

# Sparsity and Deep Neural Networks: a Match Made in Heaven

Pier Luigi Dragotti, Imperial College London

7 February 2022

---

## Motivation: A Theory for DL

- Deep Neural Networks achieves state-of-the-art performance in many imaging tasks
  - Fundamental questions:
    - is there a systematic way to interpret Deep Neural Networks?
    - Is there a systematic way to design the architecture of a Deep Neural Networks?
  - **Personal view:** sparse signal representation theory is much more developed and can be used to help addressing both questions
-

- Invertible Neural Networks and Wavelets
  - What are invertible Neural Networks (INN)?
  - Lifting Scheme and INN
  - Wavelet-like INN for denoising and deblurring
- Multimodal Image Processing and Unfolding
  - Multimodal Image Super-Resolution
  - Unfolding strategy for image separation in Art Investigation
- Conclusions

Common Theme: Interplay  
between sparsity and learning.



Junjie Huang (ICL, now Associate Prof. at NUDT))



Xin Deng (ICL, now Associate  
Prof. at Beihang University)



Consortium involving: UCL, ICL, Duke and National  
Gallery led by M. Rodrigues



# What are Invertible Neural Networks?

- Bijective (invertible) function approximators that have a forward mapping

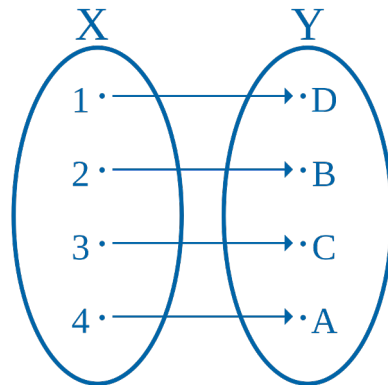
$$F_{\theta}: \mathbb{R}^d \rightarrow \mathbb{R}^l$$

$$x \mapsto z$$

- and inverse mapping

$$F_{\theta}^{-1}: \mathbb{R}^l \rightarrow \mathbb{R}^d$$

$$z \mapsto x$$



A bijective function (or  
invertible function)

# Imperial College London **What are Invertible Neural Networks?**

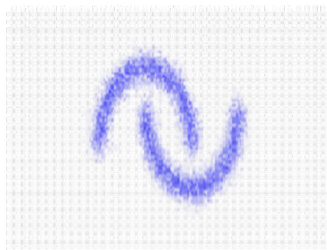
- INNs are bijective function approximators

**Inference**

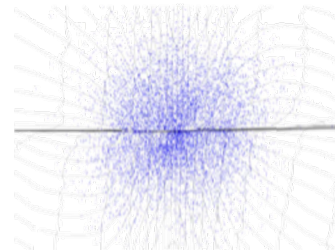
$$x \sim \hat{p}_X$$

$$z = f(x)$$

Data space  $\mathcal{X}$

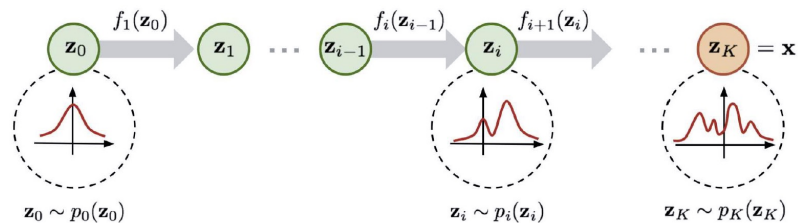


Latent space  $\mathcal{Z}$



# Imperial College London **Why Invertible Neural Networks?**

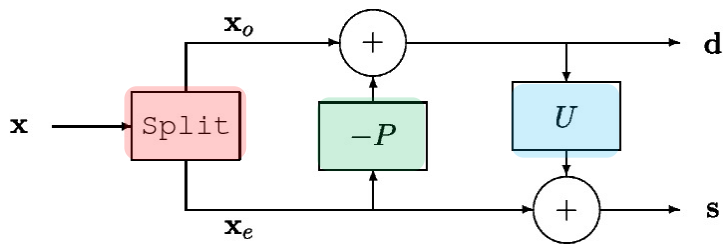
- Generative modeling
  - Tractable Jacobian, allows explicit computation of posterior probability
  - Normalizing flows



Kingma, Diederik P., and Prafulla Dhariwal. "Glow: Generative Flow with Invertible  $1 \times 1$  Convolutions." arXiv preprint arXiv:1807.03039 (2018).

# How to Achieve Invertibility?

- Invertible via lifting scheme like architectures
  - Signal splitting
  - Alternate prediction and update



$$\text{Split} \rightarrow \begin{cases} d = x_o - P(x_e) \\ s = x_e + U(d) \end{cases}$$

Forward pass

$$\begin{cases} x_o = d + P(x_e) \\ x_e = s - U(d) \end{cases} \rightarrow \text{Merge}$$

Backward pass

# Imperial College London Wavelets and Invertible Neural Networks

- In the beginning there were Wavelets 😊
- Wavelets provide sparse representations of piecewise smooth images.
- This is why they have been successfully used in many applications including denoising

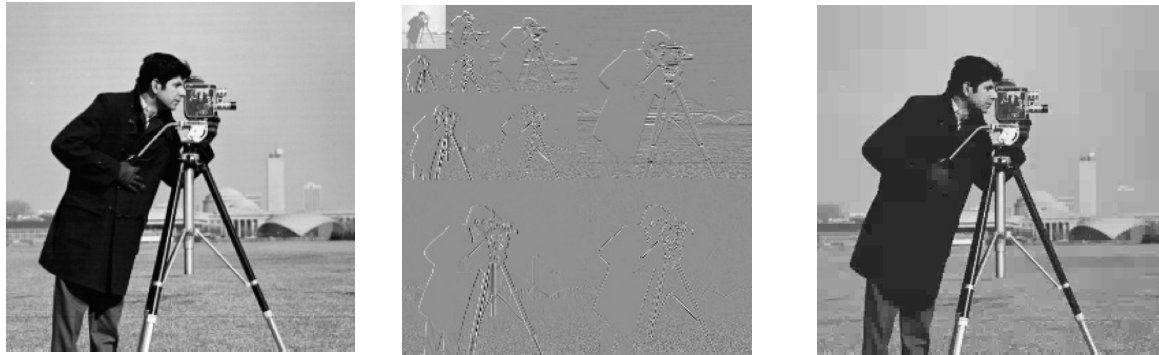
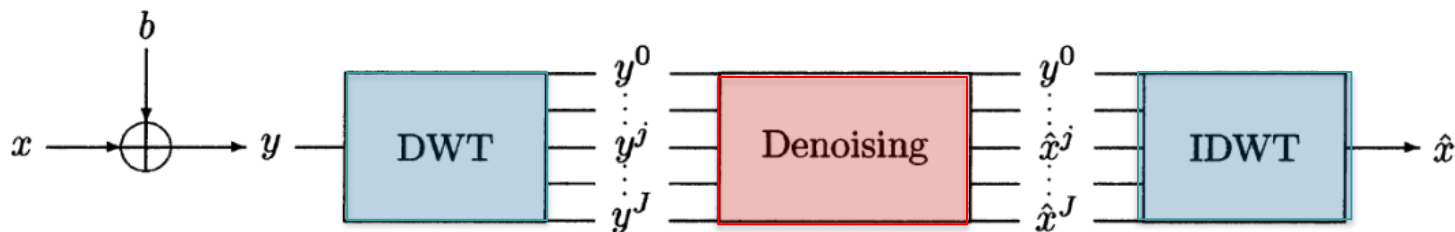


Figure: Cameraman is reconstructed using only 8% of the wavelet coefficients

- Principles of wavelet denoising:



## Wavelet transform

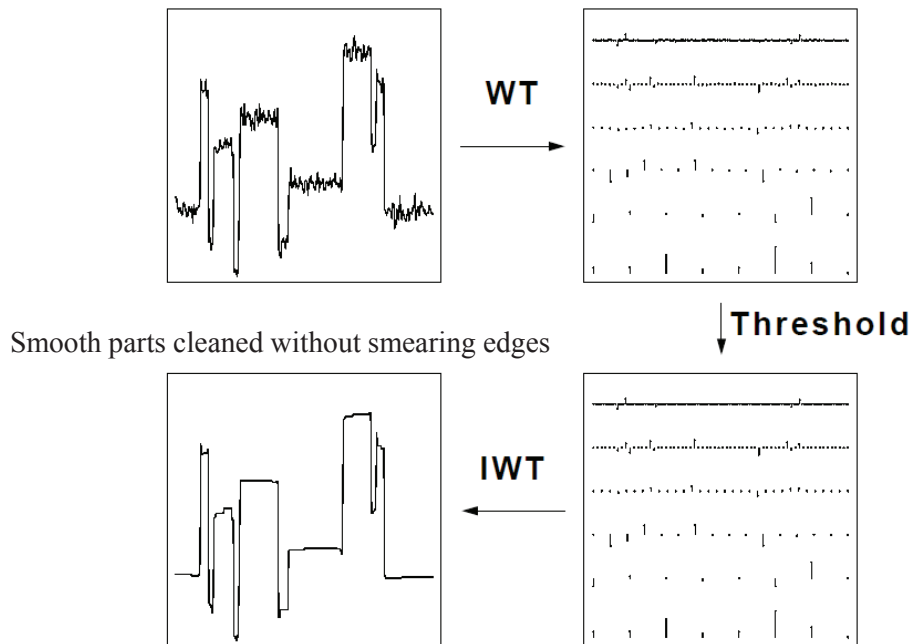
- Multi-resolution analysis
- Perfect reconstruction
- Noise is uniformly spread through the coefficients
- Image information is concentrated on small number of large coefficients

