



**IMT Atlantique**

Bretagne-Pays de la Loire

École Mines-Télécom

# Motion Magnification

# SUMMARY

## 1. Introduction to Motion Magnification

- 1.1 Context
- 1.2 Key concepts
- 1.3 Problem Statement

## 2. Phase-Based Motion Processing

- 2.1 Complex Steerable Pyramids
- 2.2 Fourier transformation and motion processing

## 3. Learning-Based Motion Magnification

- 3.1 Deep Convolutional Neural Network: Architecture
- 3.2 Minimization Problem

## 4. Results and Evaluations

- 4.1 Training Dataset
- 4.2 Comparison of Two Approaches
- 4.3 Comparison with the State-of-the-Art
- 4.4 Limitations



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# Introduction to Motion Magnification



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# Introduction to Motion Magnification

4

## 1.1 Context

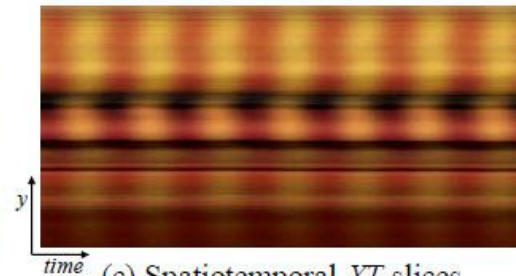
**Discerning small motions** => important applications (understanding a building's structural health, measuring a person's vital sign)



(a) Input



(b) Magnified



(c) Spatiotemporal YT slices

*Example: visualizing the human pulse, by amplifying the periodic color variation*

# Introduction to Motion Magnification

## 1.1 Context

5



*Example: effect of the wind on a crane*

**Difficulty:** How can we distinguish between small motion and noise ?

=> Previous video motion techniques suffer from **noisy outputs** and **excessive blurring**

Their **principle**: decomposing video frames into representations magnifying motion, using hand-designed filters

> Can we use better filters than hand-designed ones ?

**Goal of the paper:** learning decomposition filters using deep convolutional neural networks (CNN)

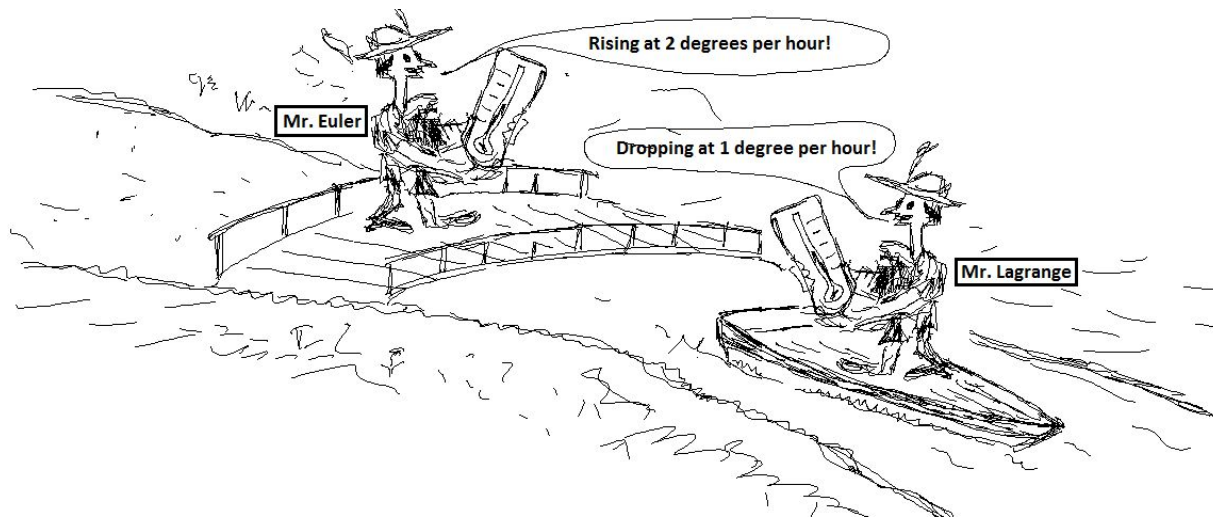
**Dataset:** a synthetic one simulating small motion

→ Design of a **network** made up of the spatial decomposition filters, the representation manipulator, and the reconstruction filters

**Training:** two-frames input, magnified difference as the target

**Contributions:** high-quality magnification, learned filters generalizing well in real videos

The Eulerian approach vs the Lagrangian approach :



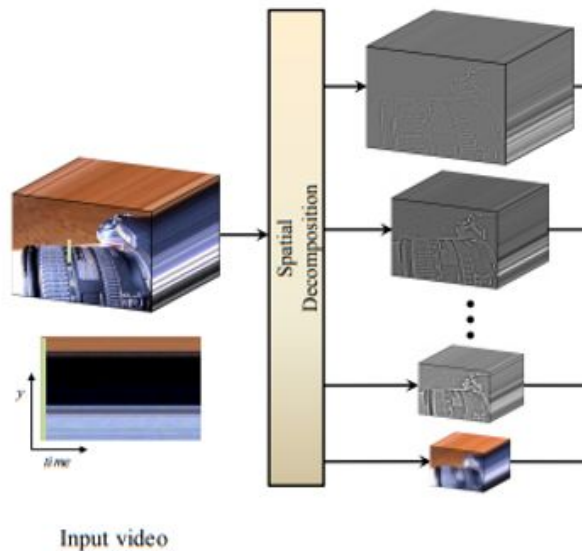
**Lagrangian method** : extraction the motion field

