty} \leq s_i\}_{i=1}^L$$



## Multi-Layer Sparse Modeling of Signals

Model

$$\mathbf{y} = \mathbf{D}_1 \boldsymbol{\gamma}_1 + \mathbf{w}, \quad \{\boldsymbol{\gamma}_{i-1} = \mathbf{D}_i \boldsymbol{\gamma}_i, \quad \|\boldsymbol{\gamma}_i\|_{0,\infty} \leq s_i\}_{i=1}^L$$



## Multi-Layer Sparse Modeling of Signals

Model

$$\mathbf{y} = \mathbf{D}_1 \boldsymbol{\gamma}_1 + \mathbf{w}, \quad \{\boldsymbol{\gamma}_{i-1} = \mathbf{D}_i \boldsymbol{\gamma}_i, \quad \|\boldsymbol{\gamma}_i\|_{0,\infty} \leq s_i\}_{i=1}^L$$

Pursuit: *Deep Coding Problem*

$$\begin{aligned} \text{find } \{\boldsymbol{\gamma}_i\}_{i=1}^L \quad \text{s.t.} \quad & \|\mathbf{y} - \mathbf{D}_1 \boldsymbol{\gamma}_1\|_2^2 \leq \varepsilon_0, & \|\boldsymbol{\gamma}_1\|_{0,\infty} \leq s_1 \\ & \|\boldsymbol{\gamma}_1 - \mathbf{D}_2 \boldsymbol{\gamma}_2\|_2^2 \leq \varepsilon_1, & \|\boldsymbol{\gamma}_2\|_{0,\infty} \leq s_2 \\ & \vdots & \vdots \\ & \|\boldsymbol{\gamma}_{L-1} - \mathbf{D}_L \boldsymbol{\gamma}_L\|_2^2 \leq \varepsilon_{L-1}, & \|\boldsymbol{\gamma}_L\|_{0,\infty} \leq s_L, \end{aligned}$$

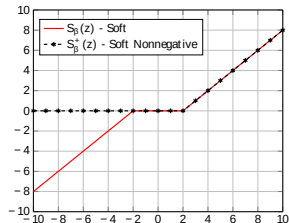
## Multi-Layer Thresholding (a.k.a Forward Pass)

$$\hat{\gamma}_1 = \text{ReLU}(\mathbf{D}_1^T \mathbf{y} + \mathbf{b}_1)$$

$$\hat{\gamma}_2 = \text{ReLU}(\mathbf{D}_2^T \hat{\gamma}_1 + \mathbf{b}_2)$$

$$\vdots$$

$$\hat{\gamma}_L = \text{ReLU}(\mathbf{D}_L^T \hat{\gamma}_{L-1} + \mathbf{b}_L)$$



## Stability

If a set of solutions  $\{\gamma_i\}_{i=1}^L$  satisfy  $\|\gamma_i\|_{0,\infty} < \frac{1}{2} \left( 1 + \frac{1}{\mu(\mathbf{D}_i)} \frac{|\gamma_i^{\min}|}{|\gamma_i^{\max}|} \right) - \frac{1}{\mu(\mathbf{D}_i)} \frac{\epsilon_L^{i-1}}{|\gamma_i^{\max}|}$ , then

- $\text{Supp}(\hat{\gamma}_i) = \text{Supp}(\gamma_i)$
- $\|\hat{\gamma}_i - \gamma_i\|_{2,\infty}^p \leq \sqrt{\|\gamma_i\|_{0,\infty}} (\epsilon_L^{i-1} + \mu(\mathbf{D}_i) (\|\gamma_i\|_{0,\infty} - 1) |\gamma_i^{\max}| + \beta_i)$

[Pappyan, Romano, Elad, JMLR 2017]

## Multi-Layer Basis Pursuit

$$\hat{\gamma}_1 = \arg \min_{\gamma} \|\mathbf{y} - \mathbf{D}_1\gamma\|_2^2 + \lambda_1 \|\gamma\|_1$$

$$\hat{\gamma}_2 = \arg \min_{\gamma} \|\hat{\gamma}_1 - \mathbf{D}_2\gamma\|_2^2 + \lambda_2 \|\gamma\|_1$$

...

$$\hat{\gamma}_L = \arg \min_{\gamma} \|\hat{\gamma}_{L-1} - \mathbf{D}_L\gamma\|_2^2 + \lambda_L \|\gamma\|_1$$

- Solutions **in the model**?
- What is really the **benefit of the model constraints**?
  - What can we get from all this **in practice**?