## Denoising

- Element-wise thresholding, e.g. soft-thresholding

# Wavelet-based Denoising

1-D Example

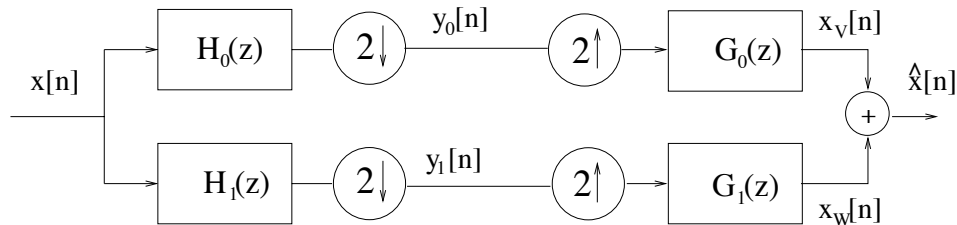


- Sparsity constraints in the wavelet domain (or in another domain) can also be used as a regularizers for different applications, e.g., deconvolution
- Iterative shrinkage:
  - $\min_{\alpha} (\|\mathbf{y} - \mathbf{H}\mathbf{W}^{-1}\alpha\|^2 + \lambda\|\alpha\|_1)$  where  $\mathbf{y}$  is the blurred image and  $\mathbf{x} = \mathbf{W}^{-1}\alpha$  is the target image
  - $\alpha_k = S_{\lambda}(\alpha_{k-1} + \mathbf{W}\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{W}^{-1}\alpha_{k-1}))$

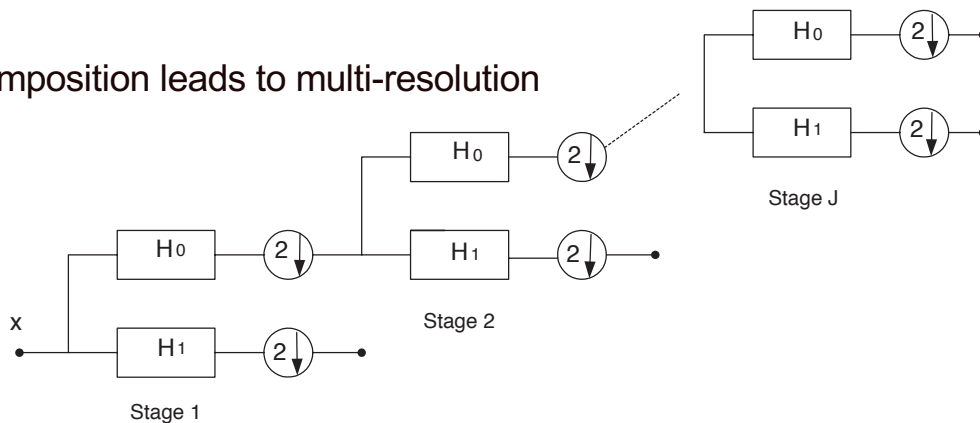


# Implementation of the Wavelet Transform

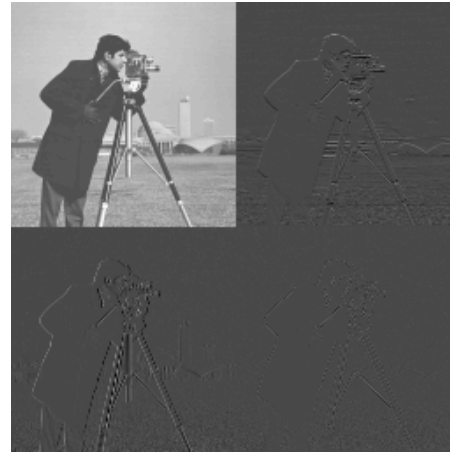
- Two-channel filter-bank



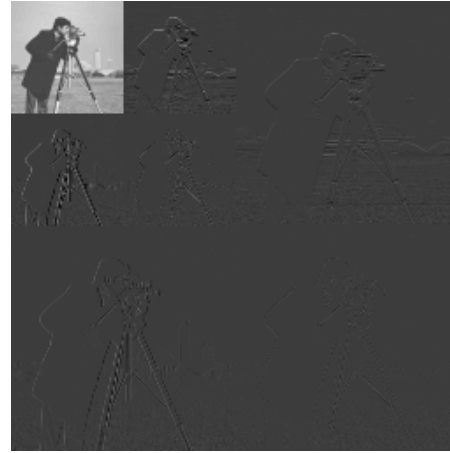
- Iterated decomposition leads to multi-resolution



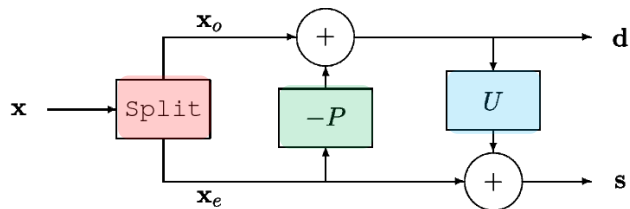
# Implementation of the 2-D Wavelet Transform



# Implementation of the 2-D Wavelet Transform



- The wavelet transform can be implemented using the lifting scheme



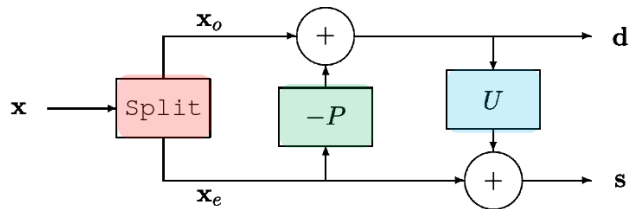
$$\text{Split} \rightarrow \begin{cases} d = x_o - P(x_e) \\ s = x_e + U(d) \end{cases} \quad \begin{cases} x_o = d + P(x_e) \\ x_e = s - U(d) \end{cases} \rightarrow \text{Merge}$$

Forward pass

Backward pass

- The **predictor** (P) predicts the odd samples using the even, the **update** (U) uses the prediction error to smooth the even samples
- Predictor/update are fixed
- The scheme is perfectly invertible

- Can we learn a wavelet-like non-linear sparsifying transform?



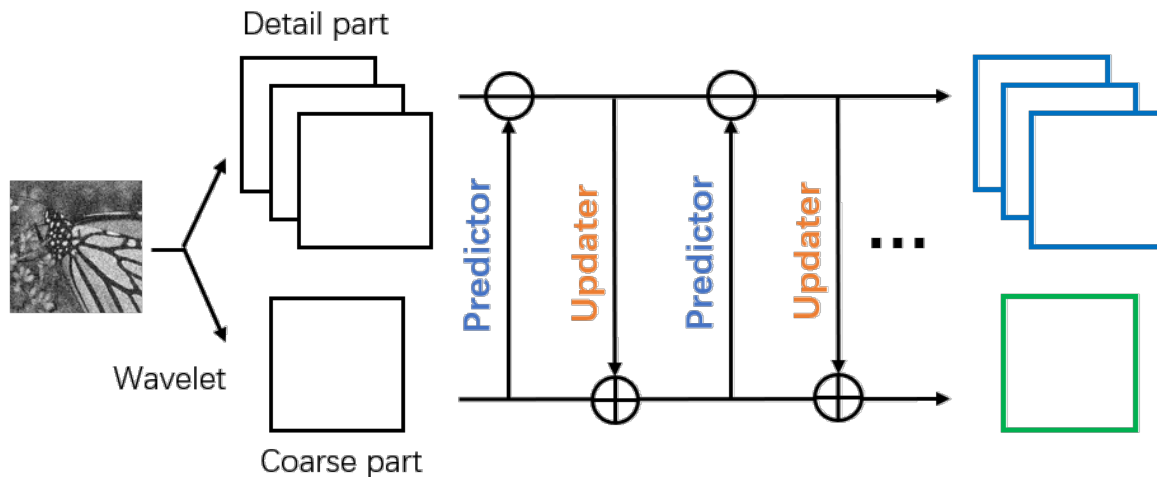
$$\text{Split} \rightarrow \begin{cases} d = x_o - P(x_e) \\ s = x_e + U(d) \end{cases} \quad \begin{cases} x_o = d + P(x_e) \\ x_e = s - U(d) \end{cases} \rightarrow \text{Merge}$$

Forward pass

Backward pass

- Approach:
  - convert the P/U operators into two deep networks and learn them
  - Use denoising as the bottleneck to impose sparsity

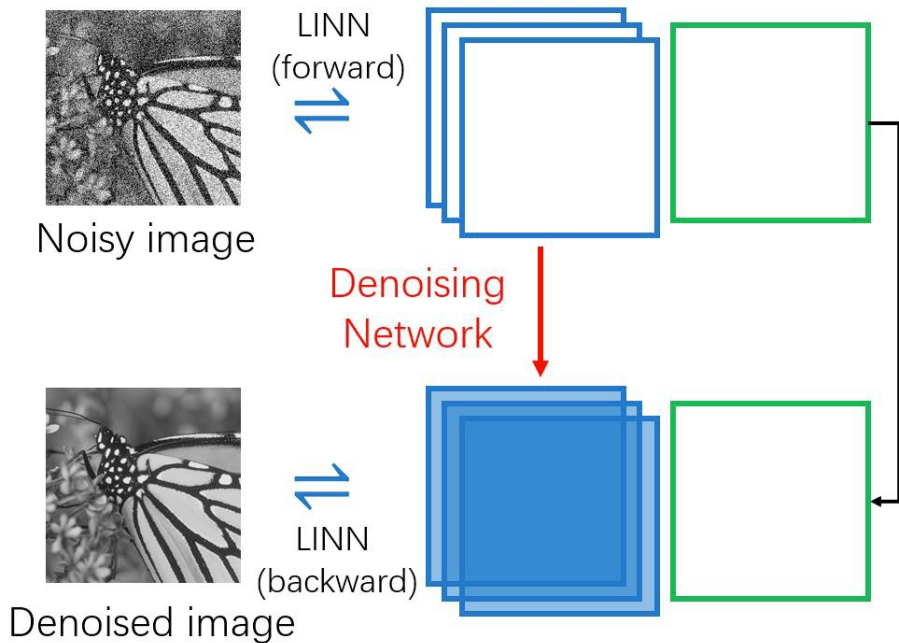
- Can we learn a wavelet-like non-linear sparsifying transform?