**Eulerian method** : representations to manipulate motions  
without needing to track the motion through space and time



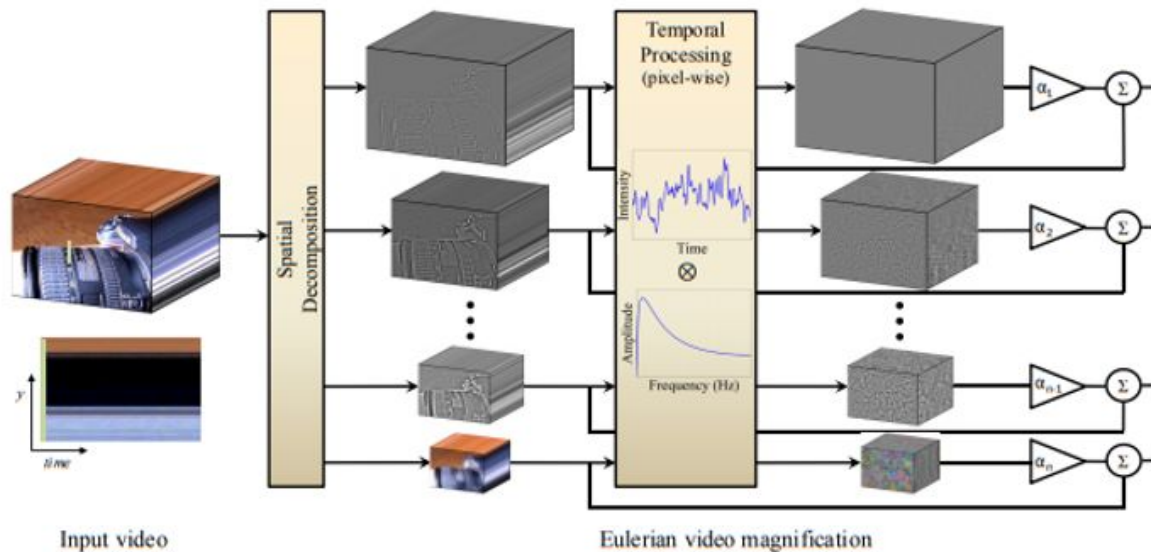
## 1.2 Key Concepts

The pipeline :



Filtering to get multiple spatial scales or frequency bands corresponding to different levels of detail

The pipeline :



Temporal filtering :  
isolate specific  
motions

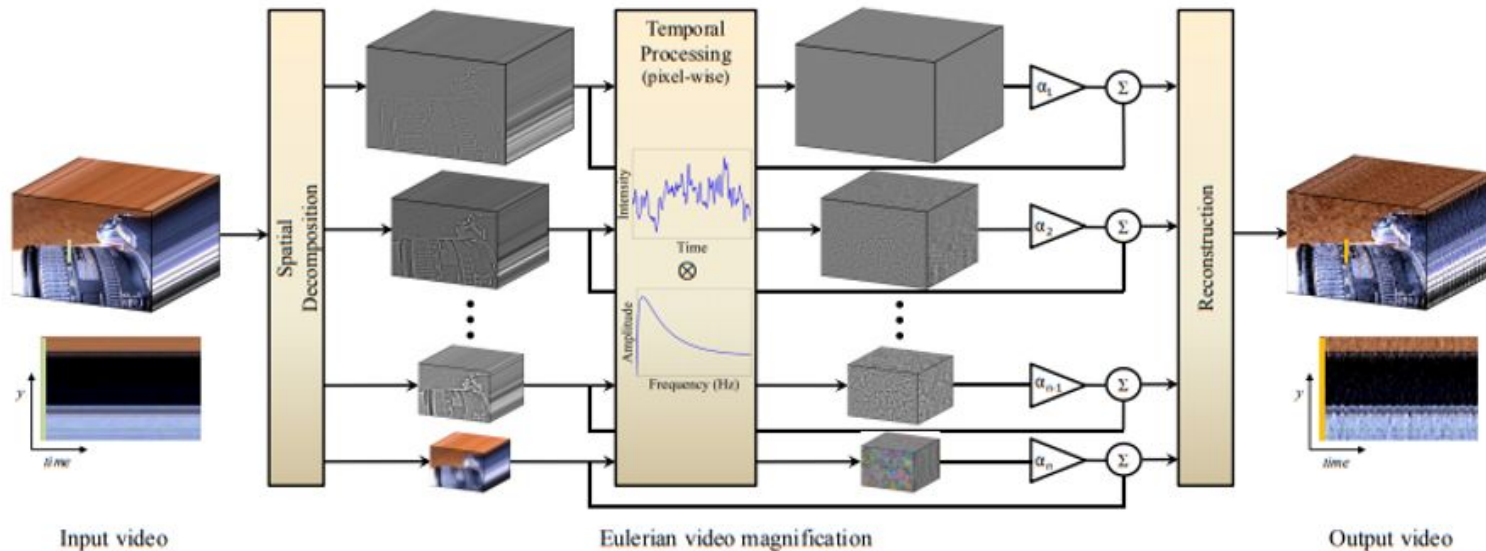
Amplification of the  
motion

# Introduction to Motion Magnification

## 1.2 Key Concepts

11

The pipeline :



### Other methods :

- . hand-designed filters like steerable filters for spatial decomposition but noise produced
- . temporal filters to isolate motion of interest and to not amplify the noise

### Our method :

- . decomposition filters are learned from examples with deep convolutional networks (much less noise)

**Original images :**

**Magnified image :**

$$I(x, t_0) = f(x) \text{ and } I(x, t_1) = f(x + \delta(t_1)) \rightarrow I'(x, t_1) = f(x + (1 + \alpha)\delta(t))$$

- .  $\delta(t)$  represents the motion field : we suppose a global translation over time
- .  $\alpha$  is the magnification factor

Magnification for a motion of interest :

$$\delta'(x, t) = T(\delta(x, t))$$

$T(\cdot)$  is a temporal bandpass filter

Motion Magnification

# Phase-Based Motion Magnification



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# 2. Phase-based Motion Magnification