## Multi-Layer Sparse Modeling of Signals

### Model

$$\mathbf{y} = \mathbf{D}_{(1,L)}\boldsymbol{\gamma}_L + \mathbf{w}, \quad \|\boldsymbol{\gamma}_L\|_0 \leq s_L, \quad \{\|\mathbf{D}_{(i,L)}\boldsymbol{\gamma}_L\|_0 \leq s_{i-1}\}_{i=1}^L$$

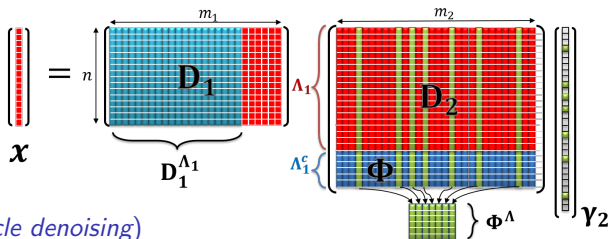
$$\mathbf{D}_{(1,L)} = \mathbf{D}_1\mathbf{D}_2 \dots \mathbf{D}_L.$$

### Pursuit

$$\min_{\boldsymbol{\gamma}_L} \|\mathbf{y} - \mathbf{D}_{(1,L)}\boldsymbol{\gamma}_L\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\gamma}_L\|_0 \leq s_L, \quad \{\|\mathbf{D}_{(i,L)}\boldsymbol{\gamma}_L\|_0 \leq s_{i-1}\}_{i=1}^L$$

$\Rightarrow$  Coupled **Synthesis** and **Analysis** priors!

## Synthesis-Analysis Interpretation

Stability (*oracle denoising*)

$$\mathbf{y} = \mathbf{D}_{(1,L)} \boldsymbol{\gamma}_L + \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

- Synthesis

$$\mathbb{E} \|\boldsymbol{\gamma}_2 - \hat{\boldsymbol{\gamma}}_2\|_2^2 \leq \frac{\sigma^2}{1 - \delta_{s_2}} s_2$$

- Synthesis-Analysis

$$\mathbb{E} \|\boldsymbol{\gamma}_2 - \hat{\boldsymbol{\gamma}}_2\|_2^2 \lesssim \frac{\sigma^2}{1 - \delta_{s_2}} (s_2 - \text{rank}(\Phi^\Lambda))$$

[Aberdam, Sulam, Elad, SIMODS 2019]

## Multi-Layer Pursuit

$$\min_{\gamma_L} \|\mathbf{y} - \mathbf{D}_{(1,L)}\gamma_L\|_2^2 \quad \text{s.t.} \quad \|\gamma_L\|_0 \leq s_L, \quad \{\|\mathbf{D}_{(i,L)}\gamma_L\|_0 \leq s_{i-1}\}_{i=1}^L$$

*Just relax:* Multi-Layer Basis Pursuit

$$(P) : \quad \min_{\gamma_2} \frac{1}{2} \|\mathbf{y} - \mathbf{D}_1 \mathbf{D}_2 \gamma_2\|_2^2 + \lambda_1 \|\mathbf{D}_2 \gamma_2\|_1 + \lambda_2 \|\gamma_2\|_1$$

## Related problems

- **Analysis Lasso** [Candes et al, 2010], **Robust Sparse Analysis Regularization** [Vaiter et al, 2012]
- **Inverse Problems with Compound Regularizers** [Figueiredo et al, 2010], [Haeffele et al, 2014]

## Multi-Layer ISTA

$$\min_{\gamma_2} F(\gamma_2) = f(\mathbf{D}_2 \gamma_2) + g_1(\mathbf{D}_2 \gamma_2) + g_2(\gamma_2)$$

*Proximal Gradient-Mapping*

$$\gamma_2^{k+1} = \text{prox}_{tg_2} \left( \gamma_2^k - t \mathbf{D}_2^T G_{1/\mu}^{f,g_1}(\gamma_1^k) \right)$$

where  $\gamma_1^k = \mathbf{D}_2 \gamma_2^k$  and

$$G_{1/\mu}^{f,g_1}(\gamma_1^k) = \frac{1}{\mu} \left[ \gamma_1^k - \text{prox}_{\mu g_1} \left( \gamma_1^k - \mu \nabla f(\gamma_1^k) \right) \right]$$

