

# Model-driven deep learning

## —— differentiable programming

# Background



Yann LeCun

8 hrs · 🌐

OK, Deep Learning has outlived its usefulness as a buzz-phrase.

Deep Learning est mort. Vive Differentiable Programming!

Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers.

But the important point is that people are now building a new kind of software by assembling networks of parameterized functional blocks and by training them from examples using some form of gradient-based optimization.

An increasingly large number of people are defining the network procedurally in a data-dependant way (with loops and conditionals), allowing them to change dynamically as a function of the input data fed to them. It's really very much like a regular program, except it's parameterized, automatically differentiated, and trainable/optimizable. Dynamic networks have become increasingly popular (particularly for NLP), thanks to deep learning frameworks that can handle them such as PyTorch and Chainer (note: our old deep learning framework Lush could handle a particular kind of dynamic nets called Graph Transformer Networks, back in 1994. It was needed for text recognition).

People are now actively working on compilers for imperative differentiable programming languages. This is a very exciting avenue for the development of learning-based AI.

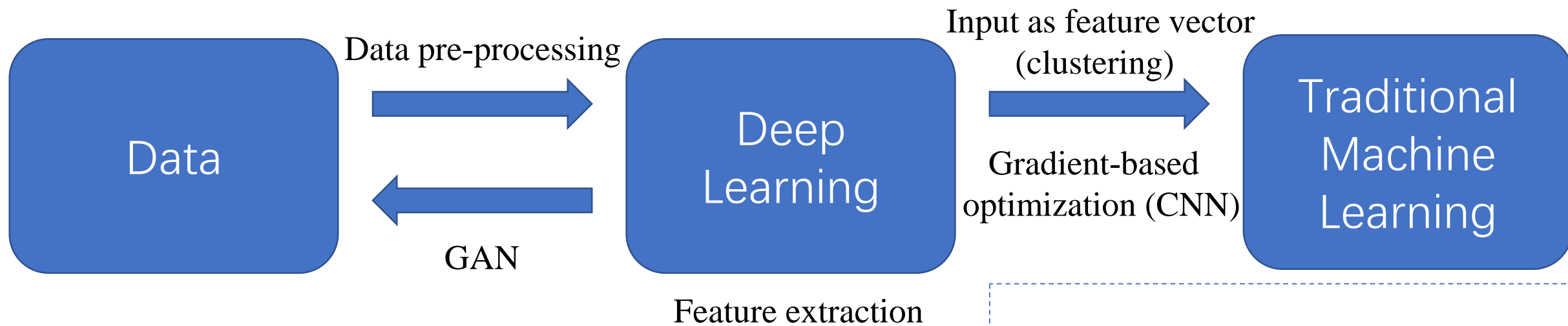
Important note: this won't be sufficient to take us to "true" AI. Other concepts will be needed for that, such as what I used to call predictive learning and now decided to call Imputative Learning. More on this later....

## Model-driven deep learning (differentiable programming):

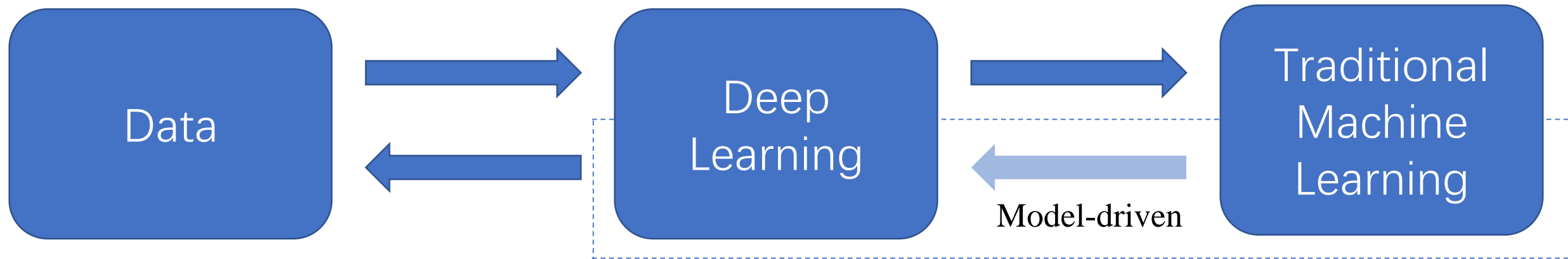
treats neural network as a language

- Deep learning: a standard **data-driven** method, which uses networks as a black box to rely on large amounts of data to solve problems
- Model-driven method: forms a cost function from **the goal, mechanism and a prior**, and then solve problems by minimizing the cost function

Black box, only data-driven



Model-driven deep learning — try to transform traditional machine learning methods into equivalent neural networks



- Strong **interpretability** of traditional methods
- Superior **fitting capabilities** of deep neural networks

The general procedure of model-driven deep learning<sup>[1]</sup>:



- Construct a **model family** based on the task backgrounds
- Design an **algorithm family** for solving the model family, and establish the corresponding convergence theory
- The **algorithm family is unfolded to a deep network** with which parameter learning is performed as a deep learning approach

---

# Learning Fast Approximations of Sparse Coding

---

Karol Gregor and Yann LeCun

{KGREGOR,YANN}@CS.NYU.EDU

Courant Institute, New York University, 715 Broadway, New York, NY 10003, USA

**ICML 2010**

## Sparse coding:

For a given input vector  $\mathbf{X} \in R^n$ , to find the optimal sparse code vector  $\mathbf{Z}^* \in R^m$  ( $m > n$ ) that minimizes an energy function that combines the **square reconstruction error** and an  **$l_1$  sparsity penalty** on the code:

$$\mathbf{Z}^* = \operatorname{argmin}_{\mathbf{Z}} \frac{1}{2} \|\mathbf{X} - \mathbf{W}_d \mathbf{Z}\|_2^2 + \alpha \|\mathbf{Z}\|_1$$

where  $\mathbf{W}_d$  denotes an  $n \times m$  dictionary matrix.



## ISTA algorithm<sup>[1]</sup>:

---

### Algorithm 1 ISTA

---

function **ISTA**( $X, Z, W_d, \alpha, L$ )

**Require:**  $L >$  largest eigenvalue of  $W_d^T W_d$ .

**Initialize:**  $Z = 0$ ,

**repeat**

$$Z = h_{(\alpha/L)}\left(Z - \frac{1}{L}W_d^T(W_d Z - X)\right)$$

**until** change in  $Z$  below a threshold

**end function**

---

$$\mathbf{Z}^{(k+1)} = h_{\theta}(\mathbf{W}_e \mathbf{X} + \mathbf{S} \mathbf{Z}^{(k)})$$

$$\mathbf{W}_e = \frac{1}{L} \mathbf{W}_d^T \quad \theta = \alpha/L$$

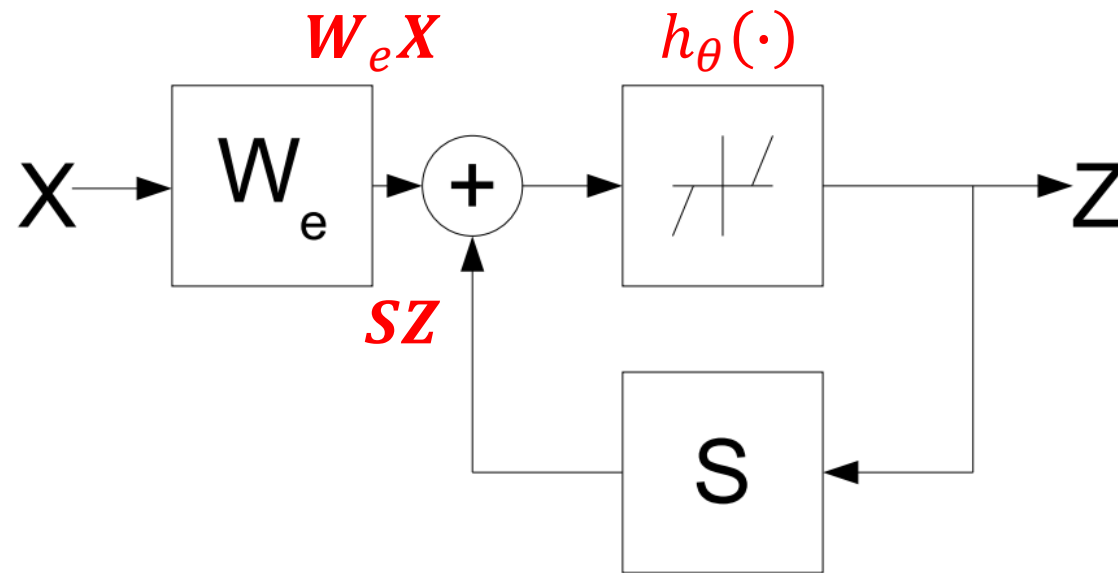
$$\mathbf{S} = \mathbf{I} - \frac{1}{L} \mathbf{W}_d^T \mathbf{W}_d$$

$$[h_{\theta}(\mathbf{V})]_i = \text{sign}(V_i) (|V_i| - \theta_i)_+$$

[1] I. Daubechies, M. Defrise, and C. D. Mol, An iterative thresholding algorithm for linear inverse problems with a sparsity constraint, *Comm. Pure Appl. Math.*, 57 (2004), pp. 1413–1457.