15

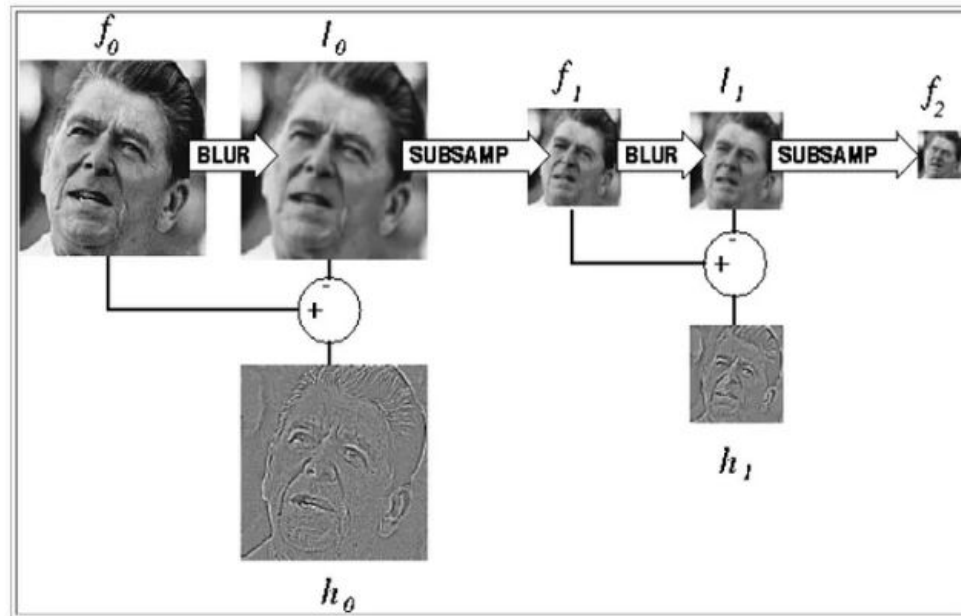
## 2.1 Complex steerable pyramids

An **Image Pyramid** is a multi-scale decomposition of an image.

- Enhanced Information Retrieval
- Scale-Sensitive Analysis
- Improved Feature Detection

A common example is the **Laplacian Pyramid**.

- Levels consist of residuals between the original and Gaussian blurred images.
- Residuals ( $h_0$ ,  $h_1$ ) act as band-pass filtered images.
- Lower scales correspond to lower frequencies.



Example of the Laplacian pyramid

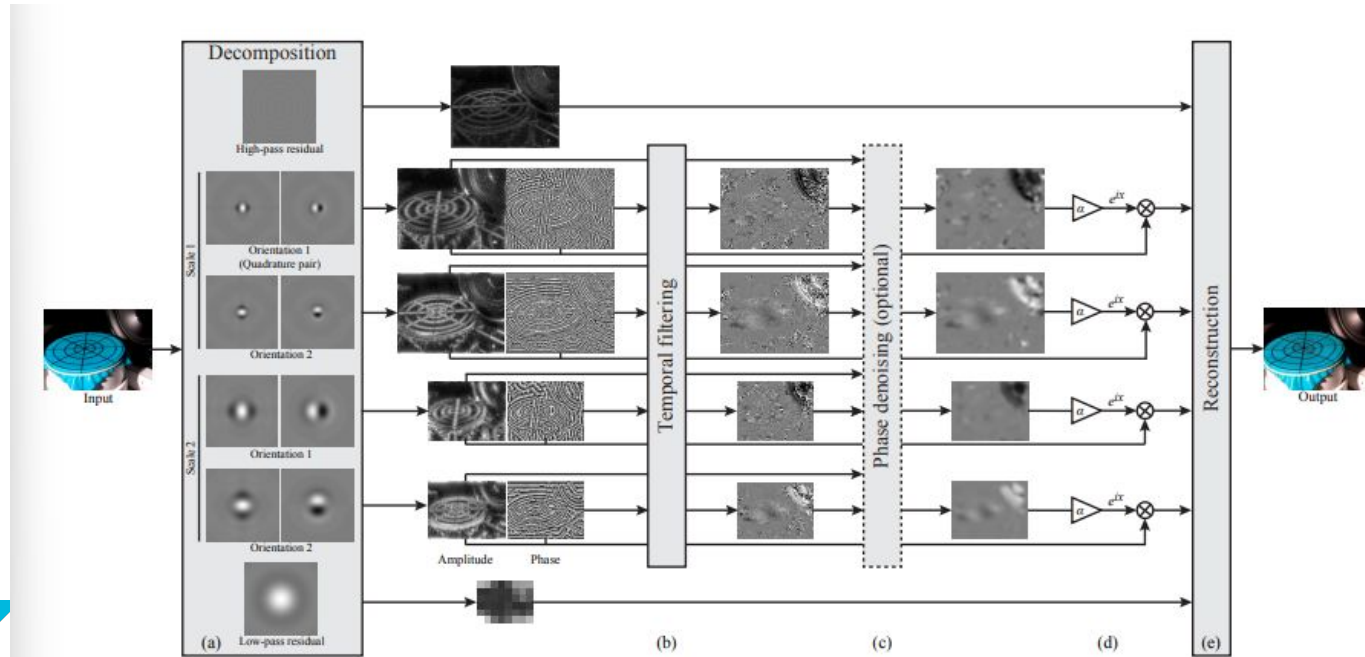
(<https://medium.com/@itberrios6/steerable-pyramids-6bfd4d23c1>)

# 2. Phase-based Motion Magnification

16

## 2.1 Complex steerable pyramids

**A Steerable Pyramid is a collection of filters** at different **sub-bands** which also include separate Low and High Pass components.





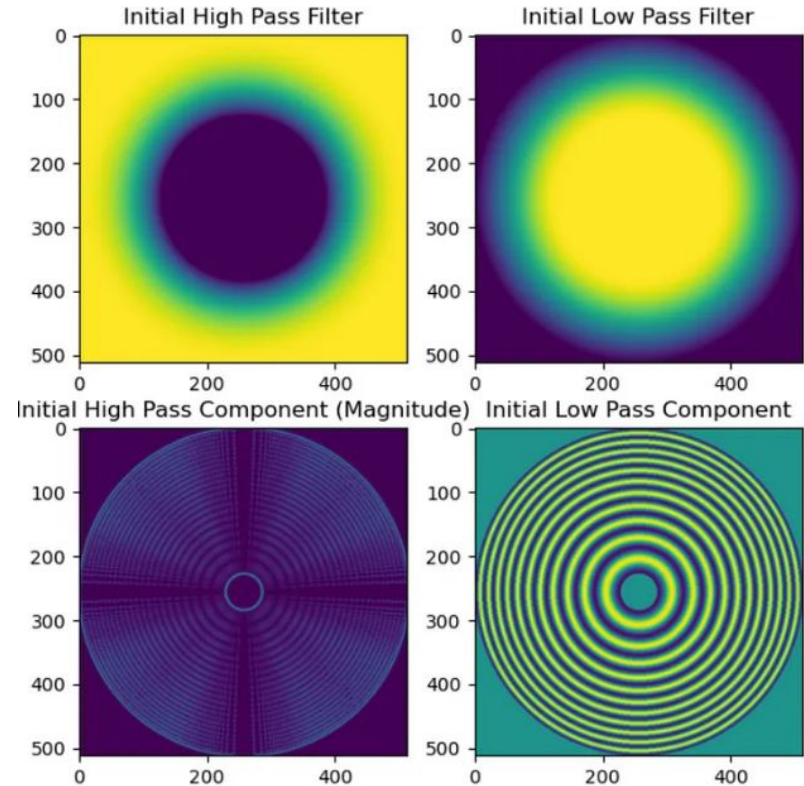
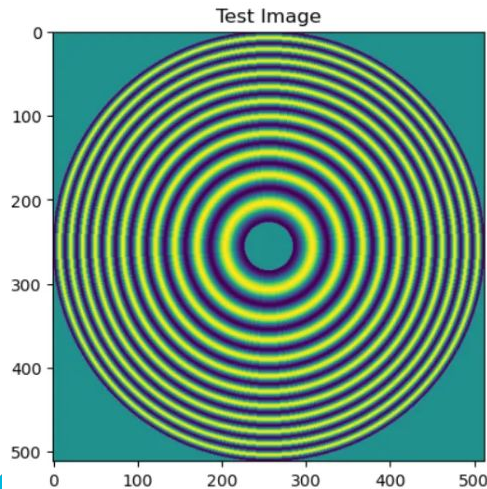
# 2. Phase-based Motion Magnification