## Multi-Layer ISTA

$$\min_{\gamma_2} F(\gamma_2) = f(\mathbf{D}_2 \gamma_2) + g_1(\mathbf{D}_2 \gamma_2) + g_2(\gamma_2)$$

## Proximal Gradient-Mapping

$$\gamma_2^{k+1} = \text{prox}_{tg_2} \left( \gamma_2^k - t \mathbf{D}_2^T G_{1/\mu}^{f, g_1}(\gamma_1^k) \right)$$

In particular:

$$\gamma_2^{k+1} = \mathcal{T}_{t\lambda_2} \left( \gamma_2^k - \frac{t}{\mu} \mathbf{D}_2^T \left( \gamma_1^k - \mathcal{T}_{\mu\lambda_1} \left( \gamma_1^k - \mu \mathbf{D}_1^T (\mathbf{D}_1 \gamma_1^k - \mathbf{y}) \right) \right) \right)$$

## ML-ISTA

Set  $\gamma_0^k = \mathbf{y} \ \forall k$  and  $\gamma_L^1 = 0$ .

**for**  $k = 1 : K$  **do**           % for each iteration

$\hat{\gamma}_i \leftarrow \mathbf{D}_{(i,L)} \gamma_L^k \quad \forall i \in [0, k-1]$

**for**  $i = 1 : L$  **do**           % for each layer

$\gamma_i^{k+1} \leftarrow \mathcal{T}_{\mu_i \lambda_i} \left( \hat{\gamma}_i - \mu_i \mathbf{D}_i^T (\mathbf{D}_i \hat{\gamma}_i - \gamma_{i-1}^{k+1}) \right)$

✓ *Nested Proximal-Gradient updates*

## ML-FISTA

Set  $\gamma_0^k = \mathbf{y} \ \forall k$  and  $\gamma_L^1 = 0$ .

**for**  $k = 1 : K$  **do**            *% for each iteration*

$\hat{\gamma}_i \leftarrow \mathbf{D}_{(i,L)} \mathbf{z} \quad \forall i \in [0, k-1]$

**for**  $i = 1 : L$  **do**            *% for each layer*

$\gamma_i^{k+1} \leftarrow \mathcal{T}_{\mu_i \lambda_i} \left( \hat{\gamma}_i - \mu_i \mathbf{D}_i^T (\mathbf{D}_i \hat{\gamma}_i - \gamma_{i-1}^{k+1}) \right)$

$\mathbf{z} \leftarrow \gamma_L^{k+1} + \rho^k (\gamma_L^{k+1} - \gamma_L^k)$

✓ *Nested Proximal-Gradient updates with momentum*

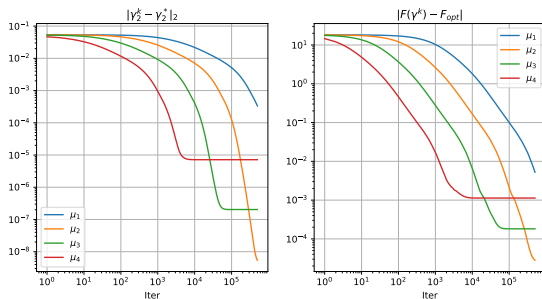
## Convergence for ML-ISTA

## Theorem

Suppose  $\{\gamma_2^k\}$  generated by ML-ISTA with  $\mu \in \left(0, \frac{1}{\|\mathbf{D}_1\|_2^2}\right)$  and  $t \in \left(0, \frac{4\mu}{3\|\mathbf{D}_2\|_2}\right)$ .  
 If  $\|\gamma_2^{k+1} - \gamma_2^k\|_2 \leq t\varepsilon$ , then

$$F(\gamma_2^{k+1}) - F_{opt} \leq \eta\varepsilon + (\beta + \kappa t)\mu,$$

where  $\eta$ ,  $\beta$  and  $\kappa$  are constants depending on  $\mathbf{D}_1$ ,  $\mathbf{D}_2$ ,  $g_1$ ,  $g_2$ .



