

ADMM DIP-TV: combining Total Variation and Deep Image Prior for image restoration

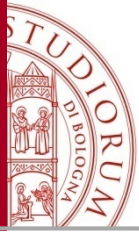
Pasquale Cascarano, *Andrea Sebastiani, Elena Loli Piccolomini,*
University of Bologna

Maria Colomba Comes, Eugenio Martinelli,
University of Rome Tor-Vergata

Warsaw, 15 December 2020

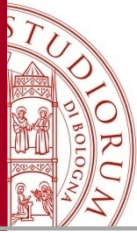
BOS/SOR2020 Conference

Session III: Advances in optimization techniques for machine learning



Outline

1. Imaging Inverse Problems
2. Variational approach
3. Tips on supervised learning
4. Deep Image Prior (DIP)
5. The DIPTV framework: ADMM-DIP TV
6. Results and final remarks



Imaging Inverse Problems

- Imaging inverse problems attempt to **retrieve an unknown data from its measurement**.
- The model we refer reads:

$$Ax + \eta = b \quad (1)$$

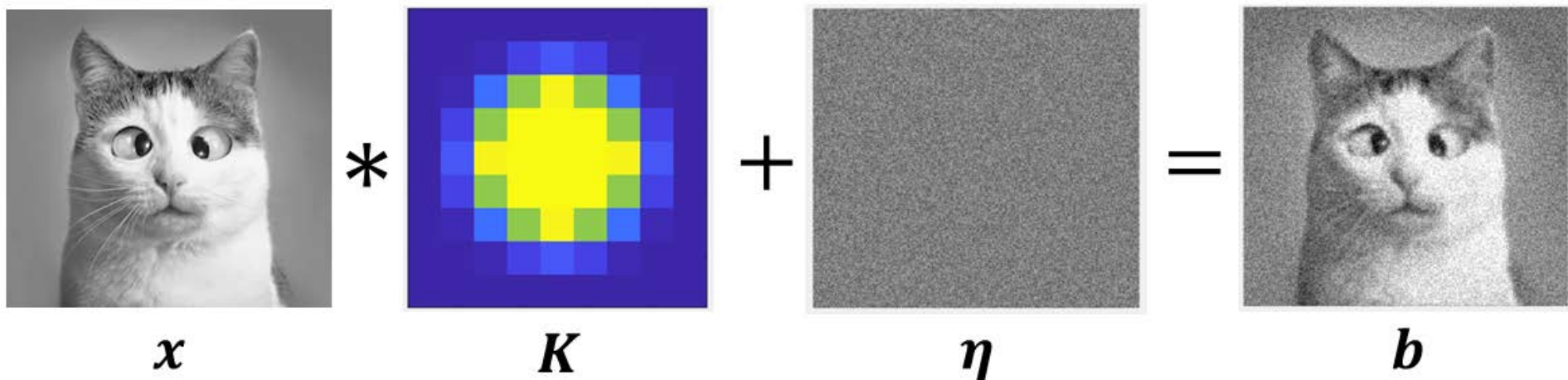
- The linear operator $A \in \mathbb{R}^{l \times n}$ is the forward operator, $x \in \mathbb{R}^n$ is the unknown data, $b \in \mathbb{R}^l$ is the measurement and η is the noisy component.
- Solving (1) means **inverting the process**:

$$b \longrightarrow x$$

- Different choices for $A \in \mathbb{R}^{l \times n}$ lead to different image restoration/reconstruction tasks.

Deblurring

- The matrix A is the discretization of a convolution.

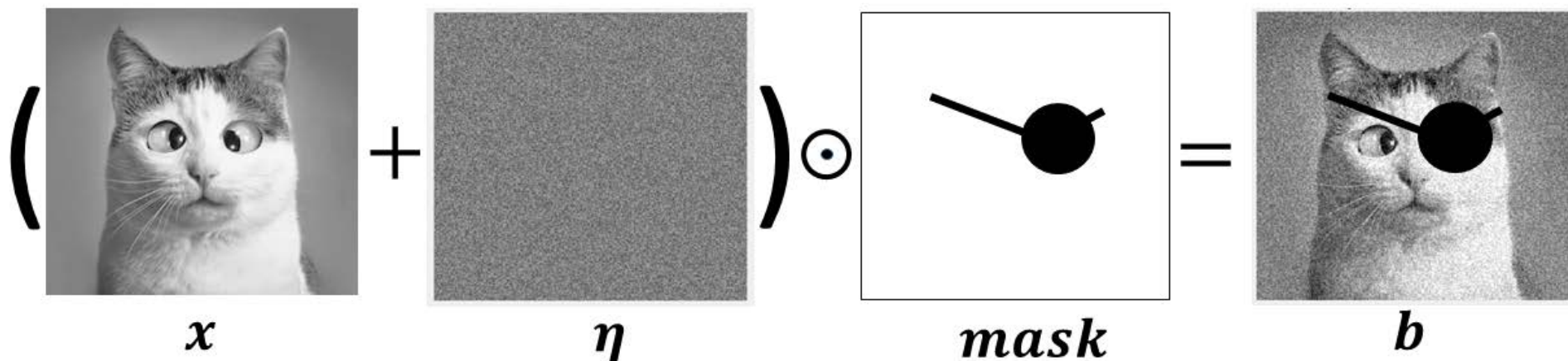


The diagram illustrates the process of deblurring. It shows a sequence of four images connected by mathematical operators:

- x : A sharp grayscale image of a cat's face.
- $*$: A convolution operator.
- K : A 2D kernel matrix represented as a color-coded grid with a central yellow and green cross on a blue background.
- $+$: An addition operator.
- η : A grayscale image of random noise.
- $=$: An equality operator.
- b : A noisy, blurred version of the original cat image.

Inpainting

- The matrix A is the discretization of the Hadamard product.



Super Resolution

- The matrix A is the discretization of the composition between a downsampling and a convolution operator.

$$S \left(\begin{array}{c} \text{cat image} \\ x \end{array} * \begin{array}{c} \text{kernel } K \\ K \end{array} \right) + \begin{array}{c} \text{noise } \eta \\ \eta \end{array} = \begin{array}{c} \text{super-resolution image } b \\ b \end{array}$$

The diagram illustrates the super-resolution process. It shows a grayscale image of a cat, labeled x , being convolved with a kernel K . The kernel K is a 7x7 grid with a central yellow pixel, surrounded by green and blue pixels. The result of the convolution is added to a noise matrix η , which is a 7x7 grid of gray noise. The final result is a super-resolution image b , which is a smaller, sharper version of the original cat image.