17

## 2.1 Complex steerable pyramids

### Performing the Decomposition

- Getting the first High and Low Pass Filters



# 2. Phase-based Motion Magnification

18

## 2.1 Complex steerable pyramids

- **Getting the Sub-Band Filters**
  - Filters in a Steerable Pyramid's sub-bands must be **polar separable**.

⇒ This entails the ability to decompose them into radial and angular components.

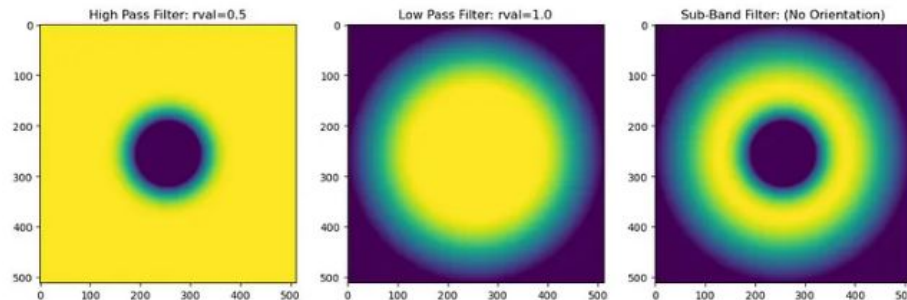
$$B_i(\vec{\omega}) = A(\theta - \theta_i)B(\omega)$$

A : the angular component w.r.t  $\theta$

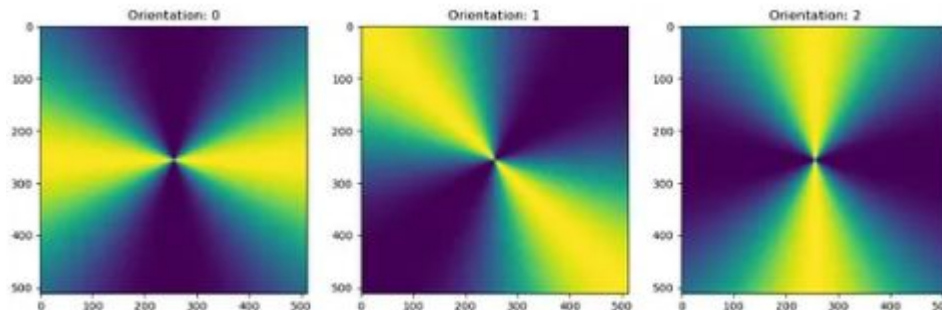
$\theta_i$  : the steering/orientation angle

B : the radial component

$\omega$  : the frequency band



*The construction of the sub-band filter  $B(w)$*

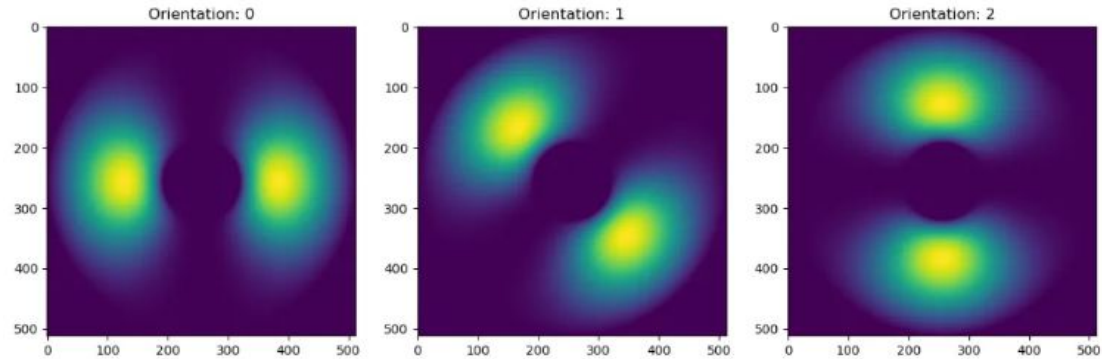
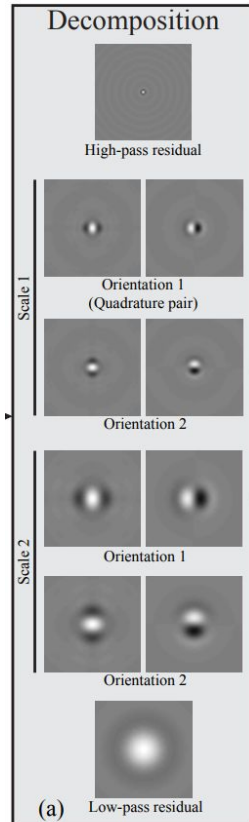


*The orientation filters (steered angles): A*

# 2. Phase-based Motion Magnification

19

## 2.1 Complex steerable pyramids



- Transfer functions (the impulse response of the filters) are applied to the discrete Fourier transform  $\hat{I}$  of an image  $I$  to create a steerable pyramid.
- This process decomposes the image into different spatial frequency bands  $S_{\omega, \theta}$  with corresponding  $\tilde{S}_{\omega, \theta}(x, y) = \tilde{I} \Psi_{\omega, \theta}$

### Conclusion:

- ⇒ Each filter isolates a continuous region in the frequency domain, resulting in a localized impulse response in space.
- ⇒ The resulting spatial frequency bands are localized in space, scale, and orientation.
- ⇒ The transfer functions of a complex steerable pyramid contain only positive frequencies.
- ⇒ This approach allows for consideration of both amplitude and phase in the response.