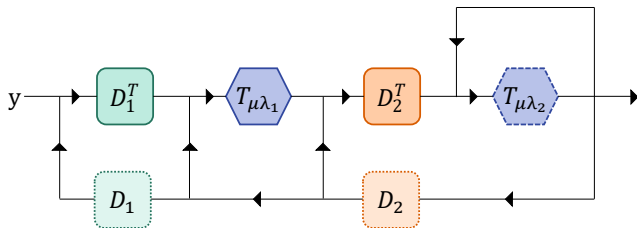
[Sulam, Aberdam, Beck, Elad, IEEE TPAMI 2019]

## Generalizing Recurrent CNNs

## Revisiting the ML-ISTA

$$\gamma_2^{(1)} = \text{ReLU}\left(\mu_2 \mathbf{D}_2^T \text{ReLU}(\mu_1 \mathbf{D}_1^T \mathbf{y} + \mathbf{b}_1) + \mathbf{b}_2\right) \quad (\text{F.P.})$$

$$\gamma_2^{(2)} = \text{ReLU}\left(\gamma_2^{(1)} - \mu_2 \mathbf{D}_2^T (\gamma_1^{(1)} - \text{ReLU}(\mu_1 \mathbf{D}_1^T (\mathbf{D}_1 \gamma_1^{(1)} - \mathbf{y}) + \mathbf{b}_1)) + \mathbf{b}_2\right)$$

$$\vdots$$


## Supervised Learning Formulation

$$\min_{\theta, \{\mathbf{D}_i, \lambda_i\}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(h_i, \zeta_{\theta}(\boldsymbol{\gamma}^*)) \quad \text{s.t.}$$

$$\boldsymbol{\gamma}^* = \arg \min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{D}_{(1,L)}\boldsymbol{\gamma}\|_2^2 + \sum_{i=1}^{L-1} \lambda_i \|\mathbf{D}_{(i+1,L)}\boldsymbol{\gamma}\|_1 + \lambda_L \|\boldsymbol{\gamma}\|_1.$$

$$\Downarrow$$

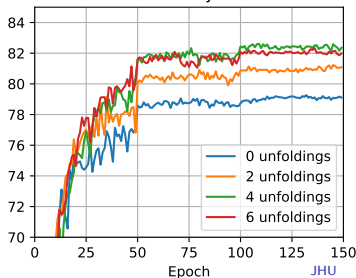
$$\min_{\theta, \{\mathbf{D}_i, \lambda_i\}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(h_i, \zeta_{\theta}(\boldsymbol{\gamma}^k))$$

$\boldsymbol{\gamma}^k$ :  $k^{\text{th}}$  iteration of ML-ISTA

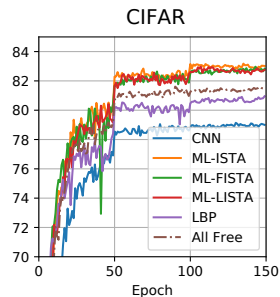
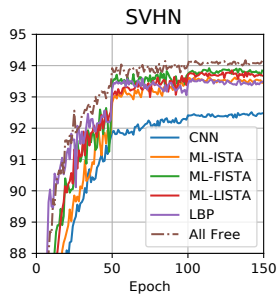
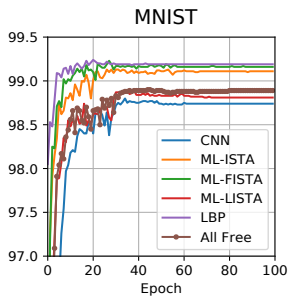
- if  $k = 1 \rightarrow$  Feed Forward CNN
- if  $k > 1 \rightarrow$  Recurrent (ML-ISTA) CNN

*But same number of parameters*

Test Accuracy - CIFAR



## Image Classification



## Take-Aways

- Multi-Layer Sparse Model puts forward a generative sparse model and framework for analysis
- Deep learning architectures can be understood analyzed as inference algorithms under this model
- ML-ISTA generalizes feed-forward networks, with improved performance

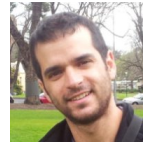
## Joint work with



Aviad Aberdam  
*Technion*



Amir Beck  
*Tel Aviv University*



Yaniv Romano  
*Stanford*



Vardan Papyan  
*Stanford*



Miki Elad  
*Technion*

*“Essentially, all models are wrong, but some are useful”*

George E.P. Box