ISTA block diagram:

$$\mathbf{z}^{(k+1)} = h_{\theta}(\mathbf{W}_e \mathbf{X} + \mathbf{S} \mathbf{z}^{(k)})$$

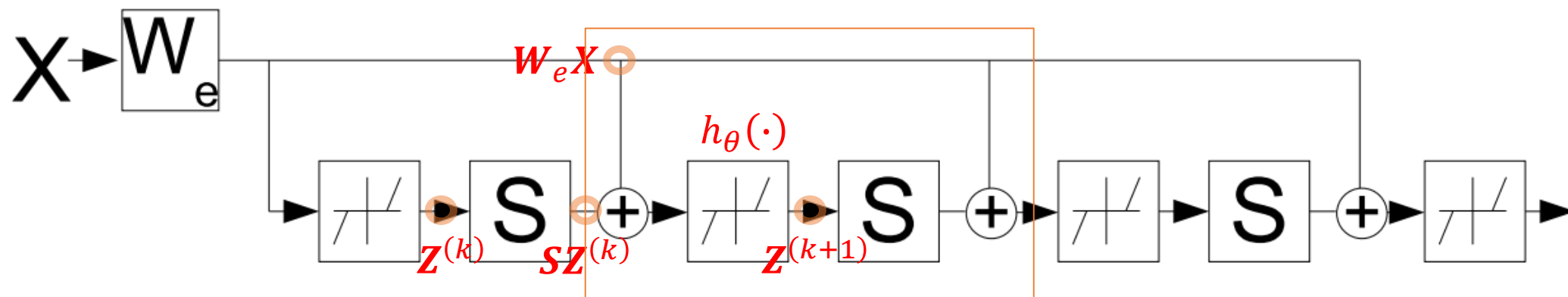


How to transform the ISTA to a neural network?

LISTA:

$$\mathbf{z}^{(k+1)} = h_{\theta}(\mathbf{W}_e \mathbf{X} + \mathbf{S} \mathbf{z}^{(k)})$$

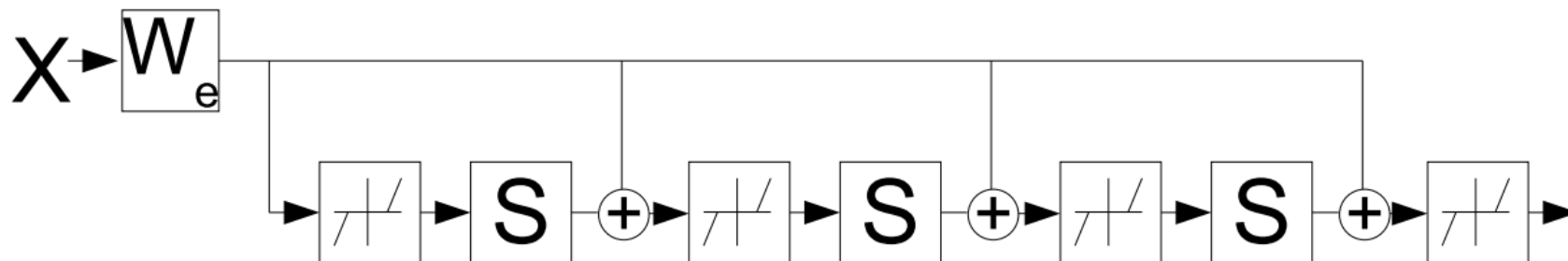
Time-unfolded version of the ISTA block diagram, truncated to a fixed number of iterations



Learnable parameters:  $W_e$ ,  $S$  and  $\theta$

LISTA:

$$\mathbf{z}^{(k+1)} = h_{\theta}(\mathbf{W}_e \mathbf{X} + \mathbf{S} \mathbf{z}^{(k)})$$



Let  $\mathbf{W} = (\mathbf{W}_e, \mathbf{S}, \theta)$ , and the LISTA as  $\mathbf{Z} = f_e(\mathbf{X}, \mathbf{W})$ .

Loss function: the squared error between the predicted code and the optimal code averaged over a training set  $(\mathbf{X}^1, \dots, \mathbf{X}^P)$

$$L(\mathbf{W}) = \frac{1}{P} \sum_{p=1}^P L(\mathbf{W}, \mathbf{X}^p) \quad \text{with}$$

$$L(\mathbf{W}, \mathbf{X}^p) = \frac{1}{2} \|\mathbf{z}^{*p} - f_e(\mathbf{W}, \mathbf{X}^p)\|_2^2$$

## Numerical results:

*Table 1.* Prediction error (squared error between the optimal code and the predicted codes) for different nonlinearities of the baseline encoder (7), and for LISTA and FISTA. LISTA produces a much better estimate after one iteration than FISTA.

NON-LINEARITY	100 UNITS	400 UNITS
$D \tanh(x)$	8.6	10.7
$D(\tanh(x + u) + \tanh(x - u))$	3.33	4.62
$h_\alpha(x)$	3.29	4.82
LISTA 1 ITERATION	1.50	2.45
LISTA 3 ITERATIONS	0.98	2.12
LISTA 7 ITERATIONS	0.52	1.62
FISTA 1 ITERATION	21.	22.

Baseline encoder:

$$Z = g(W_e X) \quad (7)$$

More interestingly LISTA is stupendously better than FISTA: It takes 18 iterations of FISTA to reach the same error obtained with 1 iteration of LISTA for  $m = 100$ , and 35 iteration for  $m = 400$ .

---

# Deep ADMM-Net for Compressive Sensing MRI

---

**Yan Yang**

Xi'an Jiaotong University  
yangyan92@stu.xjtu.edu.cn

**Jian Sun**

Xi'an Jiaotong University  
jiansun@mail.xjtu.edu.cn

**Huibo Li**

Xi'an Jiaotong University  
huibinli@mail.xjtu.edu.cn

**Zongben Xu**

Xi'an Jiaotong University  
zbxu@mail.xjtu.edu.cn

**NIPS 2016**

## Compressive sensing magnetic resonance imaging (MRI):

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{D}_l \mathbf{x}) \right\}$$

$\mathbf{x} \in C^N$ : the MRI image to be reconstructed

$\mathbf{y} \in C^{N'}$  ( $N' < N$ ): the under-sampled k-space data

$\mathbf{A} = \mathbf{P}\mathbf{F} \in C^{N' \times N}$ : the measurement matrix ( $\mathbf{P}$  the under-sampling matrix, and  $\mathbf{F}$  the Fourier transform matrix)

$\mathbf{D}_l$ : the transform matrix for a filtering operation, e.g., discrete wavelet transform, discrete cosine transform, etc.

$g(\cdot)$ : the regularization function derived from the data prior, e.g.,  $l_1$ -norm,  $l_p$ -norm

$\lambda_l$ : the regularization parameter

General CS-MRI model:

$$\hat{\mathbf{x}} = \operatorname{argmin}_x \left\{ \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{D}_l \mathbf{x}) \right\}$$

Challenge

- How to choose an optimal image transform domain/subspace and the corresponding sparse regularization?
- How to determine the optimal parameters?



ADMM algorithm<sup>[1]</sup>: 
$$\hat{\mathbf{x}} = \operatorname{argmin}_x \left\{ \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{D}_l \mathbf{x}) \right\}$$

Introducing auxiliary variables  $\mathbf{z} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_L\}$

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{z}} \quad & \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{z}_l) \\ \text{s.t.} \quad & \mathbf{z}_l = \mathbf{D}_l \mathbf{x}, \quad \forall l \in [1, 2, \dots, L] \end{aligned}$$

$\boldsymbol{\alpha} = \{\boldsymbol{\alpha}_l\}$  the Lagrangian multipliers

$\rho = \{\rho_l\}$  the penalty parameters

Its augmented Lagrangian function is

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{z}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{z}_l) - \sum_{l=1}^L \langle \boldsymbol{\alpha}_l, \mathbf{z}_l - \mathbf{D}_l \mathbf{x} \rangle + \sum_{l=1}^L \frac{\rho_l}{2} \|\mathbf{z}_l - \mathbf{D}_l \mathbf{x}\|_2^2$$

[1] Boyd S, Parikh N, Chu E. Distributed optimization and statistical learning via the alternating direction method of multipliers [M]. Now Publishers Inc, 2011.