Variational approach

- Imaging inverse problems are usually **ill-posed problems** which means that the properties of existence, uniqueness and stability of the solution are not all verified
- They are **re-casted as an optimization problem** of the following form:

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \{f(x) + \lambda \phi(x)\} \quad (2)$$

where f encodes the information on the noise, ϕ reflects prior information on the desired solution and λ is a trade-off parameter.

Under **Gaussian Noise assumptions** and if no prior is defined, (2) reads:

$$\arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2$$



b

$$(A^T A)^{-1} A^T$$



\hat{x}

Total Variation

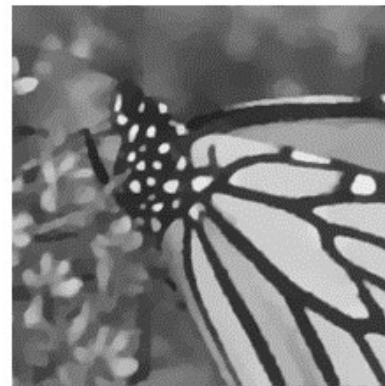
- The prior is usually defined by assuming some geometric and smoothness properties.
- One of the most famous prior is the Total Variation defined as follows:

$$TV(x) = \sum_{i=1}^n \sqrt{(D_h x)_i^2 + (D_v x)_i^2} \longrightarrow \text{images are piecewise constant}$$

$$\hat{x}_{TV} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 + \lambda TV(x)$$



b



\hat{x}_{TV}

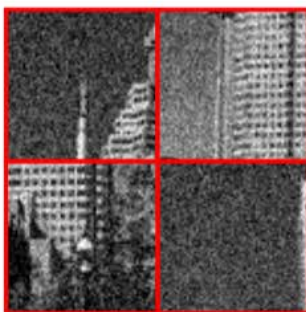
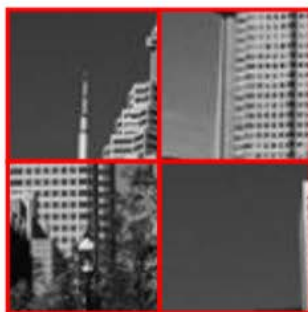
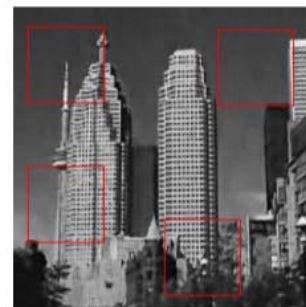
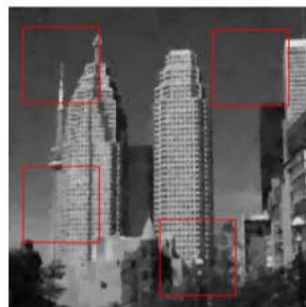
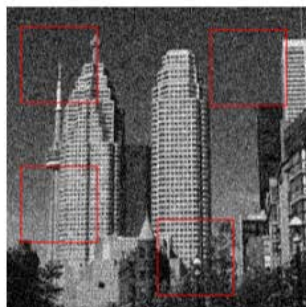
TV Prior

PROS

- Closed form.
- Strong mathematical foundations.
- Convergence guaranteed.

CONS

- Lack of flexibility.
- Complexity of image statistics is only partially reflected.

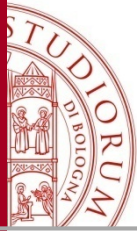


GT

noise+blur

TV

Learning based



Tips on supervised learning

Ingredients:

- A set of example pairs (*training set*):

$$(\mathbf{b}_1, \mathbf{x}_1), \dots, (\mathbf{b}_N, \mathbf{x}_N) \in \mathbf{B} \times \mathbf{X}$$

- An unknown target function which interpolates the training set:

$$\mathbf{h} \in \mathcal{H} := \{\mathbf{h} | \mathbf{h} : \mathbf{B} \rightarrow \mathbf{X}\} \quad \text{s.t.} \quad \mathbf{h}(\mathbf{b}_i) = \mathbf{x}_i$$

- A fixed parametric space (*hypothesis space*):

$$\mathcal{F}_\theta \subset \mathcal{H}$$

- A loss function deduced from a distance defined on \mathcal{H}

Tips on supervised learning

Goal:

- Approximate the target function by

$$\mathbf{NET}_{(\theta^*, \mathbf{N})} \in \mathcal{F}_\theta, \quad \mathbf{NET}_{(\theta^*, \mathbf{N})} : \mathbf{B} \rightarrow \mathbf{X}$$

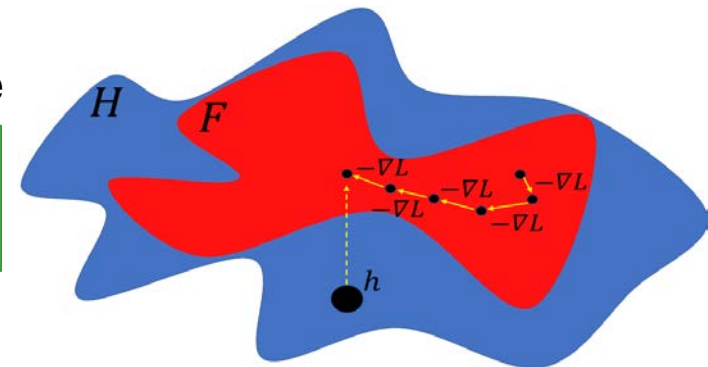
that is

$$\mathbf{b} \in \mathbf{B} \quad \mathbf{NET}_{(\theta^*, \mathbf{N})}(\mathbf{b}) \approx \mathbf{h}(\mathbf{b})$$

Approximation by gradient flow...