- Approach:
  - convert the P/U operators into two deep networks and learn them
  - Use denoising as the bottleneck to impose sparsity

## Wavelets and INN

- To make sure  $P$  acts as a sparsifying predictor:
  - Train the network with noisy/noiseless image pairs
  - Add a denoising network on the details

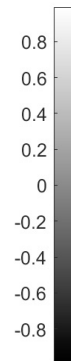
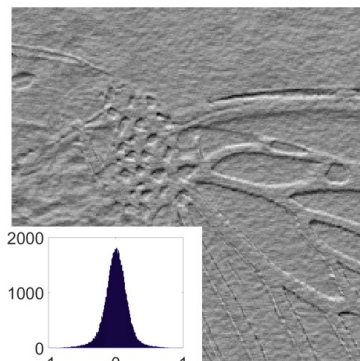


# Signal Decomposition

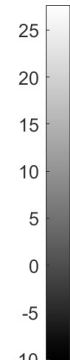
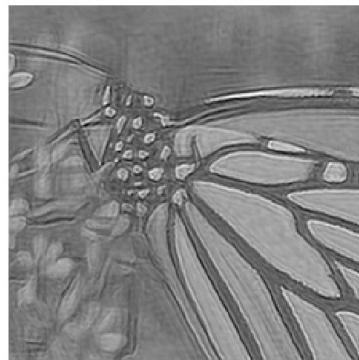
- Training with noiseless/noisy pairs leads to a sparsifying transform
- Each piece of the network is interpretable
- As for wavelets, we can now use the INN for e.g., denoising or deconvolution



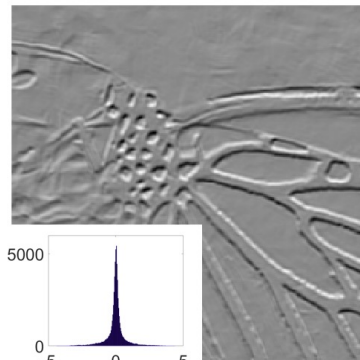
$C_1$



$d_1$



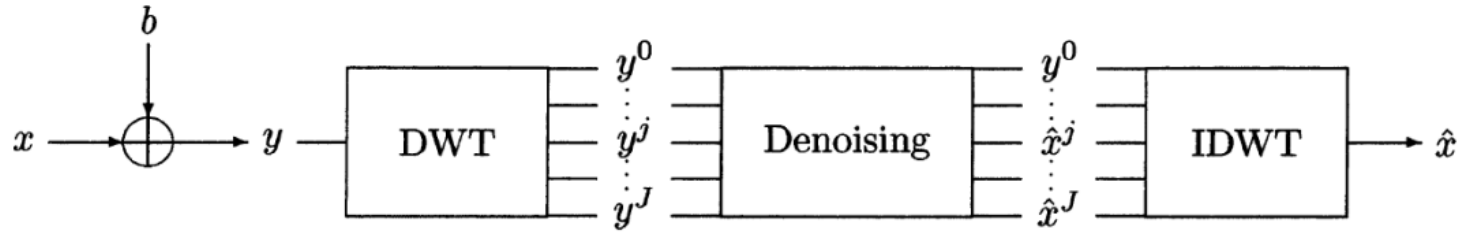
$C_2$



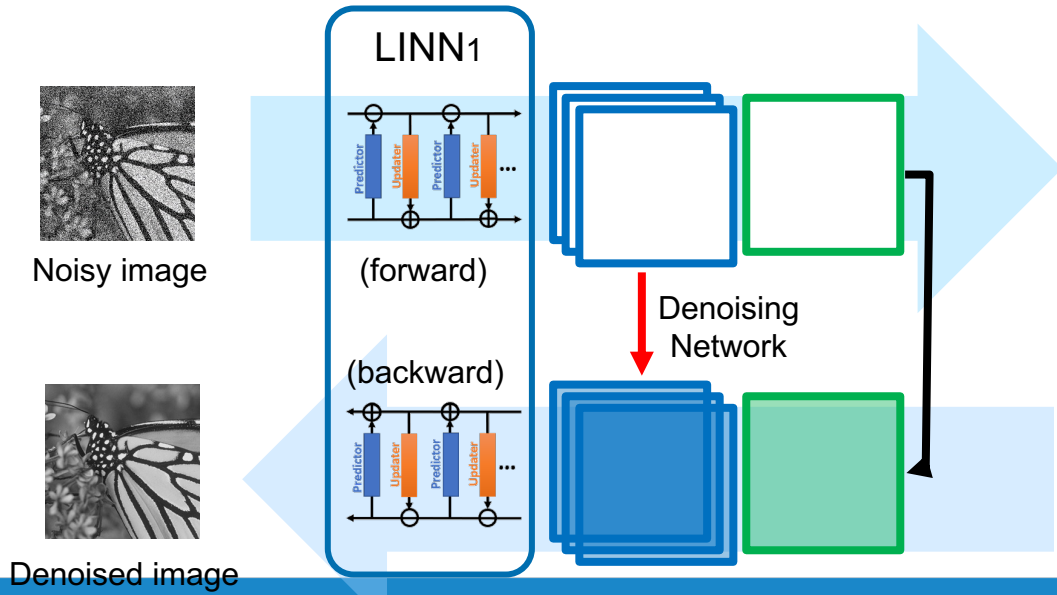
$d_2$



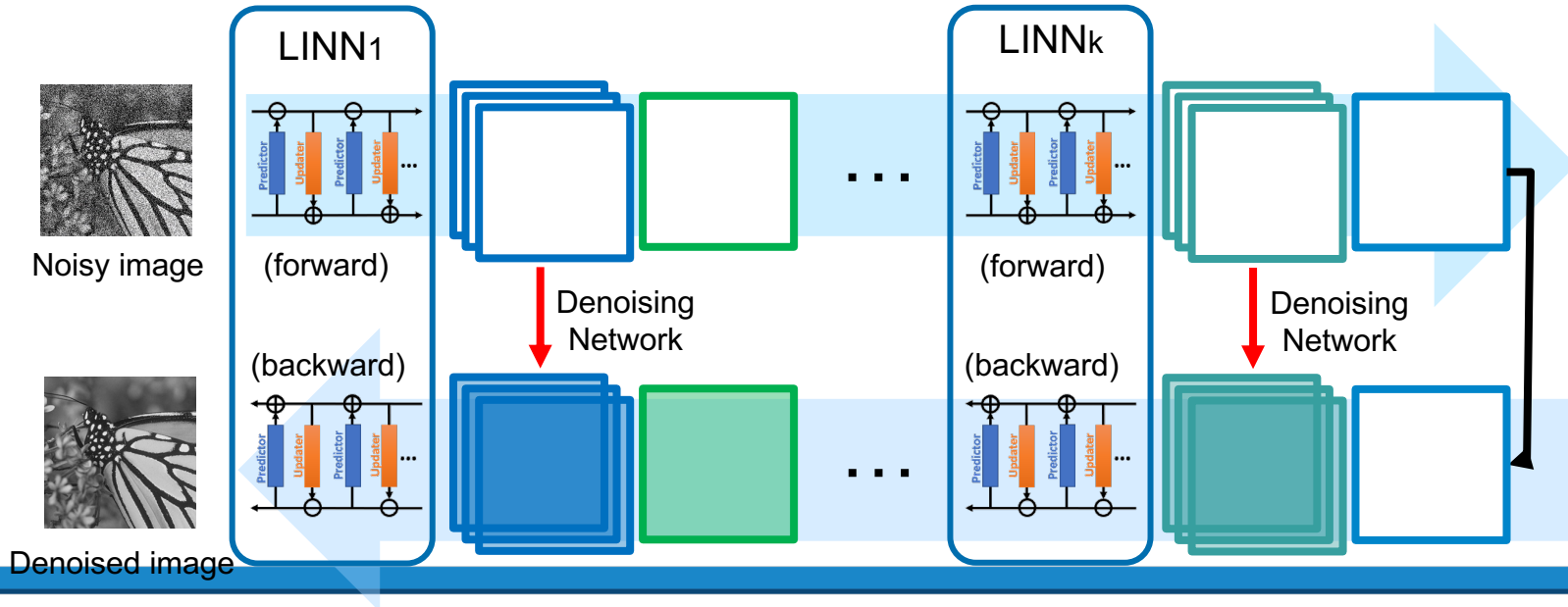
## Denoising - Overall Method



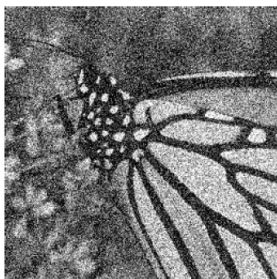
# Denoising - Overall Method



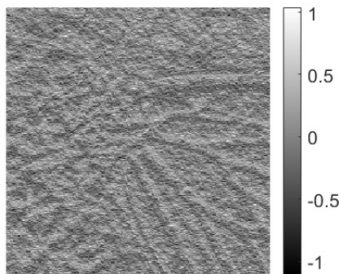
# Denoising - Overall Method



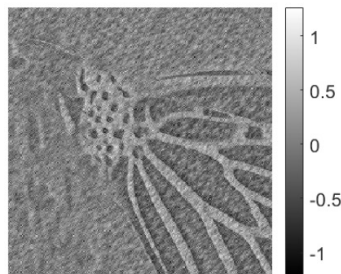
# Denoising



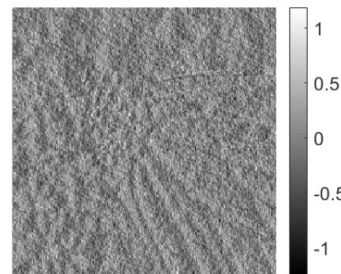
(a) Input noisy image ( $\sigma = 50$ ).