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{z}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{z}_l) - \sum_{l=1}^L \langle \boldsymbol{\alpha}_l, \mathbf{z}_l - \mathbf{D}_l \mathbf{x} \rangle + \sum_{l=1}^L \frac{\rho_l}{2} \|\mathbf{z}_l - \mathbf{D}_l \mathbf{x}\|_2^2$$

ADMM alternatively optimizes  $\{\mathbf{x}, \mathbf{z}, \boldsymbol{\alpha}\}$  by solving the three subproblems

$$\begin{cases} \mathbf{x}^{(n+1)} = \operatorname{argmin}_x \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 - \sum_{l=1}^L \langle \boldsymbol{\alpha}_l^{(n)}, \mathbf{z}_l^{(n)} - \mathbf{D}_l \mathbf{x} \rangle + \sum_{l=1}^L \frac{\rho_l}{2} \|\mathbf{z}_l^{(n)} - \mathbf{D}_l \mathbf{x}\|_2^2 \\ \mathbf{z}^{(n+1)} = \operatorname{argmin}_z \sum_{l=1}^L \lambda_l g(\mathbf{z}_l) - \sum_{l=1}^L \langle \boldsymbol{\alpha}_l^{(n)}, \mathbf{z}_l - \mathbf{D}_l \mathbf{x}^{(n+1)} \rangle + \sum_{l=1}^L \frac{\rho_l}{2} \|\mathbf{z}_l - \mathbf{D}_l \mathbf{x}^{(n+1)}\|_2^2 \\ \boldsymbol{\alpha}^{(n+1)} = \operatorname{argmin}_\alpha \sum_{l=1}^L \langle \boldsymbol{\alpha}_l, \mathbf{D}_l \mathbf{x}^{(n+1)} - \mathbf{z}_l^{(n+1)} \rangle \end{cases}$$

ADMM procedure: 
$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{D}_l \mathbf{x}) \right\}$$

$$\begin{cases} \mathbf{X}^{(n)}: \mathbf{x}^{(n)} = \mathbf{F}^T \left[ \mathbf{P}^T \mathbf{P} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T \mathbf{D}_l \mathbf{F}^T \right]^{-1} \left[ \mathbf{P}^T \mathbf{y} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T \left( \mathbf{z}_l^{(n-1)} - \boldsymbol{\beta}_l^{(n-1)} \right) \right] \\ \mathbf{Z}^{(n)}: \mathbf{z}_l^{(n)} = S \left( \mathbf{D}_l \mathbf{x}^{(n)} + \boldsymbol{\beta}_l^{(n-1)}; \lambda_l / \rho_l \right) \\ \mathbf{M}^{(n)}: \boldsymbol{\beta}_l^{(n)} = \boldsymbol{\beta}_l^{(n-1)} + \eta_l \left( \mathbf{D}_l \mathbf{x}^{(n)} - \mathbf{z}_l^{(n)} \right) \end{cases}$$

$\boldsymbol{\beta}_l = \boldsymbol{\alpha}_l / \rho_l$ : the scaled Lagrangian multiplier

$S(\cdot)$ : the nonlinear shrinkage function, e.g., soft or hard thresholding function corresponding to the sparse regularization of  $l_1$ -norm and  $l_0$ -norm, respectively

$\eta_l$ : the update rate

ADMM procedure: 
$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{D}_l \mathbf{x}) \right\}$$

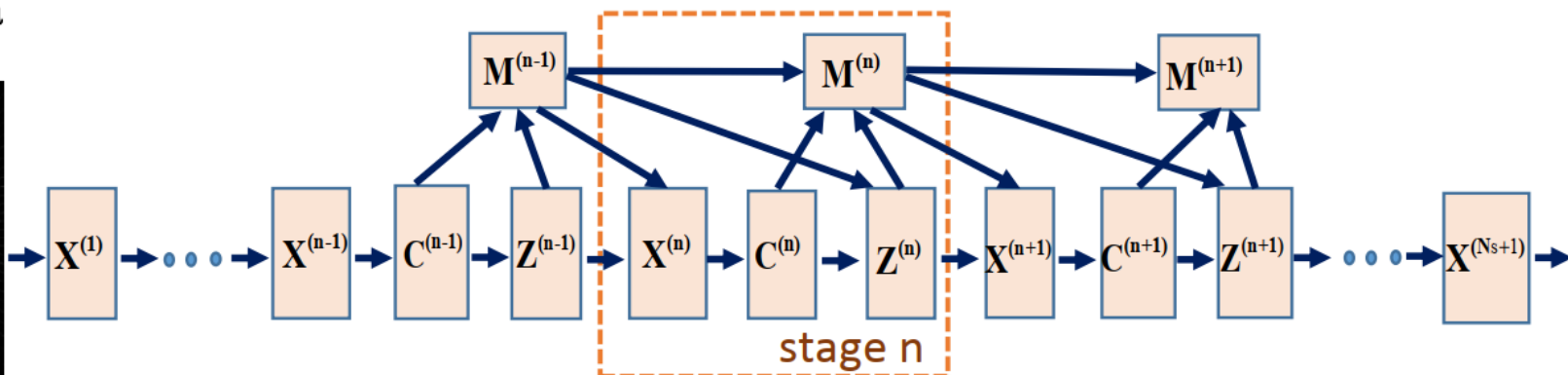
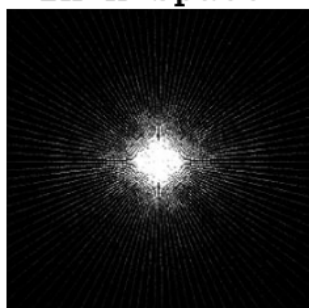
$$\begin{cases} \mathbf{X}^{(n)}: \mathbf{x}^{(n)} = \mathbf{F}^T \left[ \mathbf{P}^T \mathbf{P} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T \mathbf{D}_l \mathbf{F}^T \right]^{-1} \left[ \mathbf{P}^T \mathbf{y} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T \left( \mathbf{z}_l^{(n-1)} - \boldsymbol{\beta}_l^{(n-1)} \right) \right] \\ \mathbf{Z}^{(n)}: \mathbf{z}_l^{(n)} = S \left( \mathbf{D}_l \mathbf{x}^{(n)} + \boldsymbol{\beta}_l^{(n-1)}; \lambda_l / \rho_l \right) \\ \mathbf{M}^{(n)}: \boldsymbol{\beta}_l^{(n)} = \boldsymbol{\beta}_l^{(n-1)} + \eta_l \left( \mathbf{D}_l \mathbf{x}^{(n)} - \mathbf{z}_l^{(n)} \right) \end{cases}$$

- It is challenging to choose the transform  $\mathbf{D}_l$  and shrinkage function  $S(\cdot)$  for general regularization function  $g(\cdot)$
- It is also not trivial to tune the parameters  $\rho_l$  and  $\eta_l$  for k-space data with different sampling ratios

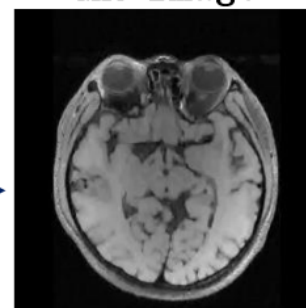
## Deep ADMM-Net:

$$\begin{cases} \mathbf{X}^{(n)}: \mathbf{x}^{(n)} = \mathbf{F}^T \left[ \mathbf{P}^T \mathbf{P} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T \mathbf{D}_l \mathbf{F}^T \right]^{-1} \left[ \mathbf{P}^T \mathbf{y} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T (\mathbf{z}_l^{(n-1)} - \boldsymbol{\beta}_l^{(n-1)}) \right] \\ \mathbf{Z}^{(n)}: \mathbf{z}_l^{(n)} = S(\mathbf{D}_l \mathbf{x}^{(n)} + \boldsymbol{\beta}_l^{(n-1)}; \lambda_l / \rho_l) \\ \mathbf{M}^{(n)}: \boldsymbol{\beta}_l^{(n)} = \boldsymbol{\beta}_l^{(n-1)} + \eta_l (\mathbf{D}_l \mathbf{x}^{(n)} - \mathbf{z}_l^{(n)}) \end{cases}$$

Sampling data  
in k-space



Reconstructed  
MR image



- Reconstruction operation  $\mathbf{X}^{(n)}$
- Convolution operation  $\mathbf{C}^{(n)}$  defined by  $\{\mathbf{D}_l \mathbf{x}^{(n)}\}_{l=1}^L$
- Nonlinear transform operation  $\mathbf{Z}^{(n)}$  defined by  $S(\cdot)$
- Multiplier update operation  $\mathbf{M}^{(n)}$

Deep ADMM-Net — Reconstruction layer  $\mathbf{X}^{(n)}$ :

$$\mathbf{x}^{(n)} = \mathbf{F}^T \left[ \mathbf{P}^T \mathbf{P} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{H}_l^T \mathbf{H}_l \mathbf{F}^T \right]^{-1} \left[ \mathbf{P}^T \mathbf{y} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{H}_l^T \left( \mathbf{z}_l^{(n-1)} - \boldsymbol{\beta}_l^{(n-1)} \right) \right]$$

$\mathbf{H}_l$ : the  $l$ -th filter

First stage update:

$$\mathbf{x}^{(1)} = \mathbf{F}^T \left( \mathbf{P}^T \mathbf{P} + \sum_{l=1}^L \rho_l^{(1)} \mathbf{F} \mathbf{H}_l^{(1)T} \mathbf{H}_l^{(1)} \mathbf{F}^T \right)^{-1} (\mathbf{P}^T \mathbf{y})$$

---

ADMM update formula:  $\mathbf{x}^{(n)} = \mathbf{F}^T \left[ \mathbf{P}^T \mathbf{P} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T \mathbf{D}_l \mathbf{F}^T \right]^{-1} \left[ \mathbf{P}^T \mathbf{y} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{D}_l^T \left( \mathbf{z}_l^{(n-1)} - \boldsymbol{\beta}_l^{(n-1)} \right) \right]$

## Deep ADMM-Net — Convolution layer $\mathbf{C}^{(n)}$ :

Transform an image into transform domain

$$\mathbf{c}_l^{(n)} = \mathbf{D}_l^{(n)} \mathbf{x}^{(n)}$$

$\mathbf{D}_l$ : a learnable filter matrix in stage  $n$

- Note that the filters  $\mathbf{H}_l$  and  $\mathbf{D}_l$  do not be constrained to be the same to increase the network capacity