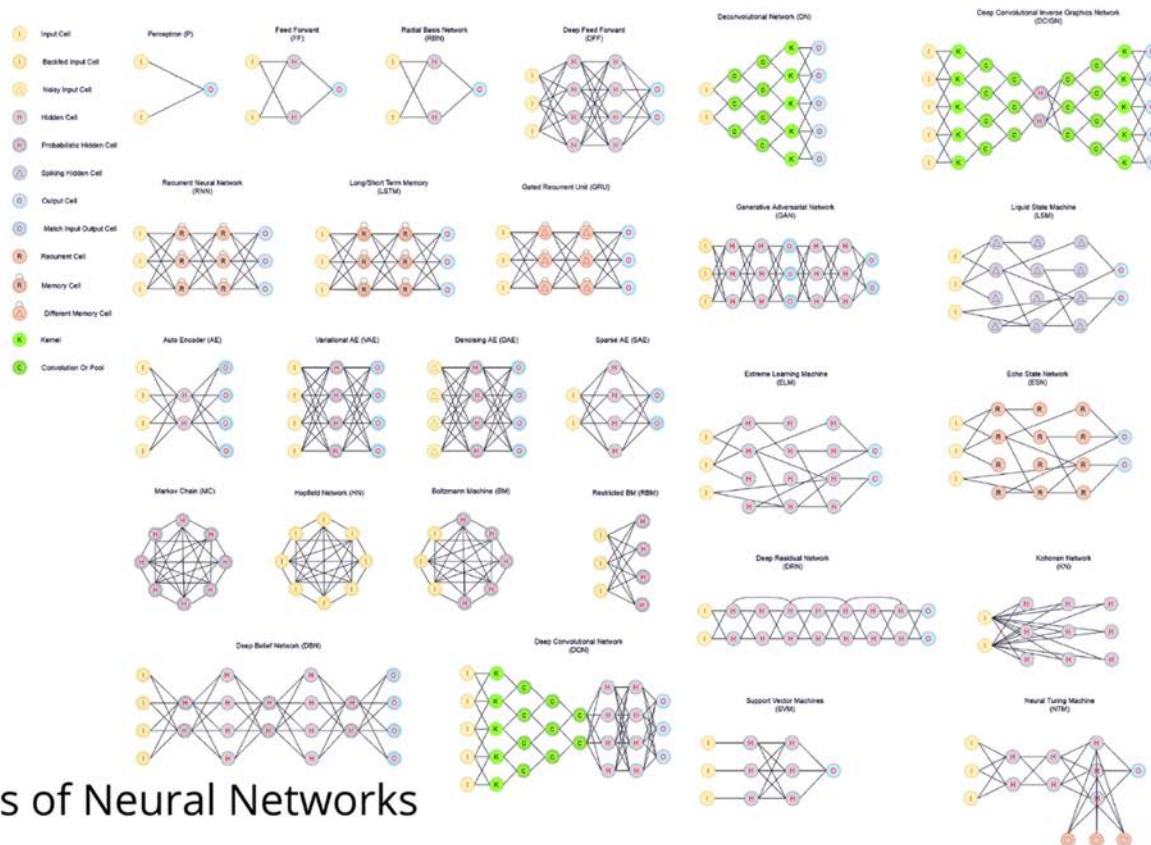
For example: Assuming the standard L_2 -norm distance

$$\mathbf{NET}_{\theta^*} \in \arg \min_{\mathbf{NET}_\theta \in \mathcal{F}_\theta} \frac{1}{N} \sum_{i=1}^N \|\mathbf{NET}_\theta(\mathbf{b}_i) - \mathbf{h}(\mathbf{b}_i)\|_2^2$$



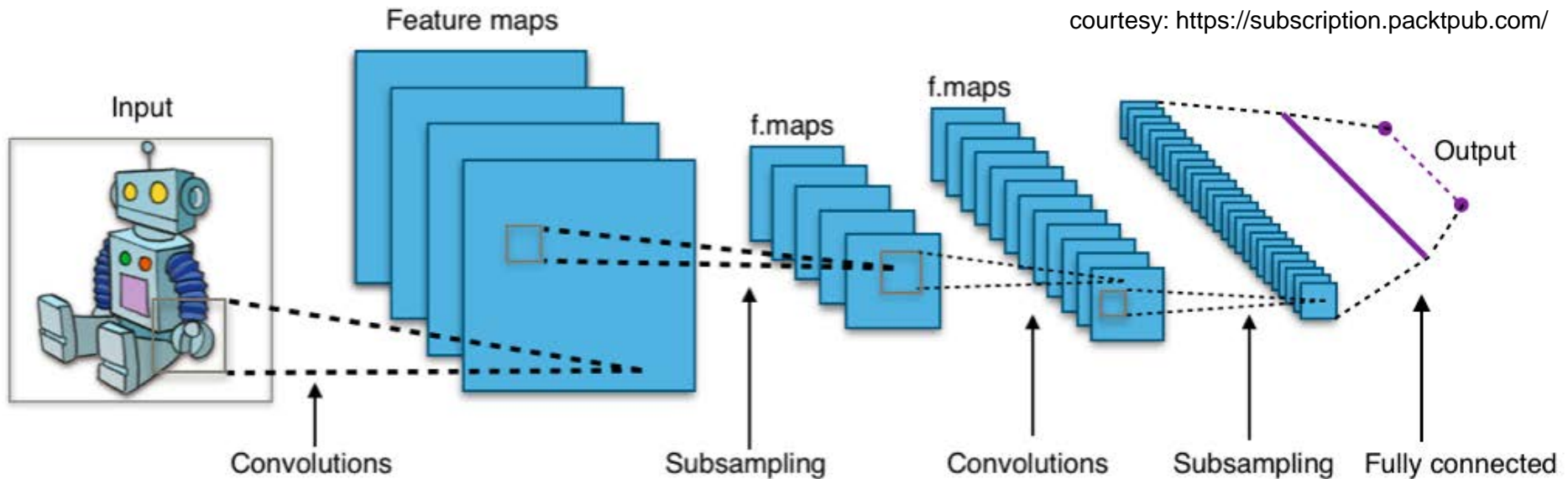
A plenty of NETs

- $\text{NET}_{\theta} = f_{\theta_1}^1 \circ f_{\theta_2}^2 \circ \dots \circ f_{\theta_t}^t$, $f_{\theta_i}^i$ is called *layer*