# 2. Phase-based Motion Magnification

## 2.2 Fourier Transformation and motion Processing

21

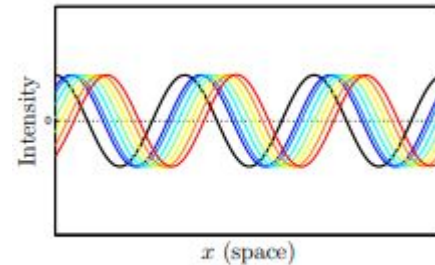
⇒ To magnify the small motion, we modify local phase variations in a complex steerable pyramid representation of the video.

A 1D image intensity profile  $f$  under global translation over time is represented as  $f(x + \delta(t))$ , where  $\delta(t)$  is some displacement function.

Goal: synthesize a sequence with modified motion,  $f(x + (1 + \alpha)\delta(t))$ , for some magnification factor  $\alpha$ .

Hint: using the Fourier decomposition,  $f(x + \delta(t)) = \sum_{\omega=-\infty}^{\infty} A_{\omega} e^{i\omega(x+\delta(t))}$

$$\Rightarrow \hat{S}_{\omega}(x, t) = S_{\omega}(x, t) e^{i\alpha B_{\omega}} = A_{\omega} e^{i\omega(x+(1+\alpha)\delta(t))}$$