(b)  $z_d^I(1)$  before denoise.



(c)  $z_d^I(2)$  before denoise.

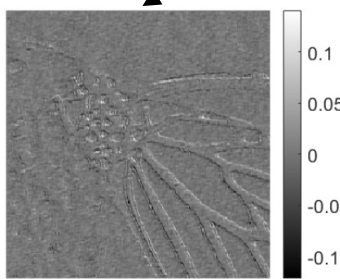


(d)  $z_d^I(3)$  before denoise.

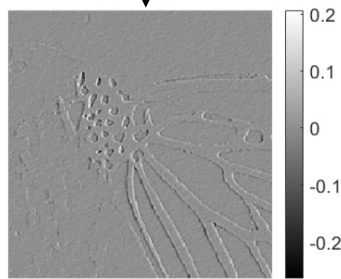
**Denoising**



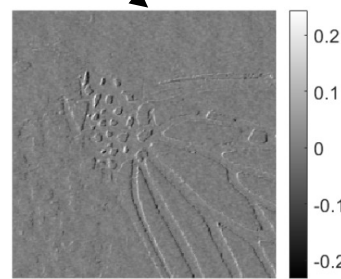
(e) The denoised image (PSNR=26.75 dB).



(f)  $z_d^I(1)$  after denoise.

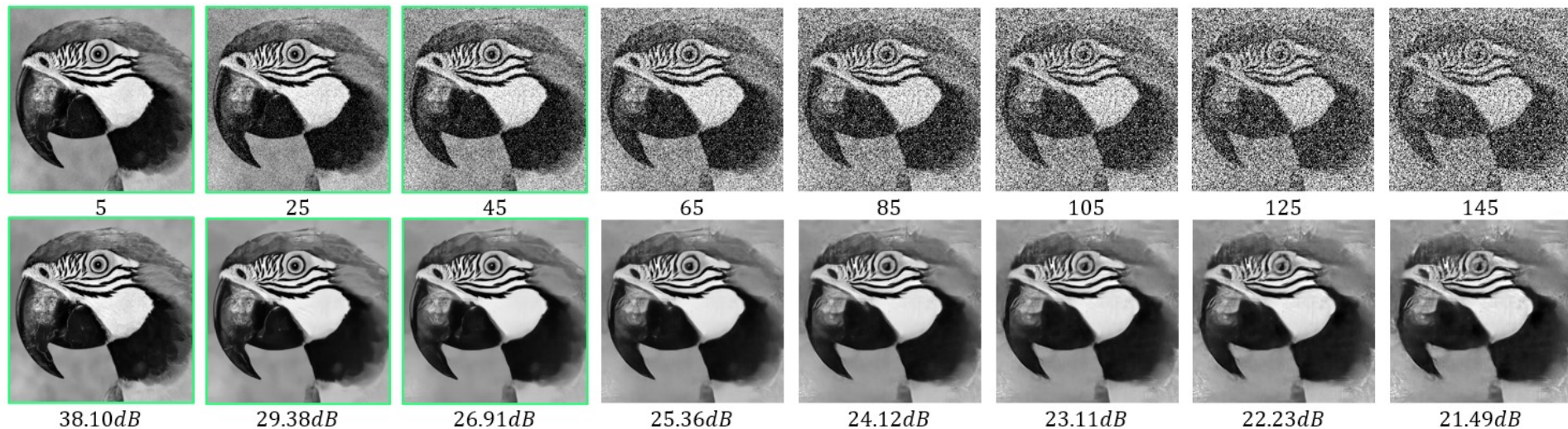


(g)  $z_d^I(2)$  after denoise.



(h)  $z_d^I(3)$  after denoise.

Denoising:



# Image Deblurring



=



+



---

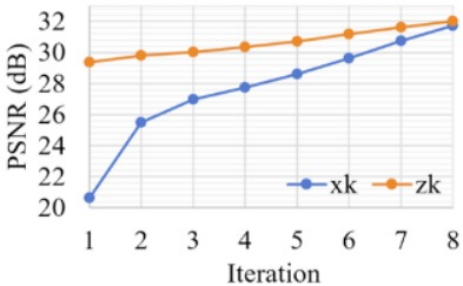
**Algorithm 1:** Plug-and-Play image deblurring with blind WINNet.

---

- 1 **Input:** Input image  $y$ , kernel  $k$ , parameter  $\lambda$ ;
  - 2 **Initialize:**  $z_0 = y$ ,  $\beta_0 = \text{NENet}(z_0)$ ,  $\beta_1 = 10 \times \beta_0$ ,  
 $k = 1$ ;
  - 3 **while**  $\beta_k > \beta_0$  **do**
    - 4  $x_k = \arg \min_x \|y - k \otimes x\|_2^2 + \frac{\lambda \beta_0^2}{\beta_k^2} \|x - z_{k-1}\|_2^2$ ;
    - 5  $\beta_{k+1} = \text{NENet}(x_k)$ ;
    - 6  $z_k = \text{WINNet}(x_k, 2\beta_{k+1})$ ;
    - 7  $k = k + 1$ ;
  - 8 **end**
  - 9 **Output:** Deblurred image  $x = z_{k-1}$ .
-



Deconvolution:



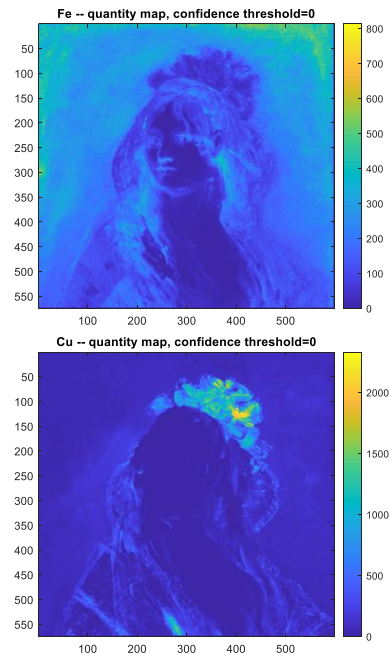


## First Set of Conclusions

- Invertible Neural Networks is an interesting new concept
  - Designing INN using wavelets/lifting leads to a more interpretable network
  - Good **generalization ability**
-

# Why Multi-modal Image Processing?

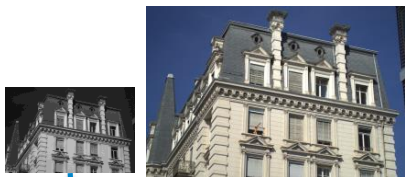
Technical Study of Old  
Master Paintings<sup>1</sup>:  
Data acquired using  
different imaging  
techniques



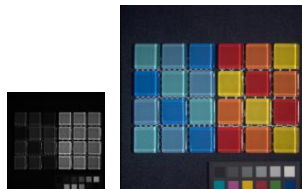
<sup>1</sup>joint project with UCL (M. Rodrigues), Duke and the National Gallery

# Why Multimodal Image Super-Resolution?

- Near-infrared (NIR)/RGB



- Multi-spectral/RGB

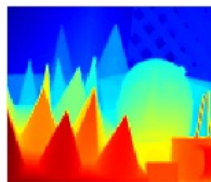


- Depth/RGB



Low resolution !