Main Types of Neural Networks

Convolutional Neural Networks

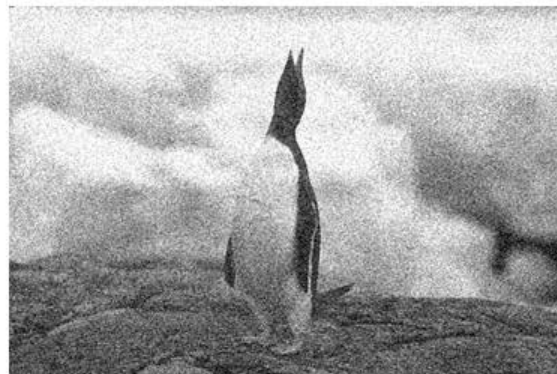


In Convolutional Neural Networks the layers are convolutions

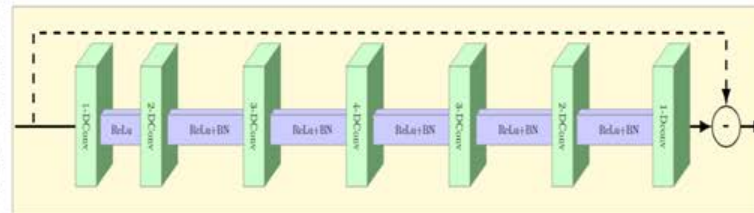
An example: image denoising

- The target is the function which denoises a noisy input.
- The training set is a set of noisy-cleaned image pairs.
- The parametric space is fixed as the set of Convolutional Neural Networks with 7 convolutions, ReLu and Batch Normalization.
- The target is approximated by solving:

$$\mathbf{NET}_{\theta^*} \in \arg \min_{\mathbf{NET}_{\theta} \in \text{CNN}_{\theta}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{NET}_{\theta}(\mathbf{b}_i) - \mathbf{x}_i\|_2^2$$



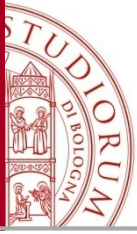
Noisy image \mathbf{b}



Neural Network Denoiser
 \mathbf{NET}_{θ^*}



Denoised image $\hat{\mathbf{x}} = \mathbf{NET}_{\theta^*}(\mathbf{b})$



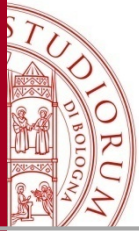
Learning Pros and Cons

PROS

- The learning process manages a huge number of training data and allows to capture a prior which reflects the complexity of image statistics.
- After training, the computation is really fast.
- Currently the state-of-the-art for several inverse problems related to images even with heavily degraded acquisition.

CONS

- Different image restoration/reconstruction tasks require different deep architectures.
- Unfortunately not all the real applications can count on a large number of training examples



Question

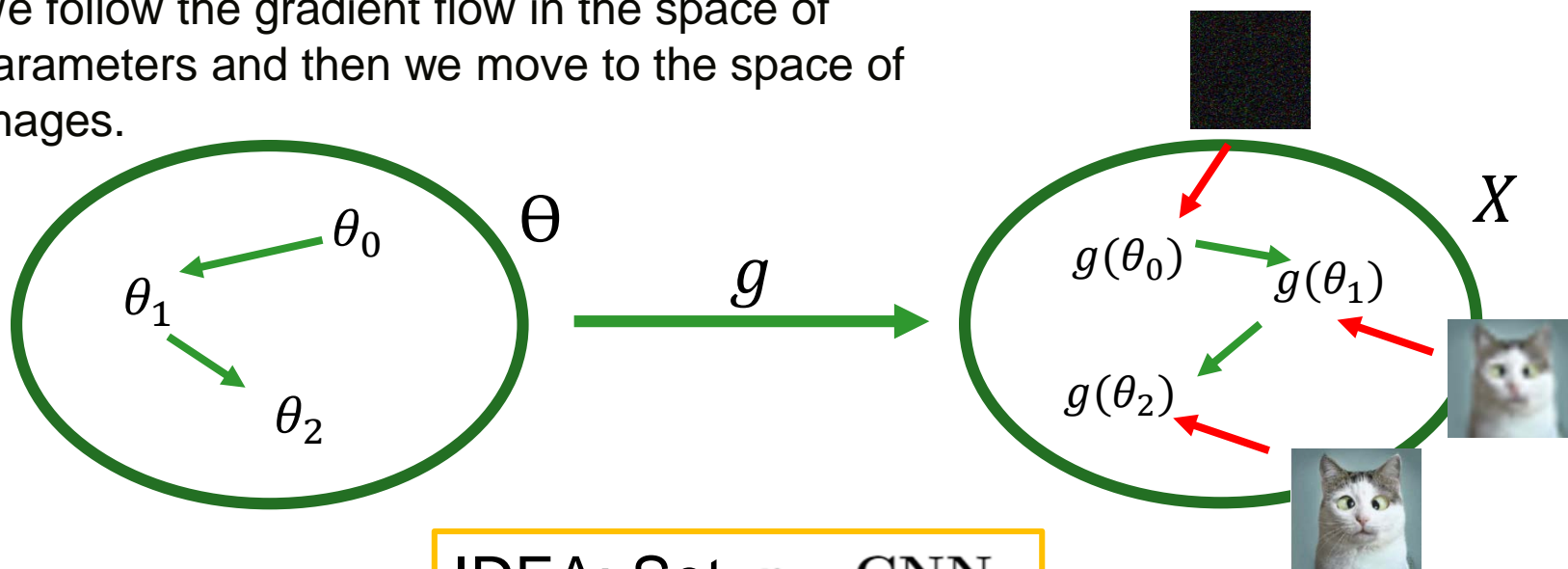
Is there a way to avoid the training set?

Parametrizing the space of images

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda R(\mathbf{x}) \quad \longrightarrow \quad \theta^* = \arg \min_{\theta \in \mathbb{R}^s} \frac{1}{2} \|\mathbf{Ag}(\theta) - \mathbf{b}\|_2^2 + \lambda R(\mathbf{g}(\theta))$$

$$\text{s.t. } \mathbf{x}^* = \mathbf{g}(\theta^*)$$

- The prior is usually defined on a transformed domain, such as wavelet, gradient or overcomplete dictionary.
- We follow the gradient flow in the space of parameters and then we move to the space of images.