(c) Phase-based

# Learning-Based Motion Magnification

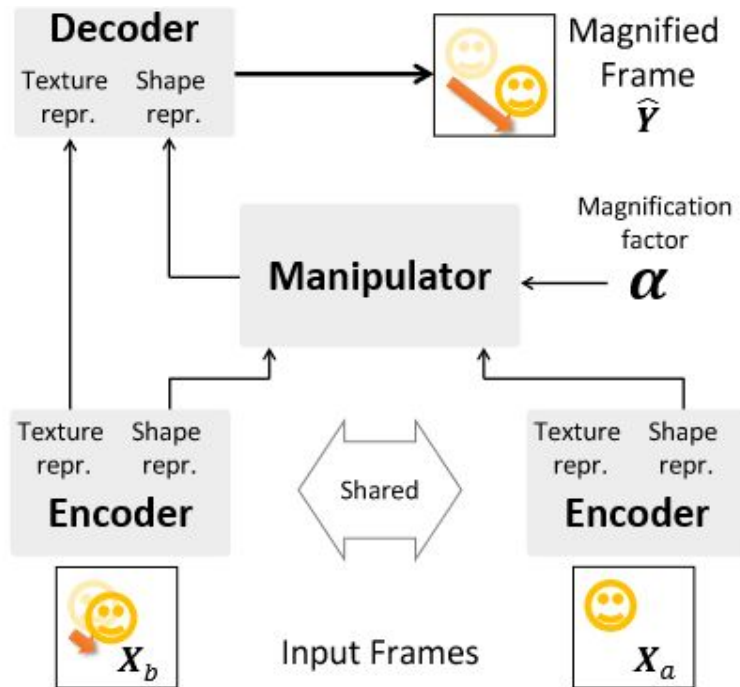


**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

# Learning-Based Motion Magnification

23

## 3.1. Deep Convolutional Neural Network - Architecture



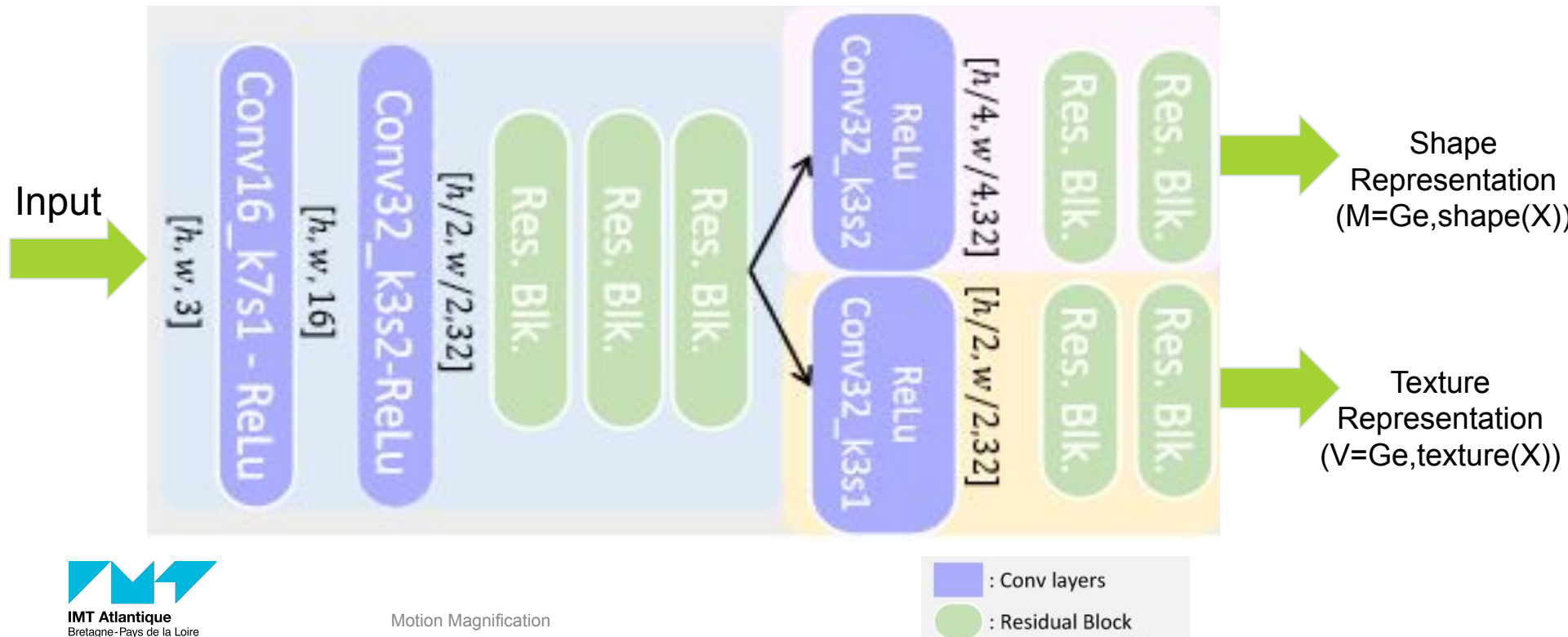


# Learning-Based Motion Magnification

## 3.1. Deep Convolutional Neural Network - Architecture

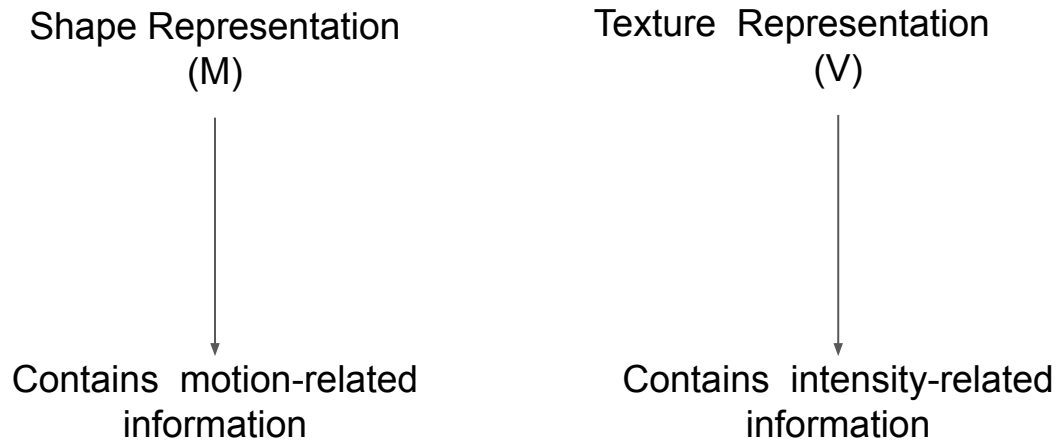
24

### Encoder

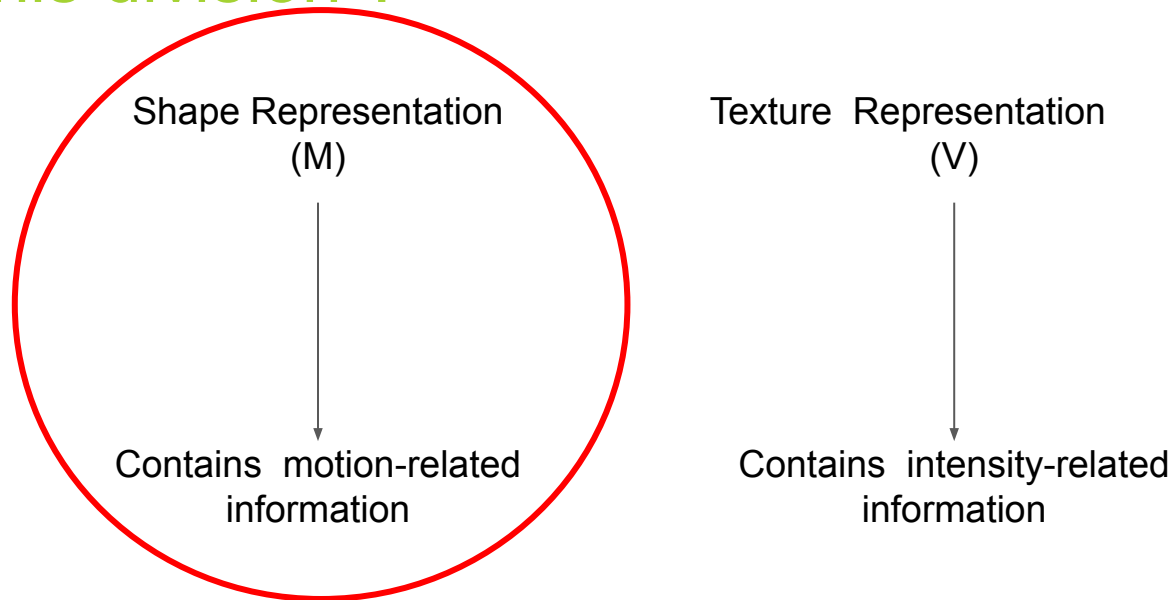




### Why this division ?



### Why this division ?



**We need to magnify M & Avoid intensity Magnification**

### Manipulator - 2 Frames

$$G_m(\mathbf{M}_a, \mathbf{M}_b, \alpha) = \mathbf{M}_a + \alpha(\mathbf{M}_b - \mathbf{M}_a).$$

Where

$\mathbf{M}_a$ : Shape representation of the first frame A

$\mathbf{M}_b$ : Shape representation of the second frame B

$\alpha$ : Magnification factor

### Manipulator - 2 Frames

$$G_m(\mathbf{M}_a, \mathbf{M}_b, \alpha) = \mathbf{M}_a + h(\alpha \cdot g(\mathbf{M}_b - \mathbf{M}_a))$$