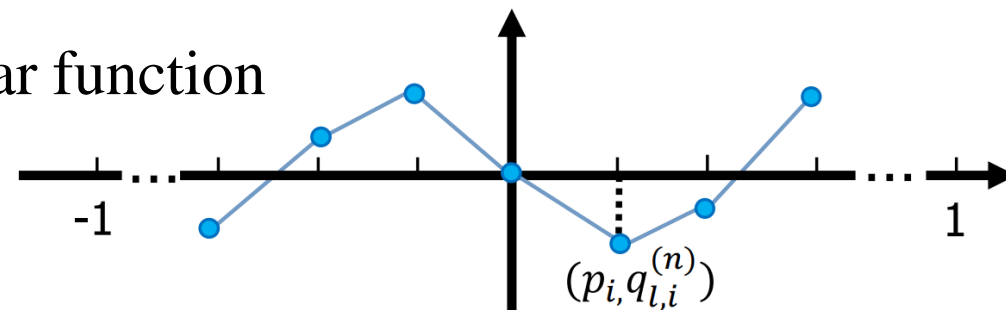
---

$$\text{ADMM update formula: } \mathbf{z}_l^{(n)} = S \left( \mathbf{D}_l \mathbf{x}^{(n)} + \boldsymbol{\beta}_l^{(n-1)}; \lambda_l / \rho_l \right)$$

## Deep ADMM-Net — Nonlinear transform layer $\mathbf{Z}^{(n)}$ :

Learn more general function using piecewise linear function

$$\mathbf{z}_l^{(n)} = S_{PLF} \left( \mathbf{c}_l^{(n)} + \boldsymbol{\beta}_l^{(n-1)}; \{p_i, q_{l,i}^{(n)}\}_{i=1}^{N_c} \right)$$



$$S_{PLF} \left( a; \{p_i, q_{l,i}^{(n)}\}_{i=1}^{N_c} \right) = \begin{cases} a + q_{l,1}^{(n)} - p_1, & a < p_1, \\ a + q_{l,N_c}^{(n)} - p_{N_c}, & a > p_{N_c}, \\ q_{l,k}^{(n)} + \frac{(a - p_k)(q_{l,k+1}^{(n)} - q_{l,k}^{(n)})}{p_{k+1} - p_k}, & p_1 \leq a \leq p_{N_c}, \end{cases}$$

$$k = \left\lfloor \frac{a - p_1}{p_2 - p_1} \right\rfloor$$

$\{p_i\}_{i=1}^{N_c}$ : the predefined positions uniformly located within  $[-1, 1]$

$\{q_{l,i}^{(n)}\}_{i=1}^{N_c}$ : the values at these positions for  $l$ -th filter in  $n$ -th stage

---


$$\text{ADMM update formula: } \mathbf{z}_l^{(n)} = S \left( \mathbf{D}_l \mathbf{x}^{(n)} + \boldsymbol{\beta}_l^{(n-1)}; \lambda_l / \rho_l \right)$$



## Deep ADMM-Net — Multiplier update layer $\mathbf{M}^{(n)}$ :

Lagrangian multiplier updating procedure

$$\boldsymbol{\beta}_l^{(n)} = \boldsymbol{\beta}_l^{(n-1)} + \eta_l^{(n)} \left( \mathbf{c}_l^{(n)} - \mathbf{z}_l^{(n)} \right)$$

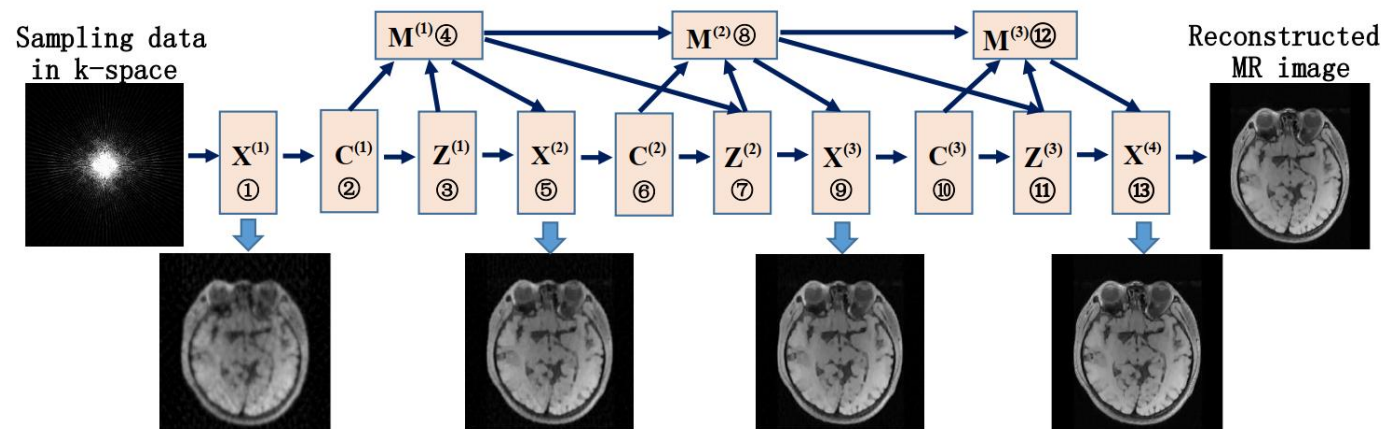
$\eta_l^{(n)}$  : a learnable parameter

---

ADMM update formula:  $\boldsymbol{\beta}_l^{(n)} = \boldsymbol{\beta}_l^{(n-1)} + \eta_l \left( \mathbf{D}_l \mathbf{x}^{(n)} - \mathbf{z}_l^{(n)} \right)$

## Network parameters and training:

$$E(\Theta) = \frac{1}{|\Gamma|} \sum_{(y, x^{gt}) \in \Gamma} \frac{\sqrt{\|\hat{x}(y, \Theta) - x^{gt}\|_2^2}}{\sqrt{\|x^{gt}\|_2^2}}$$



- Reconstruction layer  $X^{(n)}$ :  $\mathbf{x}^{(n)} = \mathbf{F}^T \left[ \mathbf{P}^T \mathbf{P} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{H}_l^T \mathbf{H}_l \mathbf{F}^T \right]^{-1} \left[ \mathbf{P}^T \mathbf{y} + \sum_{l=1}^L \rho_l \mathbf{F} \mathbf{H}_l^T \left( \mathbf{z}_l^{(n-1)} - \boldsymbol{\beta}_l^{(n-1)} \right) \right]$
- Convolution layer  $C^{(n)}$ :  $\mathbf{c}_l^{(n)} = \mathbf{D}_l^{(n)} \mathbf{x}^{(n)}$
- Nonlinear transform layer  $Z^{(n)}$ :  $\mathbf{z}_l^{(n)} = S_{PLF} \left( \mathbf{c}_l^{(n)} + \boldsymbol{\beta}_l^{(n-1)}; \left\{ p_i, q_{l,i}^{(n)} \right\}_{i=1}^{N_c} \right)$
- Multiplier update layer  $M^{(n)}$ :  $\boldsymbol{\beta}_l^{(n)} = \boldsymbol{\beta}_l^{(n-1)} + \eta_l^{(n)} \left( \mathbf{c}_l^{(n)} - \mathbf{z}_l^{(n)} \right)$

Table 1: Performance comparisons on brain data with different sampling ratios.

Method	20%		30%		40%		50%		Test time
	NMSE	PSNR	NMSE	PSNR	NMSE	PSNR	NMSE	PSNR	
Zero-filling	0.1700	29.96	0.1247	32.59	0.0968	34.76	0.0770	36.73	0.0013s
TV [2]	0.0929	35.20	0.0673	37.99	0.0534	40.00	0.0440	41.69	0.7391s
RecPF [4]	0.0917	35.32	0.0668	38.06	0.0533	40.03	0.0440	41.71	0.3105s
SIDWT	0.0885	35.66	0.0620	38.72	0.0484	40.88	0.0393	42.67	7.8637s
PBDW [6]	0.0814	36.34	0.0627	38.64	0.0518	40.31	0.0437	41.81	35.3637s
PANO [10]	0.0800	36.52	0.0592	39.13	0.0477	41.01	0.0390	42.76	53.4776s
FDLCP [8]	0.0759	36.95	0.0592	39.13	0.0500	40.62	0.0428	42.00	52.2220s
BM3D-MRI [11]	0.0674	37.98	0.0515	40.33	0.0426	41.99	0.0359	43.47	40.9114s
Init-Net <sub>13</sub>	0.1394	31.58	0.1225	32.71	0.1128	33.44	0.1066	33.95	0.6914s
ADMM-Net <sub>13</sub>	0.0752	37.01	0.0553	39.70	0.0456	41.37	0.0395	42.62	0.6964s
ADMM-Net <sub>14</sub>	0.0742	37.13	0.0548	39.78	0.0448	41.54	0.0380	42.99	0.7400s
ADMM-Net <sub>15</sub>	0.0739	37.17	0.0544	39.84	0.0447	41.56	0.0379	43.00	0.7911s

Table 2: Comparisons of NMSE and PSNR on chest data with 20% sampling ratio.