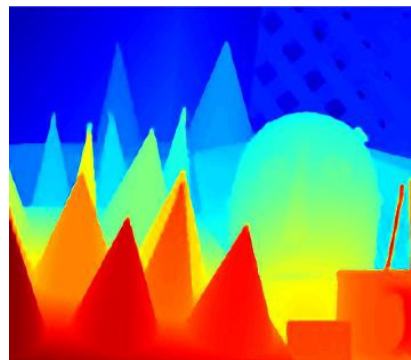
# Multimodal Image Super-Resolution (MISR)



LR Depth Image



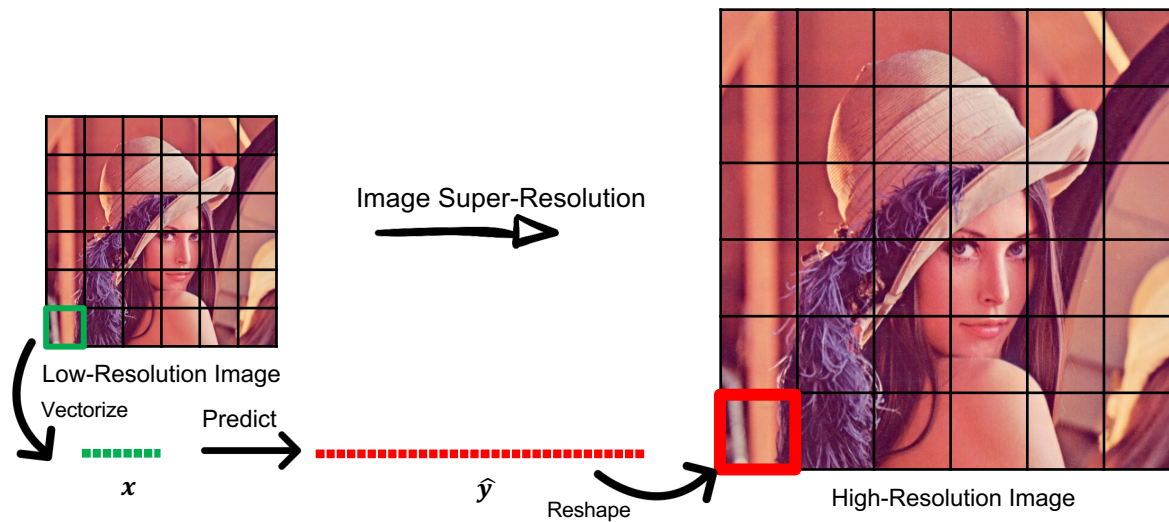
HR Color Image (guidance image)



Estimated HR Depth Image

# Single Image Super-Resolution

- Patch-based Prediction:

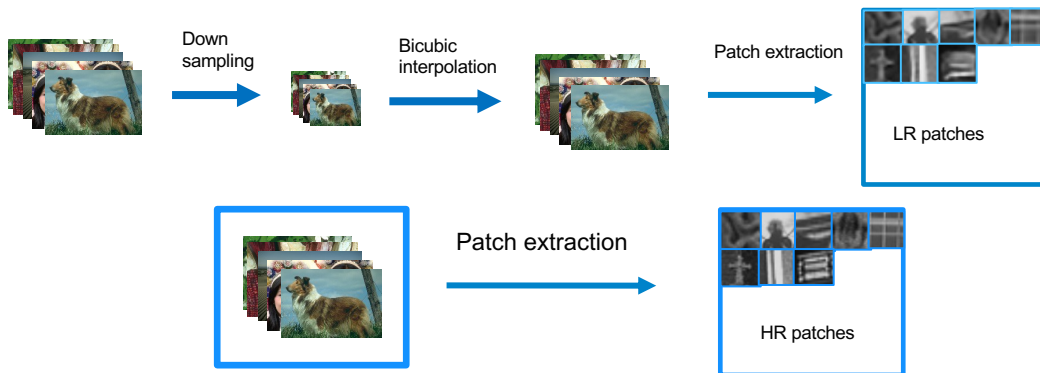


# Single Image Super-Resolution

Start with an external dataset of images (e.g., BSD 300 dataset)




Extract pairs of LR and HR patches





# Super-Resolution Model

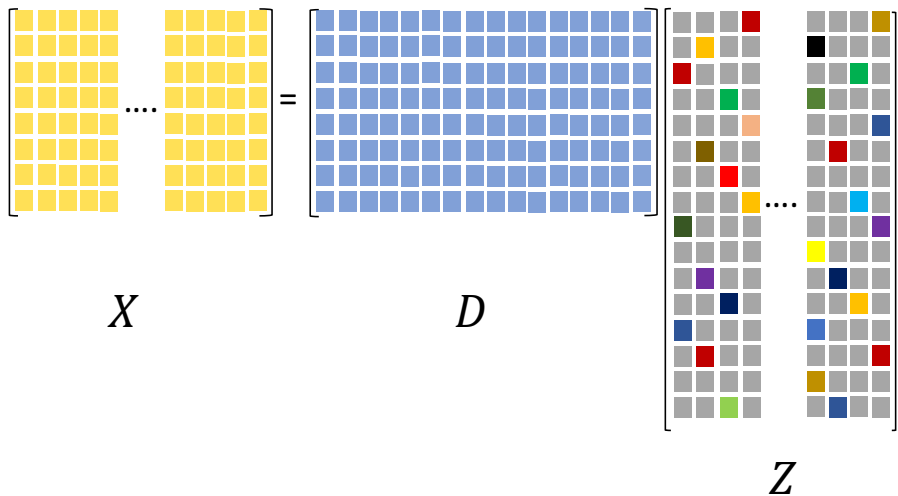
- One assumes that HR patches and LR patches admit a common sparse representation  $z_i$  :

$$\begin{aligned} \mathbf{x}_i^{LR} &= \mathbf{D}^{LR} \mathbf{z}_i \\ \mathbf{x}_i^{HR} &= \mathbf{D}^{HR} \mathbf{z}_i \end{aligned}$$




# Super-Resolution: Training

- Training:
  1. Given  $x_i^{LR}$ , learn  $D^{LR}$  and  $z_i$  using for example K-SVD or MOD



2. Given  $x_i^{HR}$  and  $z_i$  compute  $D^{HR}$  directly

# Multi-Modal Image Super-Resolution Model

- Approach:
  - Use sparse representation and dictionary learning to model dependency across modality and to drive the design of the neural network architecture through *unfolding*.

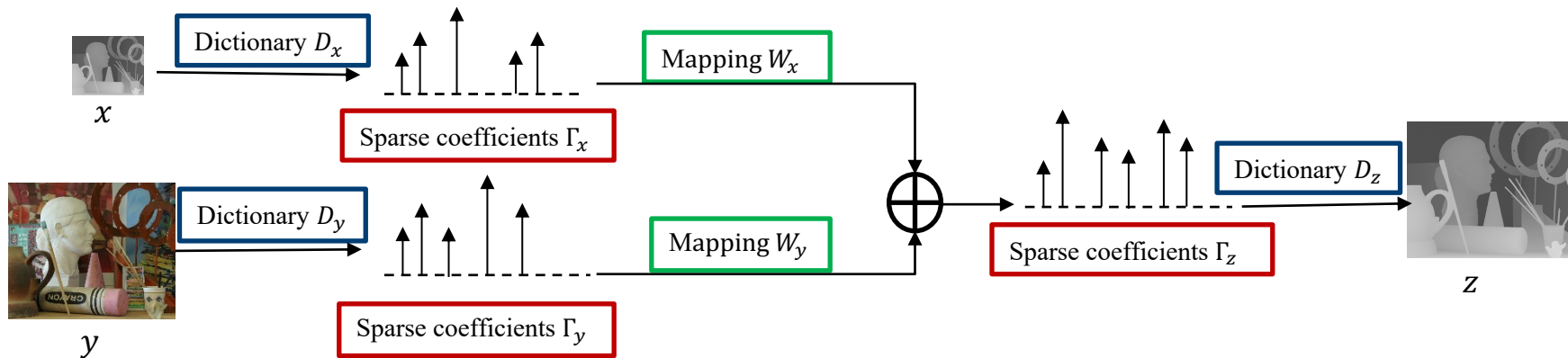
# Multi-Modal Image Super-Resolution Model

- Approach:
  - Use sparse representation and dictionary learning to model dependency across modality and to drive the design of the neural network architecture through *unfolding*.
  - **This is a trend now:** e.g Deep K-SVD [Elad et al.19], Neumann Networks [Willett-19], Deep Ultrasound [Eldar-19], Algorithm Unrolling [SP Mag, Eldar-21]

# Multi-Modal Image Super-Resolution Model

- Approach:
  - Use sparse representation and dictionary learning to model dependency across modality and to drive the design of the neural network architecture through *unfolding*.
- Model:
  - In the multimodal case the two modalities share some but not all latent features

# Joint multimodal dictionary learning (JMDDL)



- Assume patches  $x, y, z$  are sparse in learned dictionaries  $D_x, D_y$  and  $D_z$ , we can have the followings:

$$x \simeq D_x a,$$

$$y \simeq D_y b,$$

$$z \simeq D_z c,$$

where  $a, b, c$  are the sparse representations for  $x, y, z$ , respectively.

# Joint multimodal dictionary learning (JMDL)

- Since  $x, y, z$  are from the same image scene, we assume the sparse representation  $c$  of  $z$  can be inferred from the others:

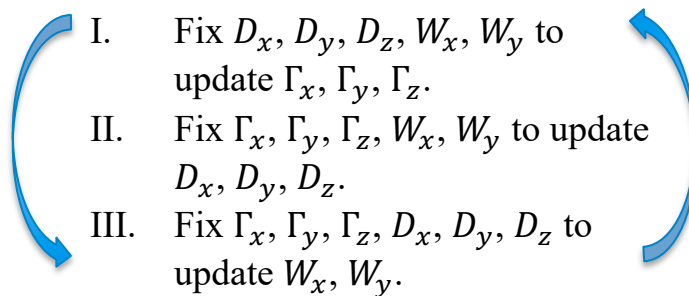
$$c = S_{\lambda_c}(\mathbf{W}_x \mathbf{a} + \mathbf{W}_y \mathbf{b})$$

where  $\mathbf{W}_x, \mathbf{W}_y$  are the transform matrices to be learned.