Where

$\mathbf{M}_a$ : Shape representation of the first frame A

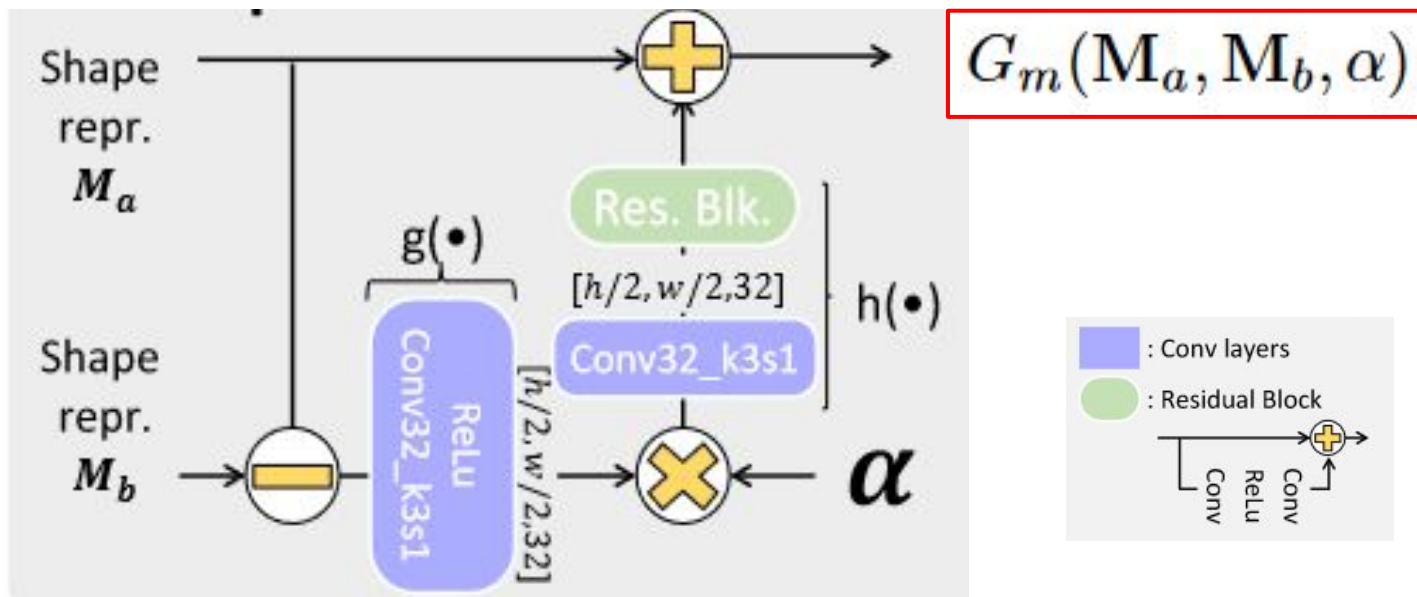
$\mathbf{M}_b$ : Shape representation of the second frame B

$\alpha$ : Magnification factor

$h(\cdot)$ : 3\*3 conv + residual block

$g(\cdot)$ : 3\*3 conv + RELU

### Manipulator - 2 Frames



### Manipulator - Temporal Filtering

$$G_{m,temporal}(\mathbf{M}(t), \alpha) = \mathbf{M}(t) + \alpha \mathcal{T}(\mathbf{M}(t)).$$