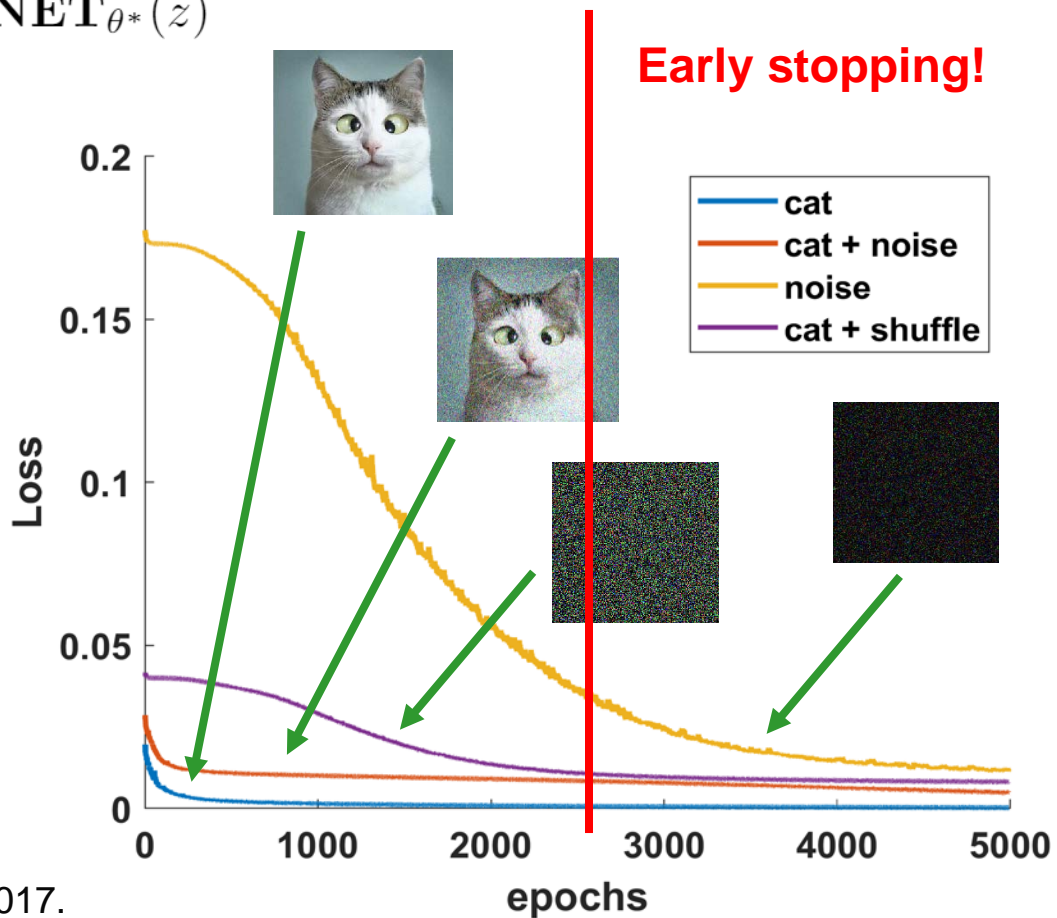
IDEA: Set $g = \text{CNN}_\theta$

Deep Image Prior (DIP)

$$\mathbf{NET}_{\theta^*} = \arg \min_{\mathbf{NET}_{\theta} \in \text{CNN}_{\theta}} \frac{1}{2} \|A \mathbf{NET}_{\theta}(z) - b\|_2^2$$

$$\text{s.t. } x^* = \mathbf{NET}_{\theta^*}(z)$$

High noise
impedance property



Deep Image Prior, Dimitri Ulyanov et al., 2017.

High noise impedance property

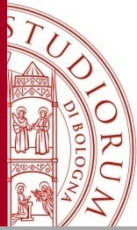


LR noisy baby



GT baby

Stopping criterion needed!



ADMM-DIP TV

- Solving the constrained formulation of DIP TV functional using ADMM.

From an unconstrained formulation...

$$\mathbf{NET}_{\theta^*} = \arg \min_{\mathbf{NET}_{\theta} \in \text{CNN}_{\theta}} \frac{1}{2} \|A \mathbf{NET}_{\theta}(z) - b\|_2^2 + \lambda \sum_{i=1}^N \sqrt{(D_h \mathbf{NET}_{\theta}(z))_i^2 + (D_v \mathbf{NET}_{\theta}(z))_i^2}$$

To a constrained formulation...

$$\begin{aligned} \mathbf{NET}_{\theta^*} = \arg \min_{\mathbf{NET}_{\theta} \in \text{CNN}_{\theta}} & \frac{1}{2} \|A \mathbf{NET}_{\theta}(z) - b\|_2^2 + \sum_{i=1}^N \|t_i\|_2 \\ \text{s.t. } & D(\mathbf{NET}_{\theta}(z)) = t \end{aligned}$$

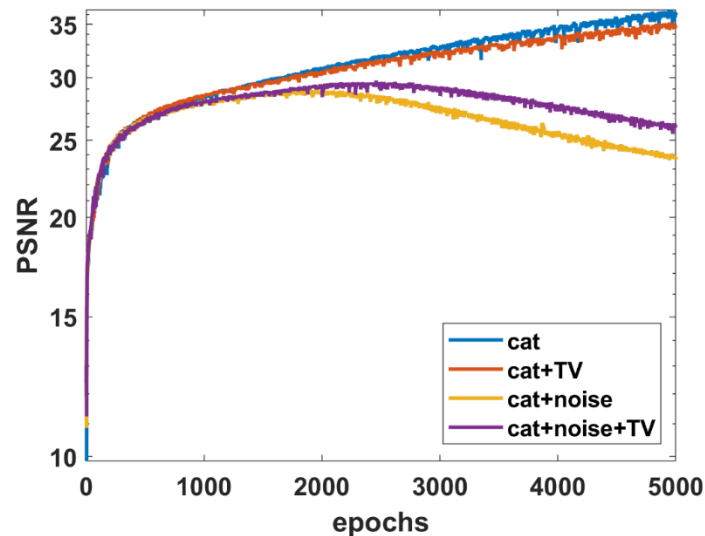
ADMM-DIP TV

$$\left\{ \begin{array}{l} \theta^{k+1} \in \arg \min_{\theta} \frac{1}{2} \|A \mathbf{NET}_{\theta}(z) - b\|_2^2 + \frac{\beta_t}{2} \|D(\mathbf{NET}_{\theta}(z)) - t^k + \frac{\lambda_t^k}{\beta_t}\|_2^2 \quad (1) \\ t^{k+1} = \arg \min_t \lambda \sum_{i=1}^N \|t_i\|_2 + \frac{\beta_t}{2} \|t - (D(\mathbf{NET}_{\theta_{k+1}}(z)) + \frac{\lambda_t^k}{\beta_t})\|_2^2 \quad (2) \\ \lambda_t^{k+1} = \lambda_t^k + \beta_t (D(\mathbf{NET}_{\theta_{k+1}}(z)) - t^{k+1}) \quad (3) \end{array} \right.$$

(1) One step of ADAM

(2) L2 proximity operator

(3) Simple update



Remark: Changing A in (1) allows to solve different image restoration tasks.