Method	TV	RecPF	PANO	FDLCP	ADMM-Net <sub>15</sub> -B	ADMM-Net <sub>15</sub>	ADMM-Net <sub>17</sub>
NMSE	0.1019	0.1017	0.0858	0.0775	0.0790	0.0775	0.0768
PSNR	35.49	35.51	37.01	37.77	37.68	37.84	37.92

Expansion:

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{D}_l \mathbf{x}) \right\}$$

ADMM-Net: introduce auxiliary variables  $\mathbf{z} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_L\}$

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{z}} \quad & \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{z}_l) \\ \text{s.t.} \quad & \mathbf{z}_l = \mathbf{D}_l \mathbf{x}, \quad \forall l \in [1, 2, \dots, L] \end{aligned}$$

ADMM-CSNet<sup>[1]</sup>: introduce auxiliary variable  $\mathbf{z} = \mathbf{x}$

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{z}} \quad & \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \sum_{l=1}^L \lambda_l g(\mathbf{D}_l \mathbf{z}) \\ \text{s.t.} \quad & \mathbf{z} = \mathbf{x} \end{aligned}$$

[1] Yang Y, Sun J, Li H, et al. ADMM-CSNet: A deep learning approach for image compressive sensing [J]. IEEE transactions on pattern analysis and machine intelligence, 2018, 42(3): 521-538.

# **ISTA-Net: Interpretable Optimization-Inspired Deep Network for Image Compressive Sensing**

Jian Zhang, Bernard Ghanem

King Abdullah University of Science and Technology (KAUST), Saudi Arabia

`jian.zhang@kaust.edu.sa`, `bernard.ghanem@kaust.edu.sa`

**CVPR 2018**

## Contributions:

- A novel ISTA-Net, which **adopts the structure of ISTA** update steps to design a **learnable deep network** manifestation
- By **exploiting CS realm in the residual domain**, an enhanced version **ISTA-Net<sup>+</sup>** is derived to further improve CS performance
- The **proximal mapping problem** associated to a nonlinear sparsifying transform is solved in a general and efficient way

CS reconstruction problem:

$$\min_x \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\Psi \mathbf{x}\|_1$$

ISTA algorithm<sup>[1]</sup> procedure:

$$\text{prox}_{\lambda\phi}(\mathbf{r}) = \underset{\mathbf{x}}{\text{argmin}} \frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_2^2 + \lambda\phi(\mathbf{x})$$

$$\mathbf{r}^{(k)} = \mathbf{x}^{(k-1)} - \rho \Phi^\top (\Phi \mathbf{x}^{(k-1)} - \mathbf{y})$$

$$\mathbf{x}^{(k)} = \underset{\mathbf{x}}{\text{argmin}} \frac{1}{2} \|\mathbf{x} - \mathbf{r}^{(k)}\|_2^2 + \lambda \|\Psi \mathbf{x}\|_1$$

[1] I. Daubechies, M. Defrise, and C. D. Mol, An iterative thresholding algorithm for linear inverse problems with a sparsity constraint, *Comm. Pure Appl. Math.*, 57 (2004), pp. 1413–1457.

$$\mathbf{x}^{(k)} = \arg \min_x \frac{1}{2} \|\mathbf{x} - \mathbf{r}^{(k)}\|_2^2 + \lambda \|\Psi \mathbf{x}\|_1$$

- Inspired by the powerful representation power of CNN and its universal approximation property

$$\mathcal{F}(\mathbf{x}) = \Psi \mathbf{x}$$

$\mathcal{F}(\cdot)$  as a combination of two linear convolutional operators separated by a ReLU function

- Then

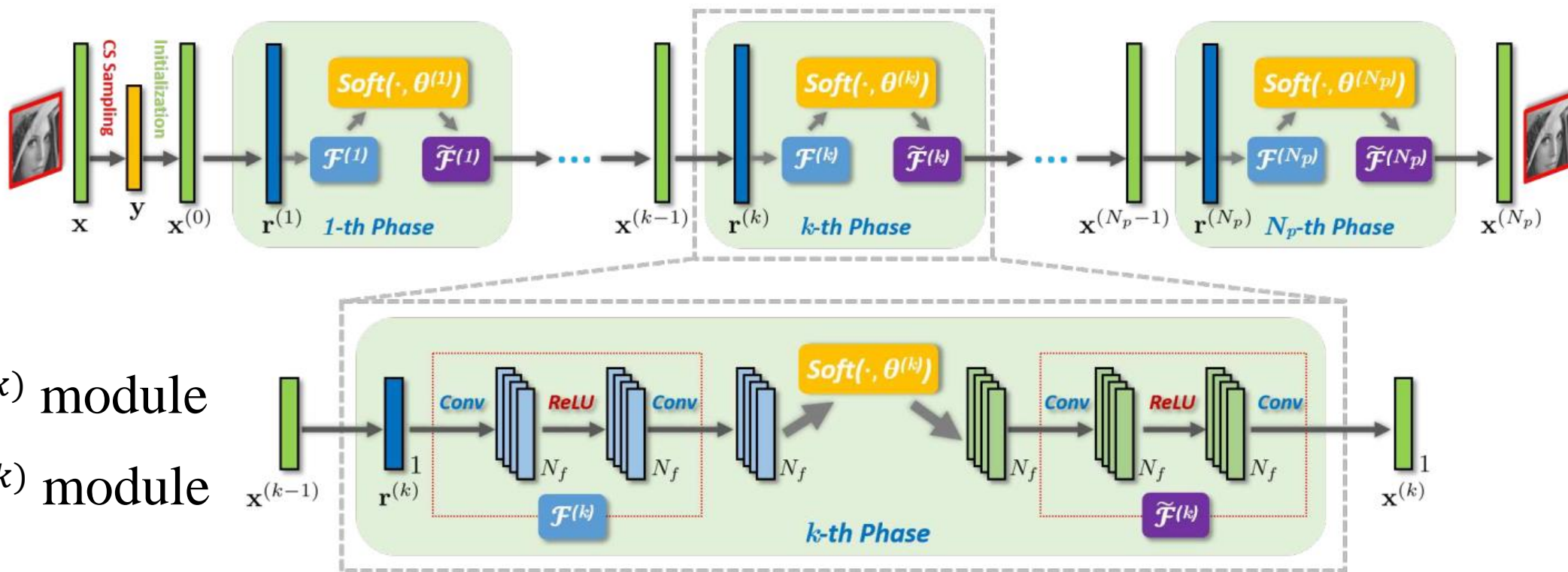
$$\mathbf{x}^{(k)} = \arg \min_x \frac{1}{2} \|\mathbf{x} - \mathbf{r}^{(k)}\|_2^2 + \lambda \|\mathcal{F}(\mathbf{x})\|_1$$



## ISTA-Net framework:

$$\mathbf{r}^{(k)} = \mathbf{x}^{(k-1)} - \rho \Phi^T (\Phi \mathbf{x}^{(k-1)} - \mathbf{y})$$

$$\mathbf{x}^{(k)} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{r}^{(k)}\|_2^2 + \lambda \|\mathcal{F}(\mathbf{x})\|_1$$



- $\mathbf{r}^{(k)}$  module
- $\mathbf{x}^{(k)}$  module

ISTA-Net —  $\mathbf{r}^{(k)}$  Module:

$$\mathbf{r}^{(k)} = \mathbf{x}^{(k-1)} - \rho \mathbf{\Phi}^\top (\mathbf{\Phi} \mathbf{x}^{(k-1)} - \mathbf{y})$$

To preserve the ISTA structure while increasing network flexibility, the step size  $\rho$  is allowed to vary across iterations

$$\mathbf{r}^{(k)} = \mathbf{x}^{(k-1)} - \rho^{(k)} \mathbf{\Phi}^\top (\mathbf{\Phi} \mathbf{x}^{(k-1)} - \mathbf{y})$$

ISTA-Net —  $\mathbf{x}^{(k)}$  Module:

$$\mathbf{x}^{(k)} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{r}^{(k)}\|_2^2 + \lambda \|\mathcal{F}(\mathbf{x})\|_1$$

- Note that  $\mathbf{r}^{(k)}$  is the immediate reconstruction result of  $\mathbf{x}^{(k)}$
- Thus,  $\mathbf{x}^{(k)} - \mathbf{r}^{(k)}$  follows an independent normal distribution with common zero mean and variance  $\sigma^2$

- Approximation (Theorem 1):

$$\|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{r}^{(k)})\|_2^2 \approx \alpha \|\mathbf{x} - \mathbf{r}^{(k)}\|_2^2$$

- Equivalent transformation

$$\mathbf{x}^{(k)} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{r}^{(k)})\|_2^2 + \theta \|\mathcal{F}(\mathbf{x})\|_1$$