- JMDL optimization problem:

$$\begin{aligned} \min_{\substack{D_x, D_y, D_z, \\ \{\Gamma_x, \Gamma_y, \Gamma_z\}, \\ \mathbf{W}_x, \mathbf{W}_y}} & \frac{1}{2} \|\mathbf{X} - \mathbf{D}_x \Gamma_x\|_F^2 + \frac{1}{2} \|\mathbf{Y} - \mathbf{D}_y \Gamma_y\|_F^2 \\ & + \frac{1}{2} \|\mathbf{Z} - \mathbf{D}_z \Gamma_z\|_F^2 + \nu_x \|\Gamma_x\|_1 + \nu_y \|\Gamma_y\|_1 \\ & + \nu_z \|\Gamma_z\|_1 + \mu_x \|\mathbf{W}_x\|_F^2 + \mu_y \|\mathbf{W}_y\|_F^2 \\ & + \alpha \|\Gamma_z - \mathbf{W}_x \Gamma_x - \mathbf{W}_y \Gamma_y\|_F^2, \\ \text{s.t.}, & \|\mathbf{d}_{x,i}\|_2^2 \leq 1, \|\mathbf{d}_{y,i}\|_2^2 \leq 1, \|\mathbf{d}_{z,i}\|_2^2 \leq 1, \forall i_j \end{aligned}$$

- Solving strategy:



- In the reconstruction phase, given  $x$  and  $y$ , we first need to calculate their sparse representations based on the learned dictionaries  $D_x$  and  $D_y$ :

$$\min_{\{\mathbf{a}, \mathbf{b}\}} \frac{1}{2} \|\mathbf{x} - D_x \mathbf{a}\|_2^2 + \frac{1}{2} \|\mathbf{y} - D_y \mathbf{b}\|_2^2 + \lambda_x \|\mathbf{a}\|_1 + \lambda_y \|\mathbf{b}\|_1 .$$

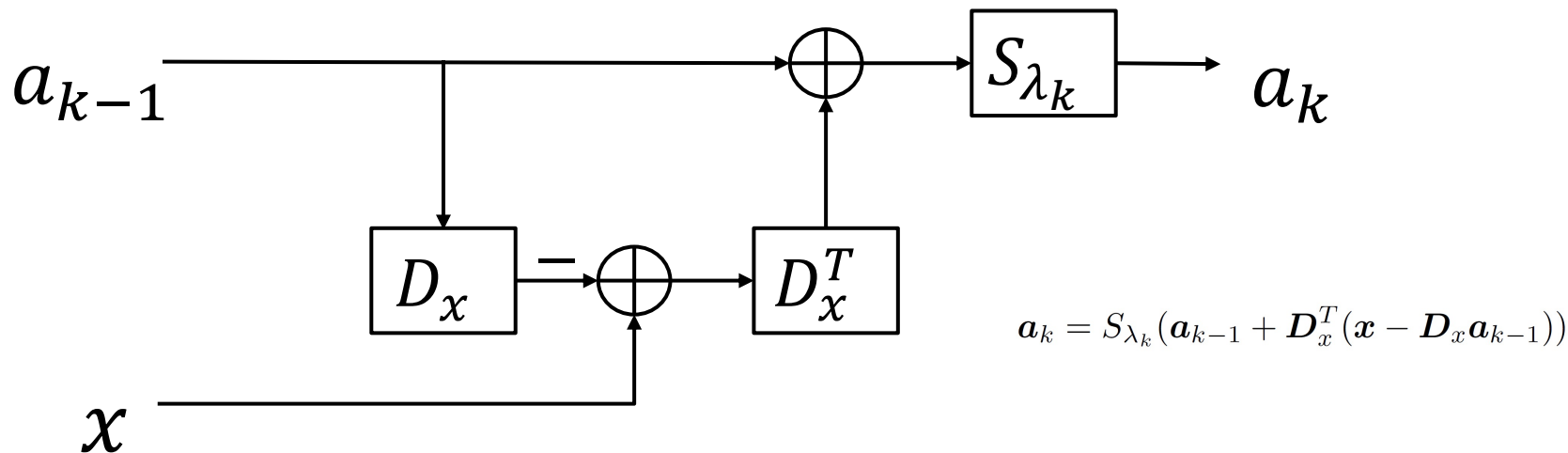
- Solved by ISTA algorithm:

$$\mathbf{a}_k = S_{\lambda_k}(\mathbf{a}_{k-1} + D_x^T(\mathbf{x} - D_x \mathbf{a}_{k-1}))$$

- Inspired by LISTA<sup>3</sup>, we “unfold” this iteration to obtain two deep networks (one per modality)

<sup>3</sup> Gregor Karol and LeCunYann, “Learning fast approximations of sparse coding”, Proceedings of the 27th International Conference on International Conference on Machine Learning, 2010

## Deep coupled ISTA network



Inspired by LISTA<sup>3</sup>, we “unfold” this iteration to obtain two deep networks (one per modality)

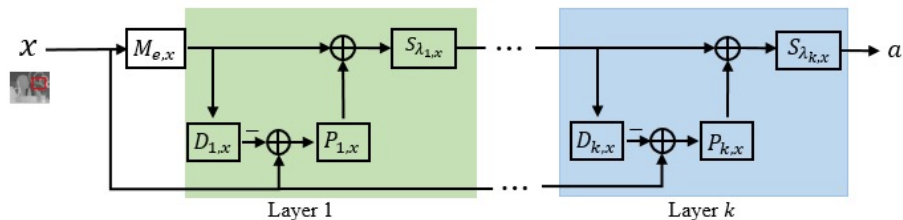
<sup>3</sup>Gregor Karol and LeCunYann, “Learning fast approximations of sparse coding”, Proceedings of the 27th International Conference on International Conference on Machine Learning, 2010



# Deep coupled ISTA network

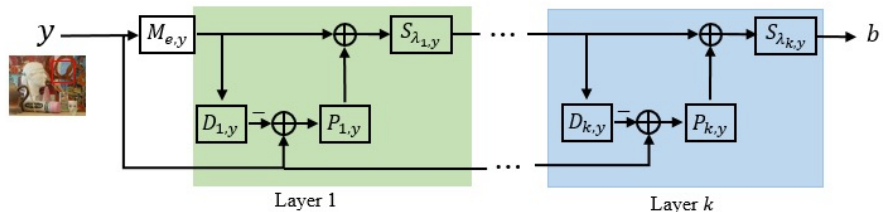
□ Solving by ISTA algorithm through unfolding:

$$\mathbf{a}_k = S_{\lambda_k}(\mathbf{a}_{k-1} + \mathbf{D}_x^T(\mathbf{x} - \mathbf{D}_x \mathbf{a}_{k-1}))$$

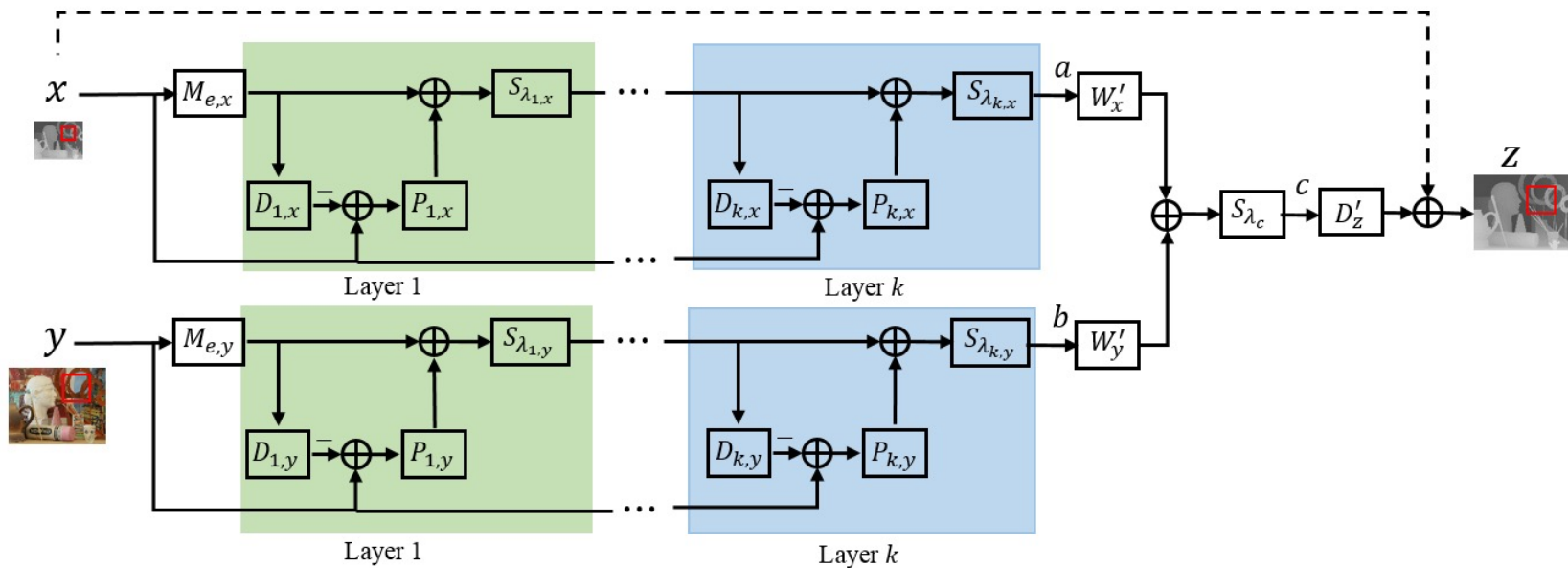


and

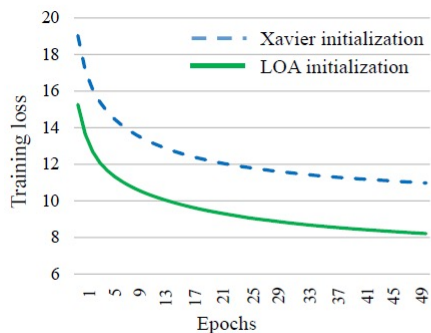
$$\mathbf{b}_k = S_{\lambda_k}(\mathbf{b}_{k-1} + \mathbf{D}_y^T(\mathbf{y} - \mathbf{D}_y \mathbf{b}_{k-1}))$$