Boosting the standard DIP

ADMM-DIP TV



SSIM=0.960, PSNR=29.978

SSIM=0.943, PSNR=28.212

SSIM=0.891, PSNR=24.568

DIP



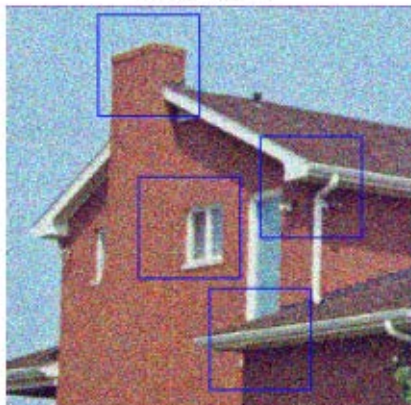
SSIM=0.946, PSNR=28.378

SSIM=0.926, PSNR=27.638

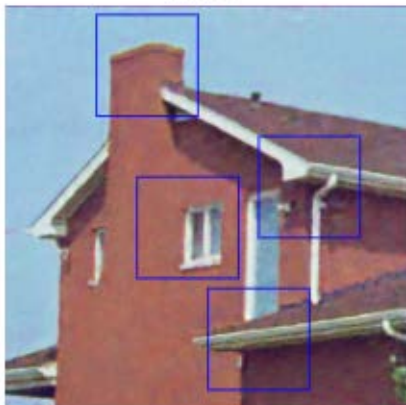
SSIM=0.864, PSNR=24.185

ADMM-DIP TV vs DIP

NOISY



DIP



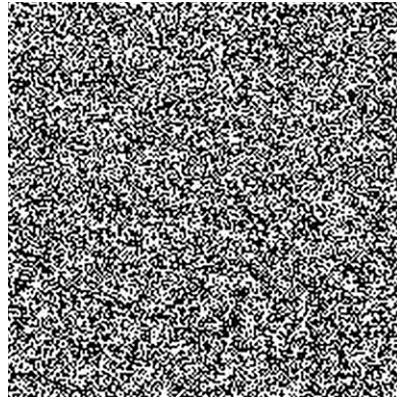
ADMM-DIPTV



Inpainting with ADMM-DIP TV



gt



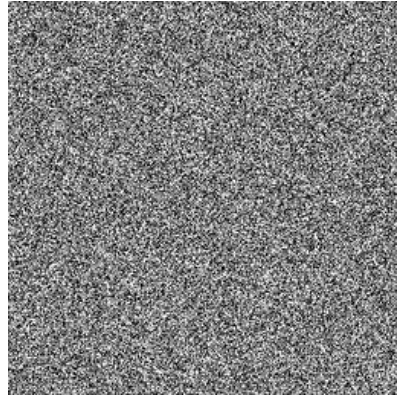
Mask (50%)



gt+mask

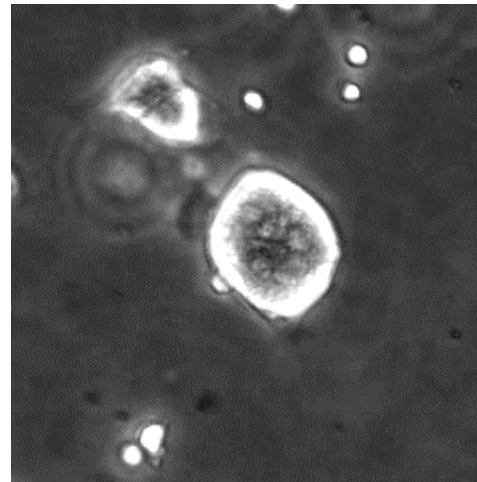


ADMM-DIP TV

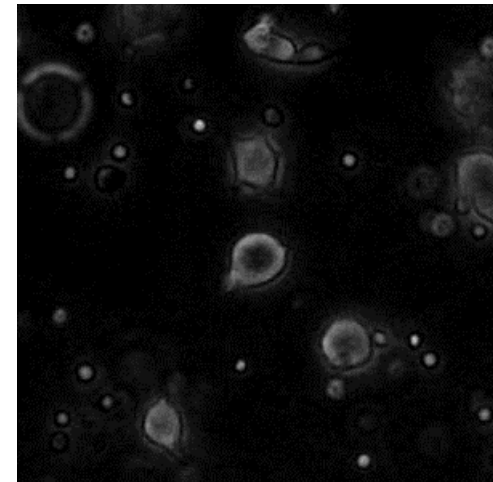


Time-Lapse Microscopy Video

- Time-Lapse Microscopy Videos are sequences of frame with biological contents whose temporal frame rate goes from seconds to minutes, acquired by a TLM microscope.
- They are used combined with automated software to track the cells.
- The way in which the cells move, indeed, has been discovered meaningful to understand wound healing, morphogenesis, cancer growth and spread of metastasis.
- For example, some motility descriptors are extracted to uncover and quantitatively evaluate the response to target therapeutic agents