**Theorem 1** *Let  $X_1, \dots, X_n$  be independent normal random variables with common zero mean and variance  $\sigma^2$ . If  $\vec{X} = [X_1, \dots, X_n]^\top$  and given any matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{s \times m}$ , define a new random variable  $\vec{Y} = \mathbf{B} \text{ReLU}(\mathbf{A}\vec{X}) = \mathbf{B} \max(\mathbf{0}, \mathbf{A}\vec{X})$ . Then,  $\mathbb{E}[\|\vec{Y} - \mathbb{E}[\vec{Y}]\|_2^2]$  and  $\mathbb{E}[\|\vec{X} - \mathbb{E}[\vec{X}]\|_2^2]$  are linearly related, i.e.  $\mathbb{E}[\|\vec{Y} - \mathbb{E}[\vec{Y}]\|_2^2] = \alpha \mathbb{E}[\|\vec{X} - \mathbb{E}[\vec{X}]\|_2^2]$ , where  $\alpha$  is only a function of  $\mathbf{A}$  and  $\mathbf{B}$ . (Please refer to the **supplementary material** for the proof and more details.)*

ISTA-Net —  $\mathbf{x}^{(k)}$  Module: 
$$\mathbf{x}^{(k)} = \arg \min_x \frac{1}{2} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{r}^{(k)})\|_2^2 + \theta \|\mathcal{F}(\mathbf{x})\|_1$$

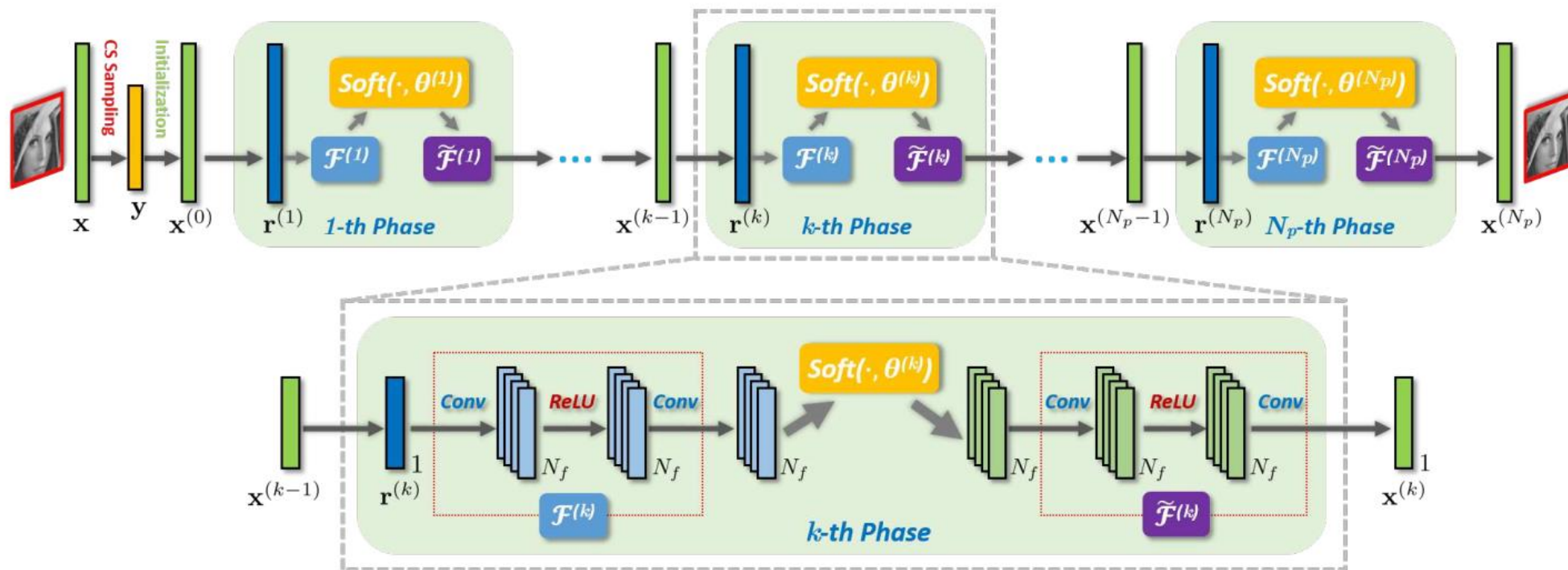
- Then, we get a closed-form version of  $\mathcal{F}(\mathbf{x}^{(k)})$

$$\mathcal{F}(\mathbf{x}^{(k)}) = \text{soft}(\mathcal{F}(\mathbf{r}^{(k)}), \theta)$$

- Motivated by the invertible characteristics of the wavelet transform, the left inverse of  $\mathcal{F}(\cdot)$  may be introduced, such that  $\tilde{\mathcal{F}} \circ \mathcal{F} = I$ , yielding

$$\mathbf{x}^{(k)} = \tilde{\mathcal{F}}(\text{soft}(\mathcal{F}(\mathbf{r}^{(k)}), \theta))$$

## Parameters in ISTA-Net:



- $r^{(k)}$  module:  $r^{(k)} = x^{(k-1)} - \rho^{(k)} \Phi^T (\Phi x^{(k-1)} - y)$
- $x^{(k)}$  module:  $x^{(k)} = \tilde{\mathcal{F}}(\text{soft}(\mathcal{F}(r^{(k)}), \theta))$

## Loss function design:

- The training data pairs  $\{(\mathbf{y}_i, \mathbf{x}_i)\}_{i=1}^{N_b}$ , and  $\gamma = 0.01$

$$\mathcal{L}_{total}(\Theta) = \mathcal{L}_{discrepancy} + \gamma \mathcal{L}_{constraint}$$

$$\text{with: } \begin{cases} \mathcal{L}_{discrepancy} = \frac{1}{N_b N} \sum_{i=1}^{N_b} \|\mathbf{x}_i^{(N_p)} - \mathbf{x}_i\|_2^2 & \text{Reconstruction error} \\ \mathcal{L}_{constraint} = \frac{1}{N_b N} \sum_{i=1}^{N_b} \sum_{k=1}^{N_p} \|\tilde{\mathcal{F}}^{(k)}(\mathcal{F}^{(k)}(\mathbf{x}_i)) - \mathbf{x}_i\|_2^2 & \text{Invertible error} \end{cases}$$

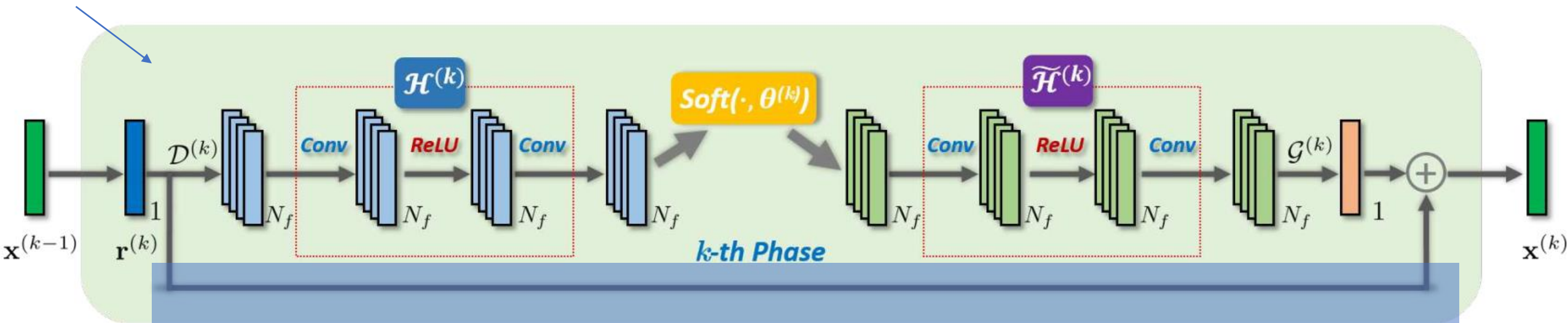
Enhanced version ISTA-Net<sup>+</sup>:  $\mathbf{x}^{(k)} = \arg \min_x \frac{1}{2} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{r}^{(k)})\|_2^2 + \theta \|\mathcal{F}(\mathbf{x})\|_1$

- Based on the fact that the residuals of natural images and videos are more compressible
- Assume that  $\mathbf{x}^{(k)} = \mathbf{r}^{(k)} + \mathbf{w}^{(k)} + \mathbf{e}^{(k)}$ , where  $\mathbf{e}^{(k)}$  denotes some noise
- $\mathbf{w}^{(k)}$  denotes some missing high-frequency component, which may be extracted by a linear operator  $\mathbf{w}^{(k)} = \mathcal{R}(\mathbf{x}^{(k)}) = \mathcal{G} \circ \mathcal{D}(\mathbf{x}^{(k)})$
- By modeling  $\mathcal{F}(\cdot) = \mathcal{H} \circ \mathcal{D}$

$$\min_x \frac{1}{2} \|\mathcal{H}(\mathcal{D}(\mathbf{x})) - \mathcal{H}(\mathcal{D}(\mathbf{r}^{(k)}))\|_2^2 + \theta \|\mathcal{H}(\mathcal{D}(\mathbf{x}))\|_1 \Rightarrow \mathbf{x}^{(k)} \Rightarrow \mathbf{x}^{(k)} = \mathbf{r}^{(k)} + \mathcal{R}(\mathbf{x}^{(k)})$$

$$\mathbf{x}^{(k)} = \mathbf{r}^{(k)} + \mathcal{G} \left( \tilde{\mathcal{H}} \left( \text{soft} \left( \mathcal{H} \left( \mathcal{D}(\mathbf{r}^{(k)}) \right), \theta \right) \right) \right)$$