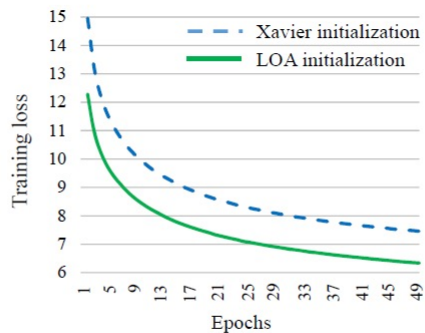
# Deep coupled ISTA network



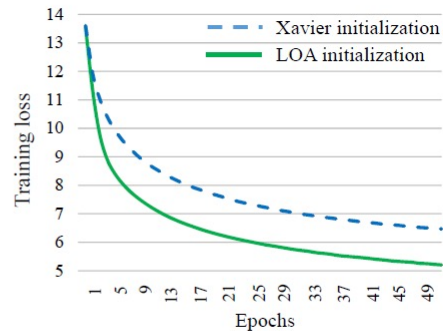
- Effectiveness of LOA



(a) Case I:  $m = 256, K = 2$



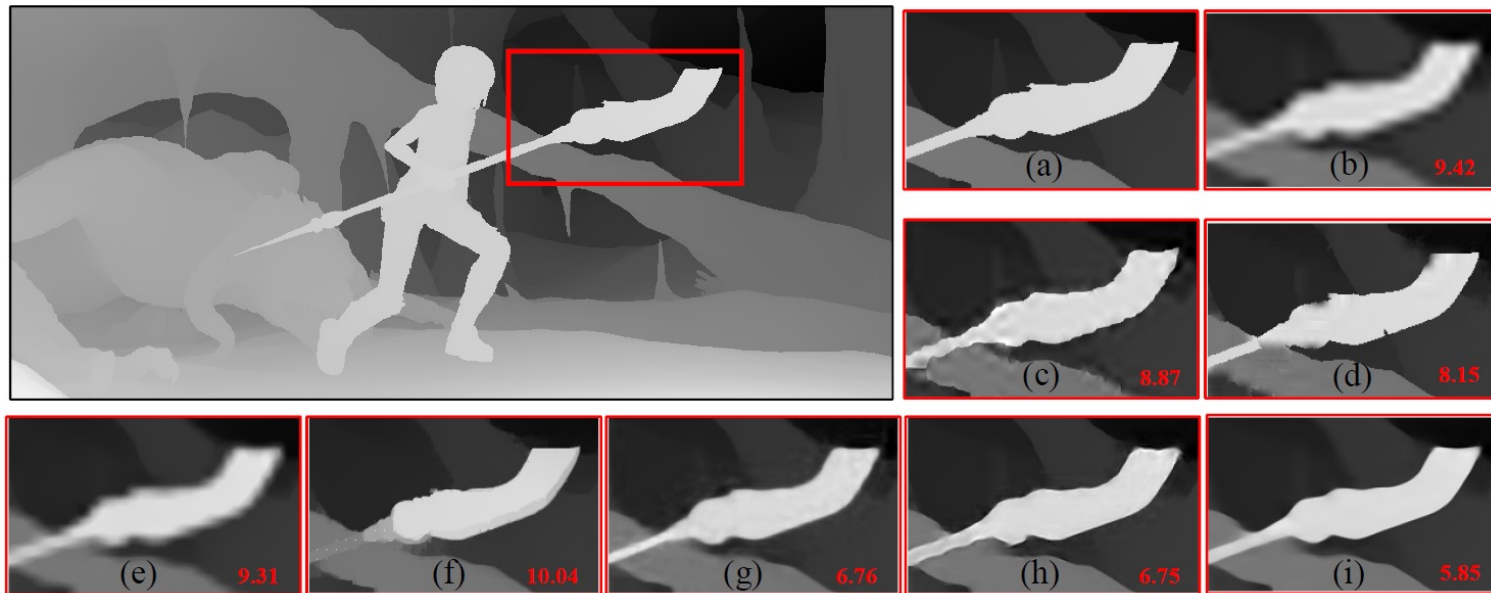
(b) Case II:  $m = 512, K = 2$



(c) Case III:  $m = 512, K = 4$

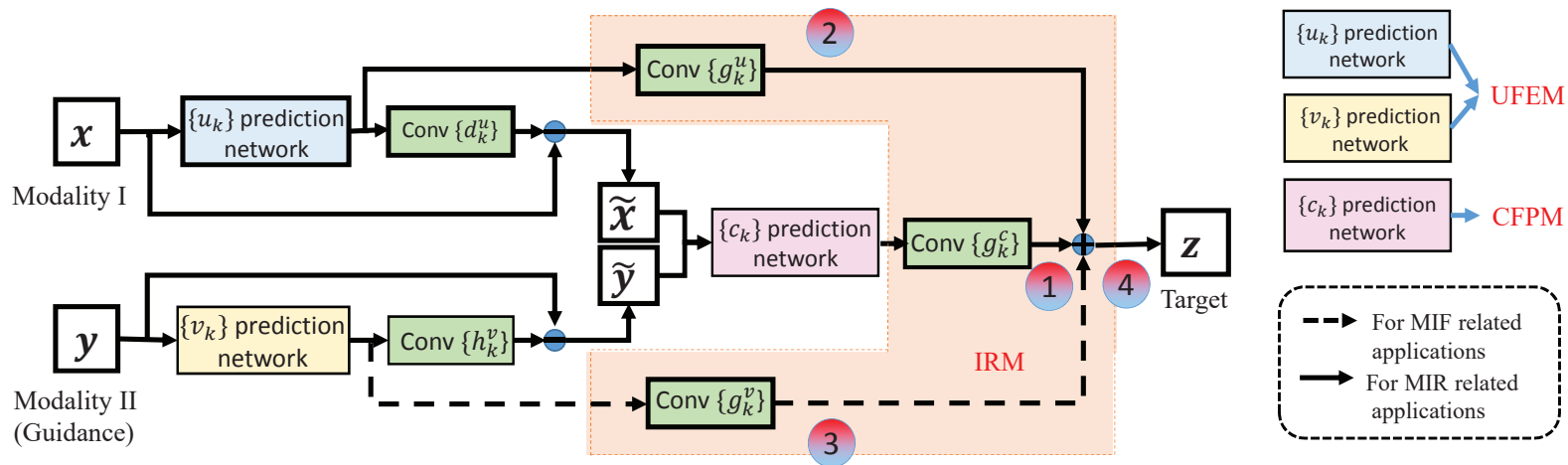
The training loss across 50 epochs with LOA and Xavier initialization methods for different settings of dictionary size  $m$  and number of layers  $K$ .

## Visual comparisons

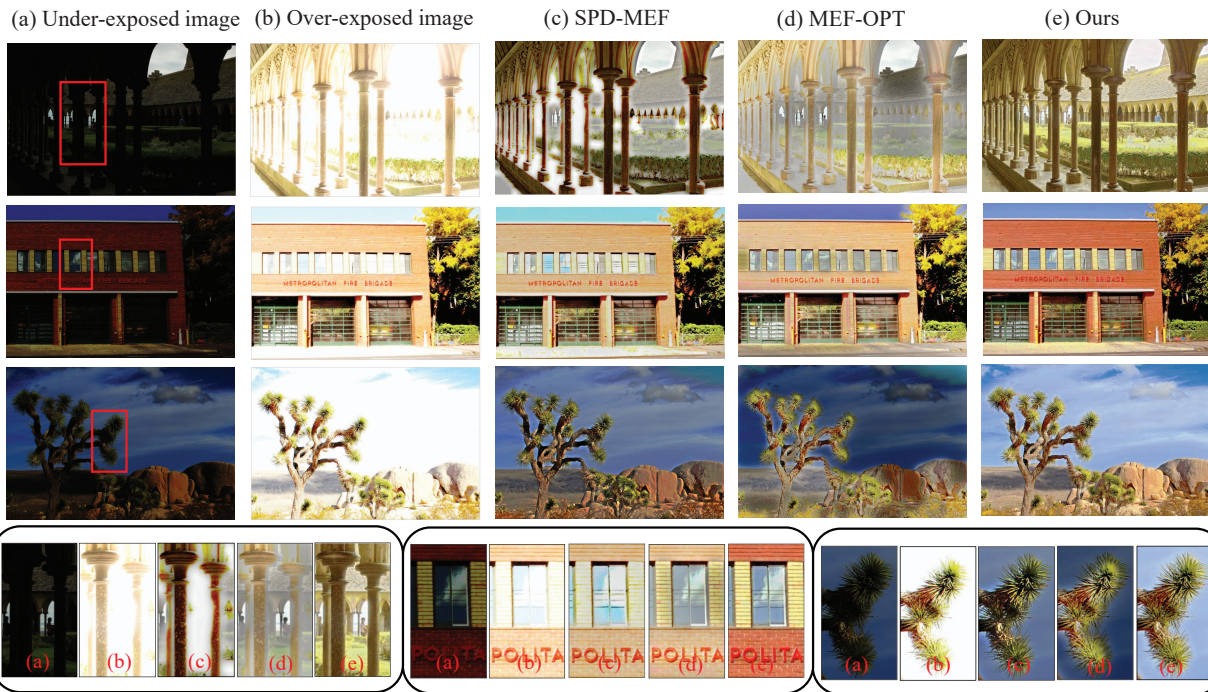


Visual comparisons of *Cave* in Sintel dataset with upscaling factor = 8 using different methods. (a) Ground truth, (b) Bicubic, (c) Park et al. [61], (d) Lu et al. [62], (e) Gu et al. [30], (f) Ferstl et al. [58], (g) SCN [12], (h) VDSR [13], (i) Ours. The numbers in red indicate the RMSE values.