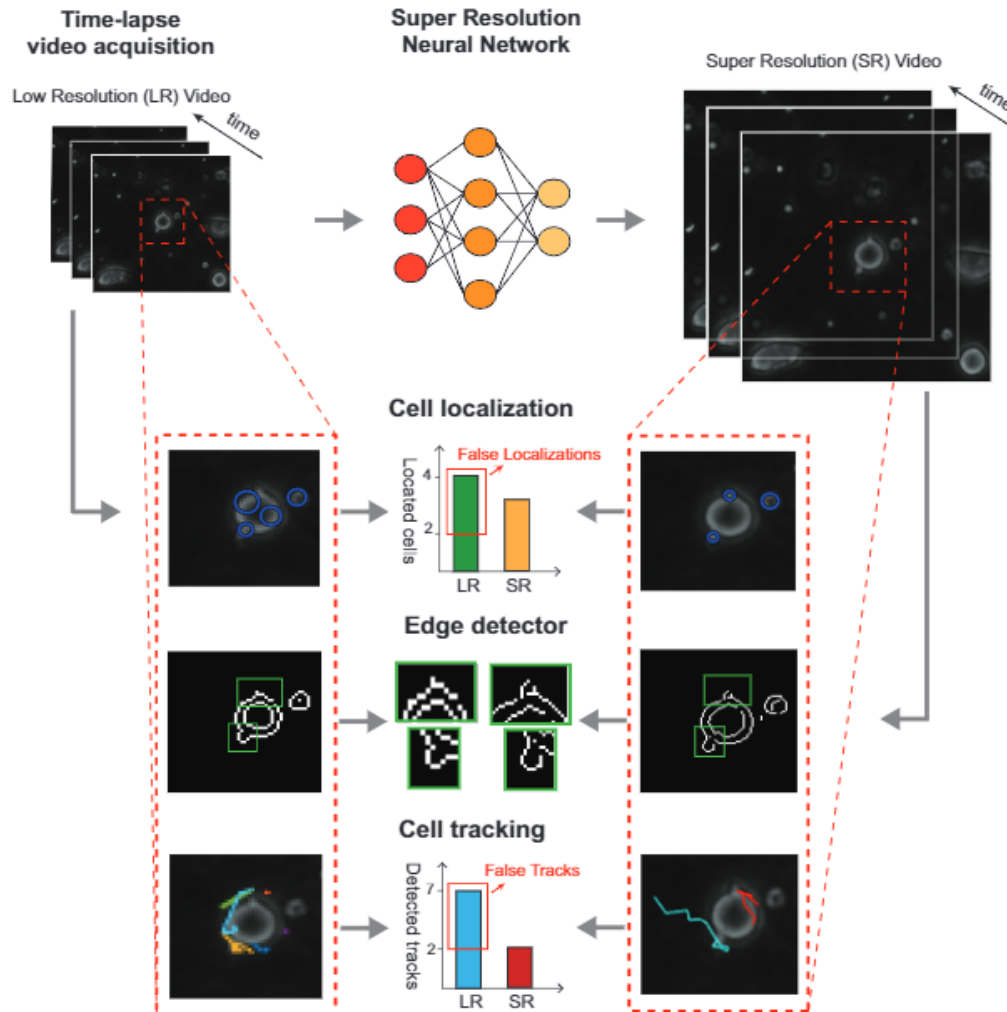


Cell Apoptosis



Breast cancer and immune cells interaction

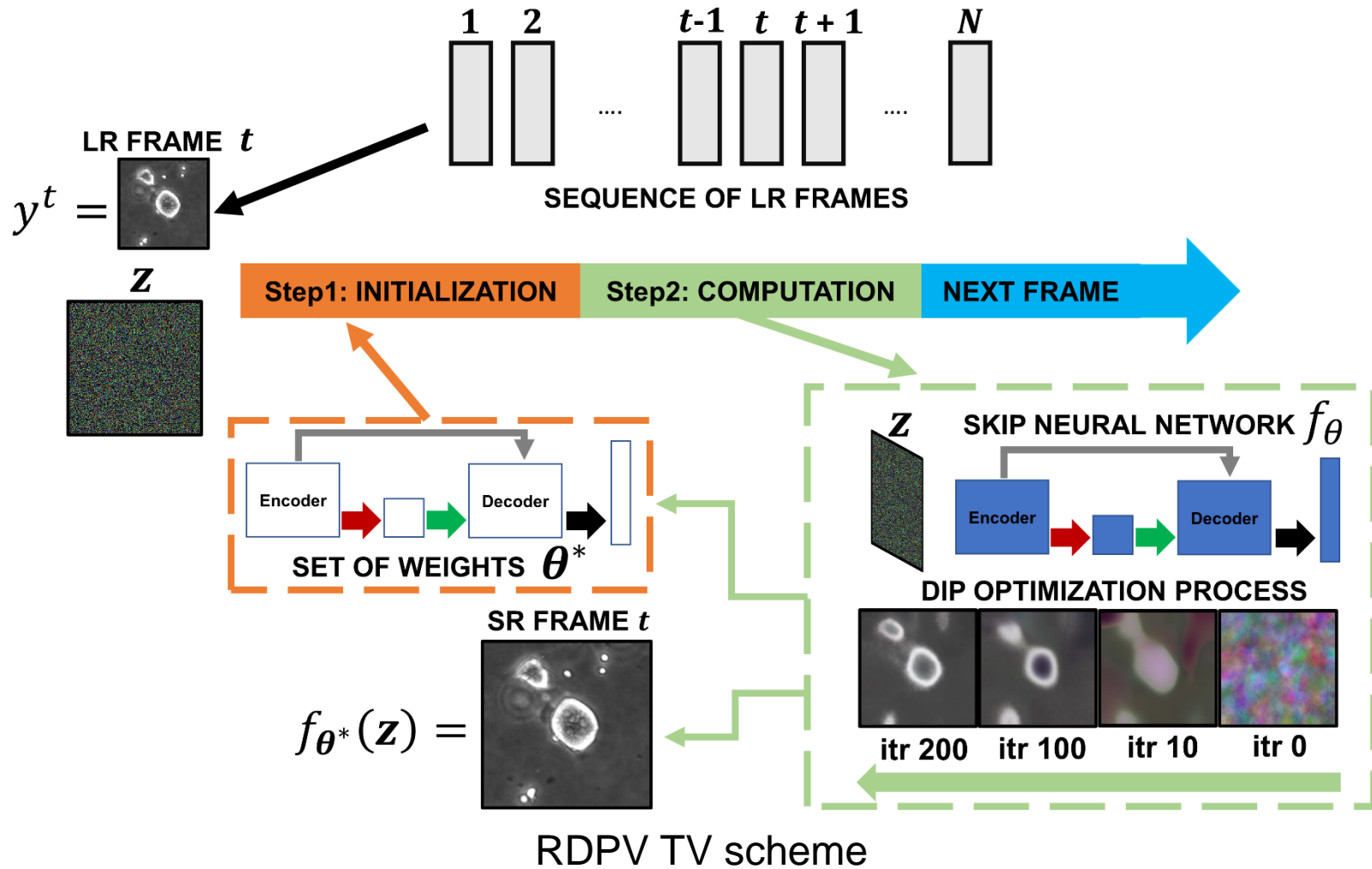
Time-Lapse Microscopy Video



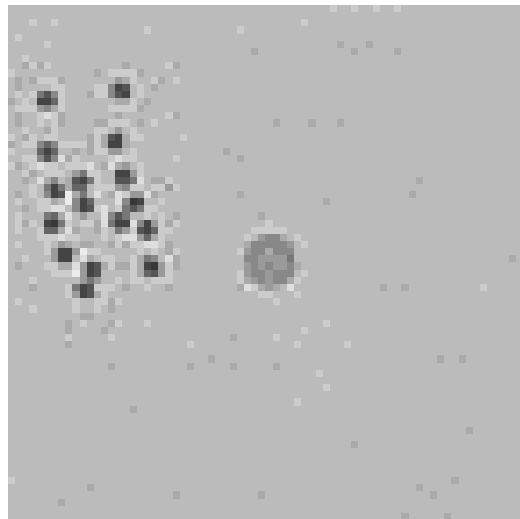
- The resolution of a TLM microscope is limited.
- An high spatial resolution positively affects the trustworthiness of the cell tracking.
- For this application it is not easy to acquire a representative dataset.