## Parameters in ISTA-Net<sup>+</sup>:



skip connections (coincide with Res-Net)

- $\mathbf{r}^{(k)}$  module:  $\mathbf{r}^{(k)} = \mathbf{x}^{(k-1)} - \rho^{(k)} \Phi^T (\Phi \mathbf{x}^{(k-1)} - \mathbf{y})$
- $\mathbf{x}^{(k)}$  module:  $\mathbf{x}^{(k)} = \mathbf{r}^{(k)} + \mathcal{G} \left( \tilde{\mathcal{H}} \left( \text{soft} \left( \mathcal{H} \left( \mathcal{D}(\mathbf{r}^{(k)}) \right) \right) \right) \right)$



## ISTA-Net v.s. ISTA-Net<sup>+</sup>:

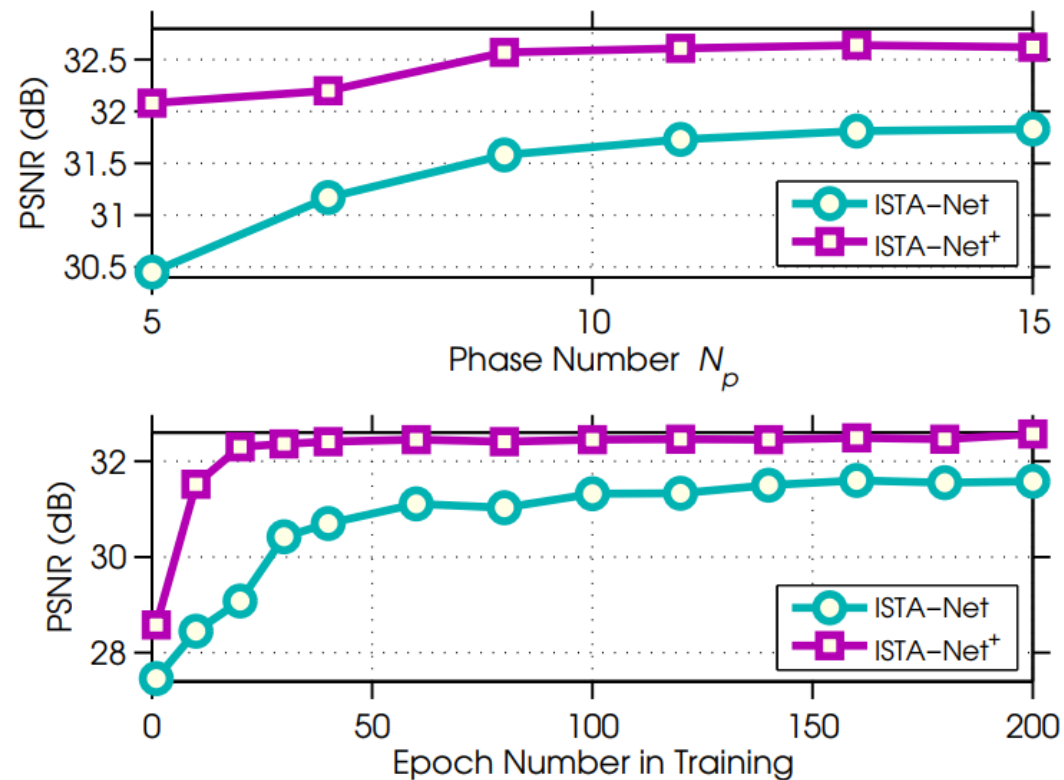


Figure 4. PSNR comparison between ISTA-Net and ISTA-Net<sup>+</sup> with various numbers of phases and epochs during training.

Table 1. Average PSNR (dB) performance comparisons on Set11 with different CS ratios. The best performance is labeled in bold and the second best is underlined. Note that the last two columns is a run-time analysis of all the competing methods, showing the average time to reconstruct a  $256 \times 256$  image and the corresponding frames-per-second (FPS).

Algorithm	CS Ratio							Time CPU/GPU	FPS CPU/GPU
	50%	40%	30%	25%	10%	4%	1%		
TVAL3 [22]	33.55	31.46	29.23	27.92	22.99	18.75	16.43	3.135s/——	0.32/——
D-AMP [28]	35.92	33.56	30.39	28.46	22.64	18.40	5.21	51.21s/——	0.02/——
IRCNN [49]	36.23	34.06	31.18	30.07	24.02	17.56	7.70	——/68.42s	——/0.015
SDA [30]	28.95	27.79	26.63	25.34	22.65	20.12	17.29	——/0.0032s	——/312.5
ReconNet [21]	31.50	30.58	28.74	25.60	24.28	20.63	17.27	——/0.016s	——/62.5
ISTA-Net	<u>37.43</u>	<u>35.36</u>	<u>32.91</u>	<u>31.53</u>	<u>25.80</u>	<u>21.23</u>	<u>17.30</u>	0.923s/0.039s	1.08/25.6
ISTA-Net <sup>+</sup>	<b>38.07</b>	<b>36.06</b>	<b>33.82</b>	<b>32.57</b>	<b>26.64</b>	<b>21.31</b>	<b>17.34</b>	1.375s/0.047s	0.73/21.3

Table 2. Average PSNR (dB) performance comparison of various network-based algorithms on the BSD68 dataset.

Algorithm	CS Ratio				
	50%	40%	30%	10%	4%
SDA [30]	28.35	27.41	26.38	23.12	21.32
ReconNet [21]	29.86	29.08	27.53	24.15	21.66
ISTA-Net	<u>33.60</u>	<u>31.85</u>	<u>29.93</u>	<u>25.02</u>	<u>22.12</u>
ISTA-Net <sup>+</sup>	<b>34.01</b>	<b>32.21</b>	<b>30.34</b>	<b>25.33</b>	<b>22.17</b>

Table 3. Average PSNR (dB) comparison between ADMM-Net [44] and our proposed ISTA-Nets for CS-MRI.

Algorithm	CS Ratio				Time GPU
	20%	30%	40%	50%	
ADMM-Net	37.17	39.84	41.56	43.00	0.9535s
ISTA-Net	<u>38.30</u>	<u>40.52</u>	<u>42.12</u>	<u>43.60</u>	0.1246s
ISTA-Net <sup>+</sup>	<b>38.73</b>	<b>40.89</b>	<b>42.52</b>	<b>44.09</b>	0.1437s

# FISTA-Net: Learning a Fast Iterative Shrinkage Thresholding Network for Inverse Problems in Imaging

Jinxi Xiang<sup>ID</sup>, *Graduate Student Member, IEEE*, Yonggui Dong<sup>ID</sup>, *Member, IEEE*,  
and Yunjie Yang<sup>ID</sup>, *Member, IEEE*

**TMI 2021**

## Contributions:

- In FISTA-Net, the gradient matrix **is learned by substituting the fixed classic gradient** throughout iterations
- The core parameters of FISTA-Net, e.g., the step size, thresholding values, **are regularized to converge properly**
- A **momentum module** is added in FISTA-Net to accelerate convergence

Sparse representation:  $\hat{\mathbf{x}}_{\mathcal{F}_{l_1}} = \underset{\mathbf{x}}{\operatorname{argmin}} \{ \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathcal{F}(\mathbf{x})\|_1 \}$

FISTA-Net:

- Gradient descent module

$$\mathbf{r}^{(k)} = \mathbf{y}^{(k)} - (\mathbf{W}^{(k)})^T (\mathbf{Ay}^{(k)} - \mathbf{b})$$

- Proximal mapping module

$$\mathbf{x}^{(k)} = \mathcal{J}_{\theta^{(k)}}(\mathbf{r}^{(k)})$$

- Momentum module

$$\mathbf{y}^{(k+1)} = \mathbf{x}^{(k)} + \rho^{(k)}(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$

FISTA algorithm<sup>[1]</sup>:

$$\mathbf{x}^{(k)} = \mathcal{J}_{\alpha} \left( \mathbf{y}^{(k)} - \mu \mathbf{A}^T (\mathbf{Ay}^{(k)} - \mathbf{b}) \right)$$

$$t^{(k+1)} = \frac{1 + \sqrt{1 + 4(t^{(k)})^2}}{2}$$

$$\mathbf{y}^{(k+1)} = \mathbf{x}^{(k)} + \left( \frac{t^{(k)} - 1}{t^{(k+1)}} \right) (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$

[1] A fast iterative shrinkage-thresholding algorithm for linear inverse problems [J]. SIAM journal on imaging sciences, 2009, 2(1): 183-202.

FISTA-Net — Gradient descent module  $\mathbf{r}^{(k)}$  :

$$\mathbf{r}^{(k)} = \mathbf{y}^{(k)} - (\mathbf{W}^{(k)})^T (\mathbf{A}\mathbf{y}^{(k)} - \mathbf{b})$$

$\mathbf{W}^{(k)}$  can be decomposed as the product of a scalar  $\mu^{(k)}$  and a matrix  $\tilde{\mathbf{W}}$  independent of layer index  $k$ <sup>[1]</sup>:  $\mathbf{W}^{(k)} = \mu^{(k)}\tilde{\mathbf{W}}$

The matrix  $\tilde{\mathbf{W}}$  is pre-computed by solving

$$\begin{aligned} \tilde{\mathbf{W}} \in \arg \min_{\mathbf{W} \in \mathbb{R}^{N \times M}} \|\mathbf{W}^T \mathbf{A}\|_F^2 \\ \text{s.t. } (\mathbf{W}_{:,m})^T \mathbf{A}_{:,m} = 1, \forall m = 1, 2, \dots, M \end{aligned}$$

[1] Liu J, Chen X. ALISTA: Analytic weights are as good as learned weights in LISTA [C]//International Conference on Learning Representations (ICLR). 2019.