# Unfolding Convolutional Dictionaries

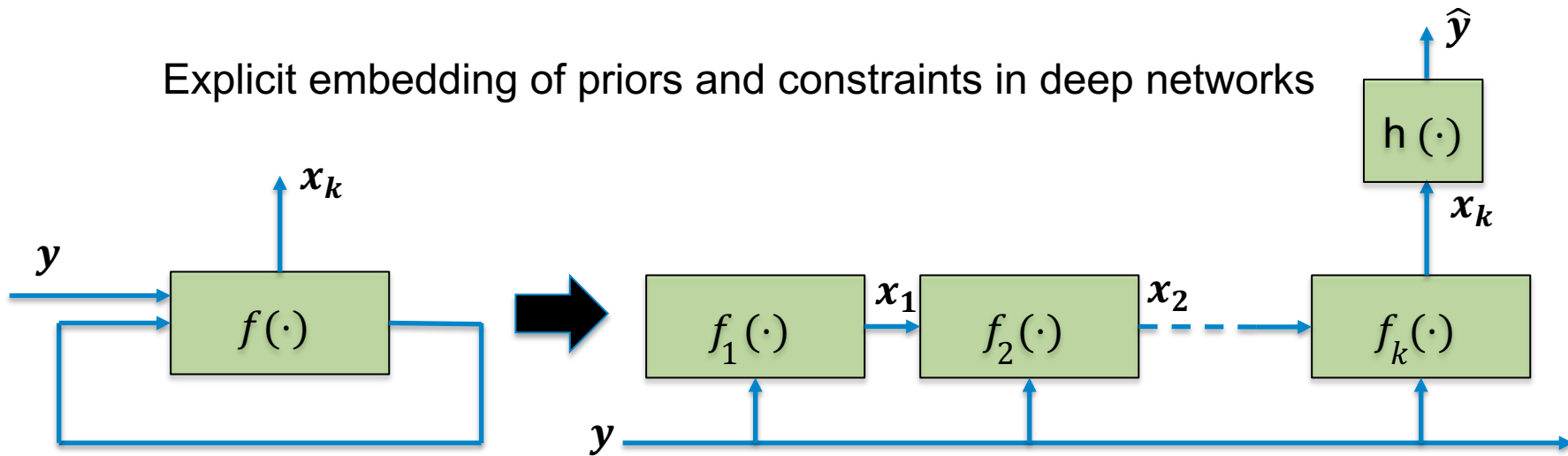


# Visual comparisons



# Unfolding Strategy

Explicit embedding of priors and constraints in deep networks



Iterative algorithm with  $y$   
as input and  $x$  as output

Unfolded version of the iterative algorithm with  
learnable parameters

Need to re-synthesize the input, if self-supervised

- Goal: Use multi-modal imaging techniques
  - for material characterization
  - to discover underdrawings and concealed design



Visible



X-ray



- Goal: we want to separate the two x-ray images
- Approach:
  - Use the visible RGB image as side information (x-ray visible similar to RGB image)
  - Exclusion loss: the “contours” of the two x-ray images should be as different as possible



Visible



X-ray

# Imperial College London X-ray Separation – Proposed Sparsity Model

$$\begin{aligned} \mathbf{x}_1 &= \sum_{k=1}^K \mathbf{\Xi}_k * \mathbf{z}_{1,k}, & \mathbf{x}_2 &= \sum_{k=1}^K \mathbf{\Xi}_k * \mathbf{z}_{2,k}, \\ \mathbf{r}_{1,s} &= \sum_{k=1}^K \Omega_{k,s} * \mathbf{z}_{1,k}, & \mathbf{x} &= \sum_{k=1}^K \mathbf{\Xi}_k * (\mathbf{z}_{1,k} + \mathbf{z}_{2,k}), \end{aligned}$$

- The visible image and the two separated X-ray images have a sparse representation in proper dictionaries
- RGB image and visible X-ray share the same sparse representation
- The two X-rays  $\mathbf{x}_1, \mathbf{x}_2$  share the same dictionary
- The measured X-ray is  $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$



Visible



X-ray

# X-ray Separation – Exclusion Loss

- Given the reconstructed X-ray images  $x_1, x_2$ , we expect that their edges are as different as possible we therefore add an “exclusion term” in the optimization

$$\begin{aligned}
 & \min_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{z}_{1,k}, \mathbf{z}_{2,k}} \|\mathbf{x} - \Psi * \mathbf{y}_1 - \Psi * \mathbf{y}_2\|_F^2 \\
 & + \tau_1 \|\mathbf{y}_1 - \sum_{k=1}^K \Theta_k * \mathbf{z}_{1,k}\|_F^2 \\
 & + \tau_2 \|\mathbf{y}_2 - \sum_{k=1}^K \Theta_k * \mathbf{z}_{2,k}\|_F^2 \\
 & + \gamma \sum_{s=1}^3 \|\mathbf{r}_{1,s} - \Phi_s * \mathbf{y}_1\|_F^2 \\
 & + \lambda_1 \sum_{k=1}^K \|\mathbf{z}_{1,k}\|_1 + \lambda_2 \sum_{k=1}^K \|\mathbf{z}_{2,k}\|_1 \\
 & + \sum_{i=1}^I \mu_i \|(\mathbf{W}_i * \mathbf{y}_1) \odot (\mathbf{W}_i * \mathbf{y}_2)\|_1,
 \end{aligned}$$



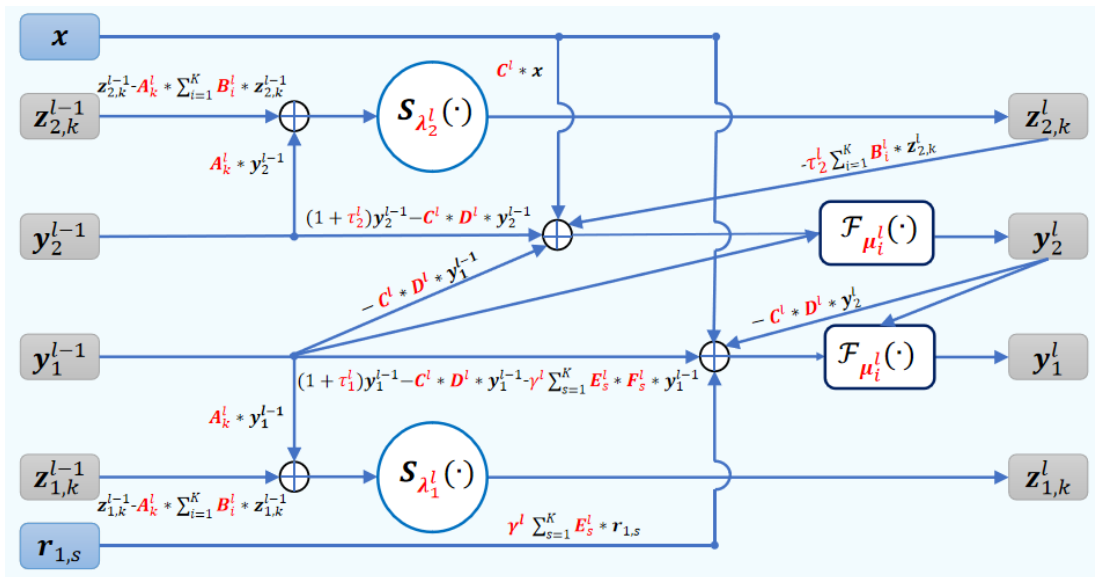
Visible



X-ray

# One Layer of the Network

- The sparsity model and the exclusion constraint leads to an iterative optimization method which leads to a network through unfolding



## Separation Results



## Conclusions

- Cross fertilization between dictionary learning/sparse representation and deep learning is fruitful
  - Dictionary Learning/sparsity useful:
    - to impose models and structure to the deep network (through sparse modelling and optimization)
    - To design wavelet-like INN
    - For better interpretability and generalization ability
-

- **J. Huang** and P.L. Dragotti, “LINN: Lifting Inspired Invertible Neural Network for Image Denoising”, in proc. of 29th European Signal Processing Conference, EUSIPCO 2021
- **J. Huang** and P.L. Dragotti, “WINNet: Wavelet-inspired Invertible Network for Image Denoising”, submitted to IEEE Transactions on Image Processing, September 2021, <https://arxiv.org/abs/2109.06381>
- **X Deng** and P. L. Dragotti, “Deep Coupled ISTA Network for Multi-modal Image Super-Resolution, IEEE Transactions on Image Processing, pp 1683-1698, vol.29, 2020
- **X Deng** and P. L. Dragotti, “Deep Convolutional Neural Network for Multi-modal Image Restoration and Fusion IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 43, October 2021
- **W. Pu, J. Huang** et al., “Mixed X-Ray Image Separation for Artworks with Concealed Designs”, submitted to IEEE Transactions on Image Processing, January 2022, <https://arxiv.org/abs/2201.09167>