IDEA: Let's use DIPTV

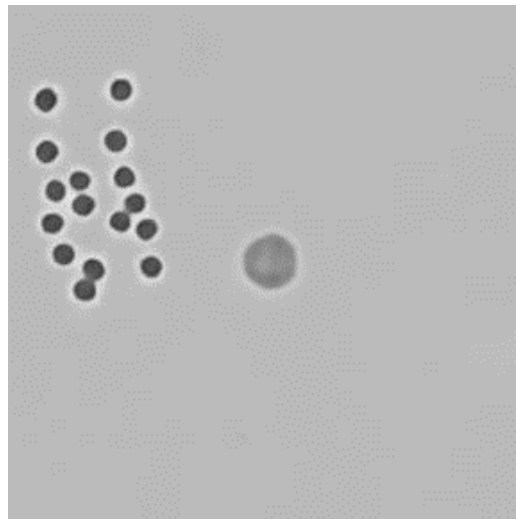
An extension of DIPTV to Video Super-Resolution



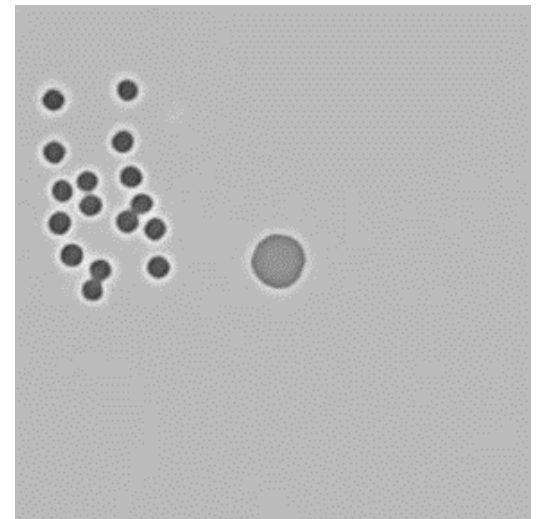
Results on synthetic Cell Video



LR TLM video



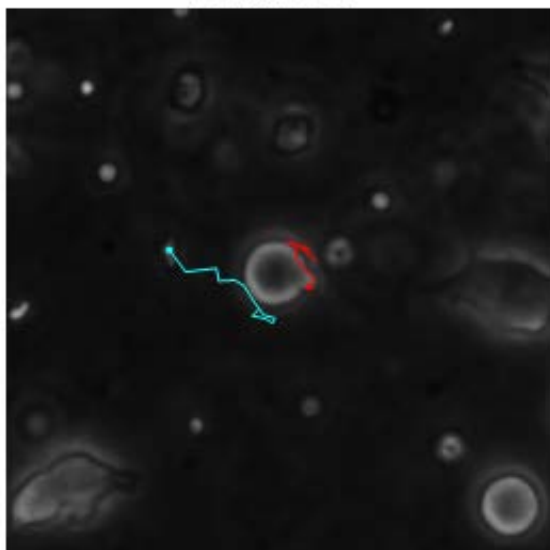
HR TLM video by DIP



HR TLM video by RDPV-TV

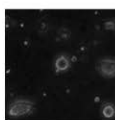
Real Cell Video

Frame 14

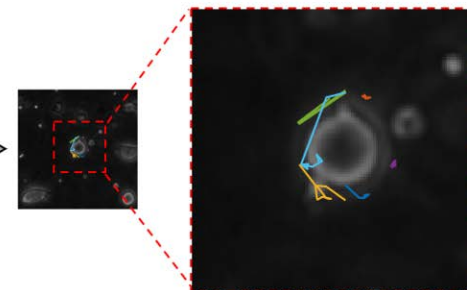


HR video trajectories

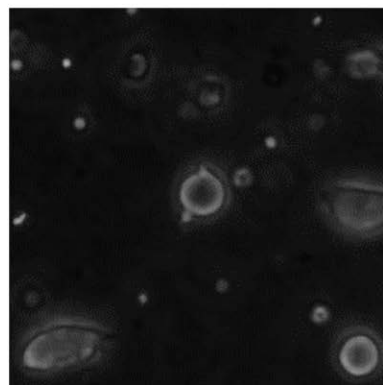
Original Video



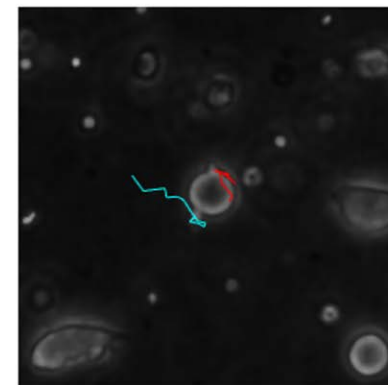
Detected trajectories



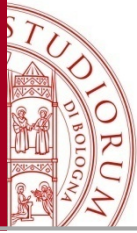
SR Video



Detected trajectories



Computed trajectories on LR and HR videos



Final remarks and future works

- The most striking property of DIP and DIPTV framework is that they do not require any training phase.
- DIP and DIPTV are useful for the solution of generic imaging inverse problems
- ADMM DIPTV leads to better results with respect to the DIP.
- DIPTV framework adapted to Time Lapse Video Microscopy Super Resolution.
- Future works will address the problem of solving the semi-convergence.

References:

1. ADMM-DIPTV: combining Total Variation and Deep Image Prior for image restoration, P.C, A. Sebastiani, M.C. Comes, arXiv.
2. Recursive Deep Prior Video: a Super Resolution algorithm for Time-Lapse Microscopy of organ-on-chip experiments, P.C, M.C. Comes et al., arXiv.



Thank for your attention