FISTA-Net — Proximal mapping module  $\mathbf{x}^{(k)}$  :

$$\mathbf{x}^{(k)} = \mathcal{J}_{\theta^{(k)}}(\mathbf{r}^{(k)})$$

- Similar to ISTA-Net
- Linear convolutional + ReLU
- Invertible operation

## FISTA-Net — Momentum module $\mathbf{y}^{(k+1)}$ :

- Accelerated by introducing a momentum term

$$\mathbf{y}^{(k+1)} = \mathbf{x}^{(k)} + \rho^{(k)}(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$



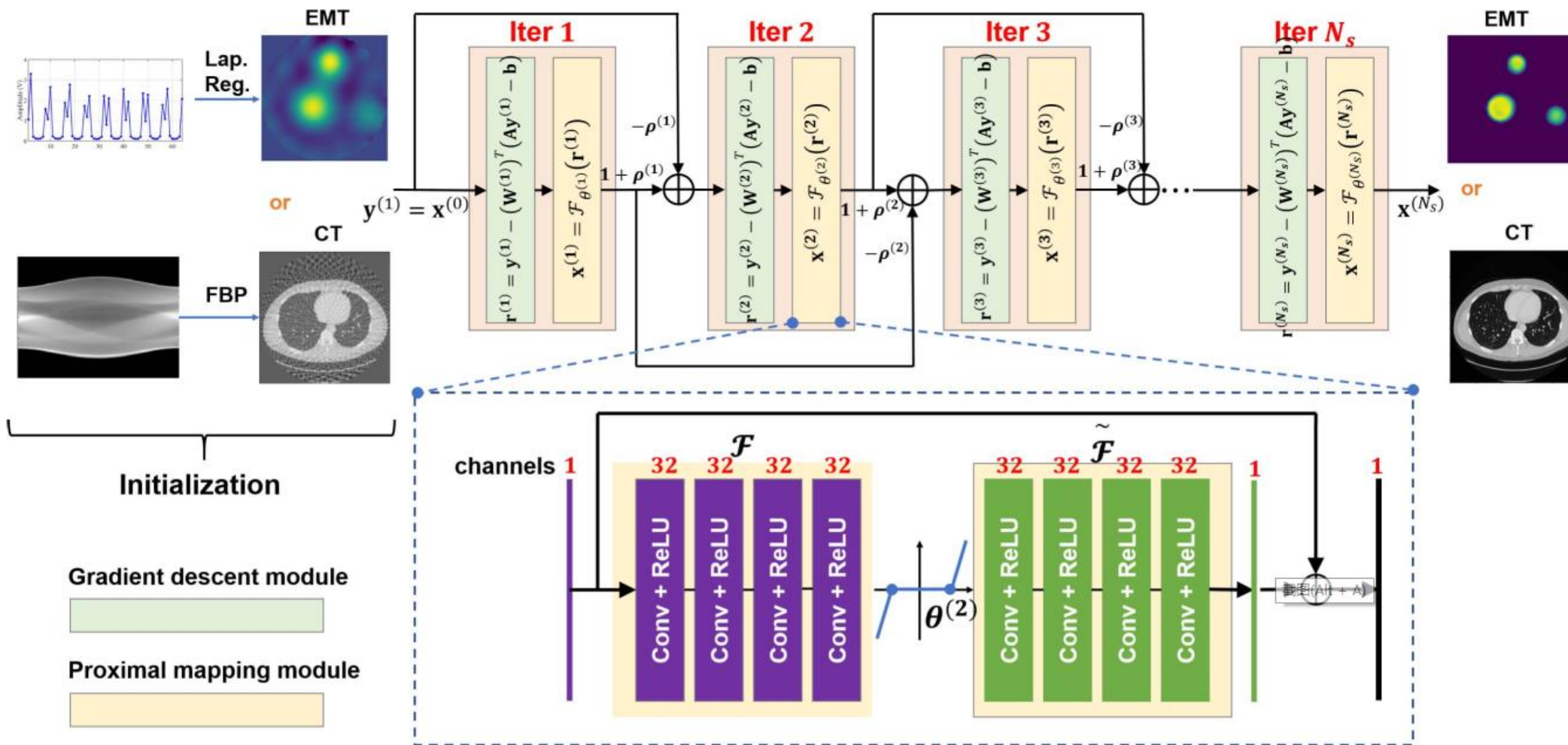
## FISTA-Net — Parameter constraints:

- ISTA-Net could possibly generate non-positive step size and thresholding values during iteration, which contradicts the definition of these variables

$$\begin{aligned}\mu^{(k)} &= sp(w_1 k + c_1), w_1 < 0 \\ \theta^{(k)} &= sp(w_2 k + c_2), w_2 < 0 \\ \rho^{(k)} &= \frac{sp(w_3 k + c_3) - sp(w_3 + c_3)}{sp(w_3 k + c_3)}, w_3 > 0\end{aligned}$$

where the softplus function  $sp(x) = \ln(1 + \exp(x))$  and  $\rho^{(k)} \in (0,1)$

# FISTA-Net framework:



## FISTA-Net Loss Function:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{mse}} + \lambda_1 \mathcal{L}_{\text{sym}} + \lambda_2 \mathcal{L}_{\text{spa}}$$

with

$$\begin{cases} \mathcal{L}_{\text{mse}} = \|\mathbf{x}_{N_s} - \mathbf{x}_{gt}\|_2^2 \\ \mathcal{L}_{\text{tsf}} = \lambda_1 \mathcal{L}_{\text{sym}} + \lambda_2 \mathcal{L}_{\text{spa}} \end{cases}$$

$$= \lambda_1 \sum_{k=1}^{N_s} \|\tilde{\mathcal{F}}(\mathcal{F}(\mathbf{r}^{(k)})) - \mathbf{r}^{(k)}\|_2^2 + \lambda_2 \sum_{k=1}^{N_s} \|\mathcal{F}(\mathbf{r}^{(k)})\|_1$$

**TABLE II**  
 QUANTITATIVE METRICS OF DIFFERENT METHODS ON EMT DATASET. THE BEST RESULTS OF  
 COMPARATIVE METHODS ARE HIGHLIGHTED IN BLUE

Testing set	Index	Lap. Reg. ( $x_0$ )	FISTA-TV [60]	FBPConvNet [21]	ISTA-Net [31]	<b>FISTA-Net</b>
	# of Pars.	1	1	482449	262094	<b>74599</b>
set 1	PSNR	16.140	17.996	<b>20.574</b>	20.074	<b>21.304</b>
	SSIM	0.556	0.661	<b>0.873</b>	0.772	<b>0.912</b>
	RMSE	0.155	0.1259	<b>0.094</b>	0.099	<b>0.086</b>
set 2	PSNR	14.975	16.508	<b>19.432</b>	19.215	<b>20.013</b>
	SSIM	0.472	0.630	<b>0.841</b>	0.725	<b>0.882</b>
	RMSE	0.178	0.149	<b>0.106</b>	0.109	<b>0.099</b>

**TABLE V**  
 QUANTITATIVE METRICS OF DIFFERENT METHODS ON MAYO CLINIC CT DATASET.  
 THE BEST RESULTS OF COMPARATIVE METHODS ARE HIGHLIGHTED IN BLUE

Testing set	Index	FBP ( $x_0$ )	FISTA-TV [60]	FBPConvNet [21]	ISTA-Net [31]	<b>FISTA-Net</b>
	# of pars	0	1	482449	262094	<b>74599</b>
60 view	PSNR	27.533	33.251	<b>33.976</b>	33.734	<b>37.319</b>
	SSIM	0.704	0.866	<b>0.942</b>	0.920	<b>0.951</b>
	RMSE	0.042	0.021	<b>0.020</b>	0.021	<b>0.013</b>
120 view	PSNR	29.152	35.217	<b>36.091</b>	35.452	<b>40.189</b>
	SSIM	0.815	0.949	<b>0.950</b>	0.948	<b>0.968</b>
	RMSE	0.035	0.017	<b>0.157</b>	0.016	<b>0.009</b>

# References

- [1] (LISTA) Gregor K, LeCun Y. Learning fast approximations of sparse coding [C]//Proceedings of the 27th international conference on international conference on machine learning. 2010: 399-406.
- [2] (ADMM-Net) Yang Y, Sun J, Li H, et al. Deep ADMM-Net for compressive sensing MRI [C]//Proceedings of the 30th International Conference on Neural Information Processing Systems. 2016: 10-18.
- [3] (ISTA-Net) Zhang J, Ghanem B. ISTA-Net: Interpretable optimization-inspired deep network for image compressive sensing [C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2018: 1828-1837.
- [4] (FISTA-Net) Xiang J, Dong Y, Yang Y. FISTA-Net: Learning A fast iterative shrinkage thresholding network for inverse problems in imaging [J]. IEEE Transactions on Medical Imaging, 2021, 40(5): 1329-1339.