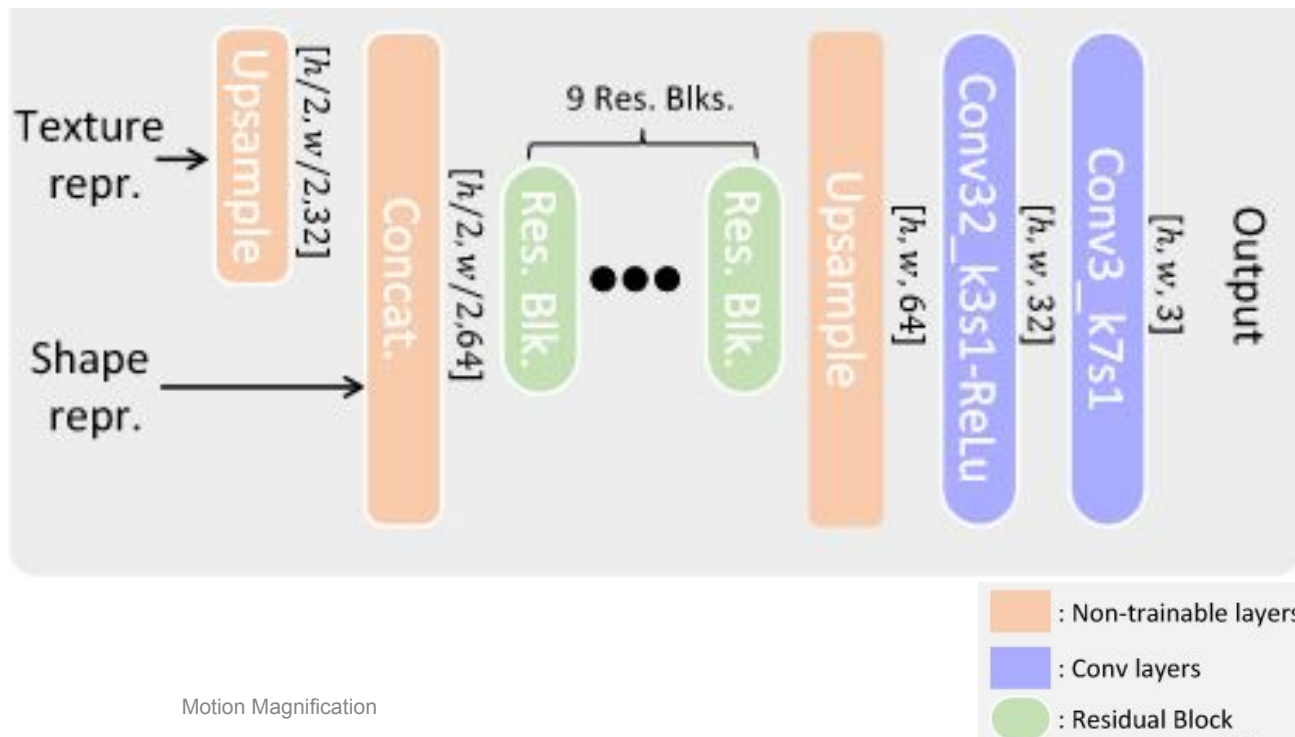
Where

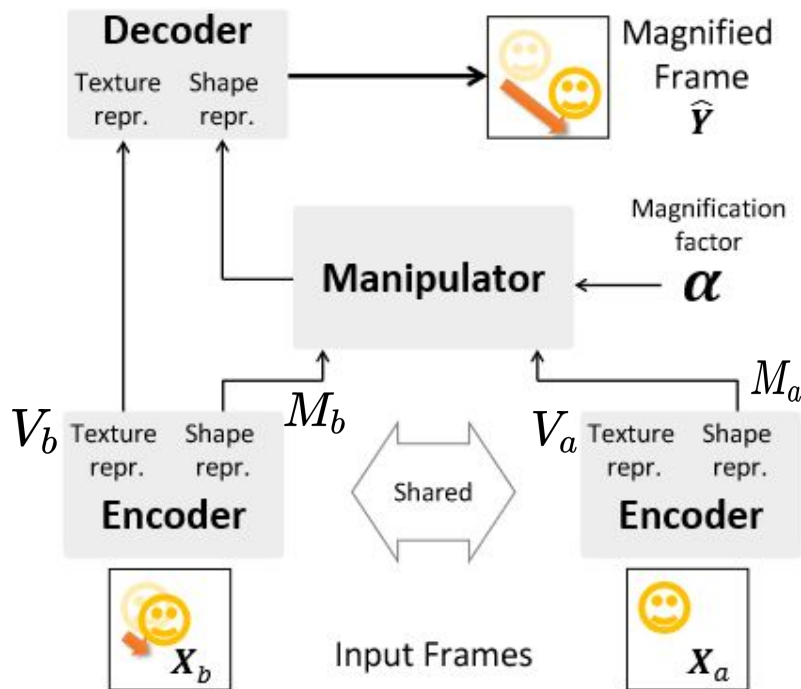
$\mathbf{M}(t)$ : Shape representation of the frame at time  $t$

$\alpha$ : Magnification factor

$\tau(.)$ : pixel-wise Temporal Filter

### Decoder

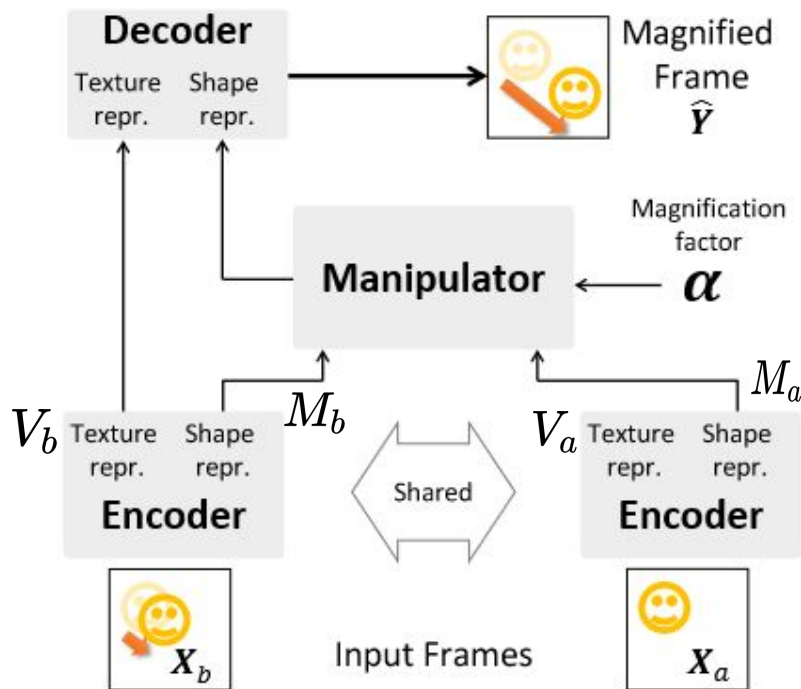




$$Loss = |\hat{Y} - Y| + \lambda(Regularisation)$$



## 3.2. Deep Convolutional Neural Network - Minimization Problem



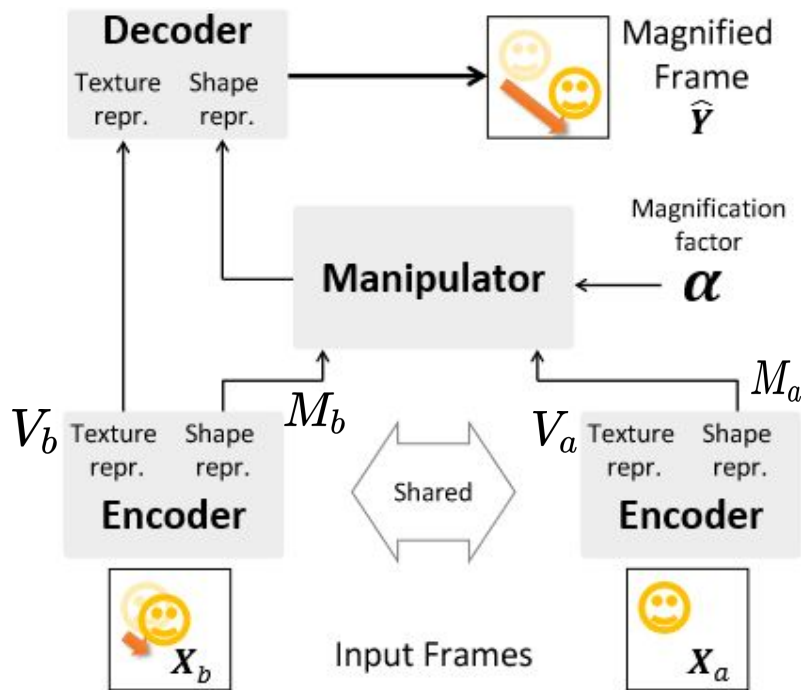
$$Loss = |\hat{Y} - Y| + \lambda(Regularisation)$$

$$Y \quad X_a \quad X_b \longrightarrow X'_a \quad X'_b \quad Y'$$

color perturbation

- Shape representation should be invariant of color perturbation
- The texture representation of B should be close of Y
- The image A and B should have close texture representation

## 3.2. Deep Convolutional Neural Network - Minimization Problem



$$Loss = |\hat{Y} - Y| + \lambda(Regularisation)$$

$$Y \quad X_a \quad X_b \xrightarrow{\text{color perturbation}} X'_a \quad X'_b \quad Y'$$

- Shape representation should be invariant of color perturbation
- The texture representation of B should be close of Y
- The image A and B should have close texture representation

$$Loss = |\hat{Y} - Y| + \lambda(|V_a - V_b| + |V'_b - V'_Y| + |M_b - M'_b|)$$

# Results and evaluations



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

## 4.1. Training dataset

Advantage of a **synthetic dataset** : large quantity

**Foreground objects and background images**: real image datasets for their realistic texture. 7,000 and 200,000 images respectively.

Each training sample: 7 to 15 foreground objects, randomly scaled from its original size. (the scaling factor limited at 2 to avoid blurry texture)

Amount and direction of motions of background and each object are randomized → learning local motions.

**Low contrast texture, global motion, and static scenes:** low performance => necessity of adding two types of example:

- the background is blurred
- only a moving background in the scene to mimic a large object

To learn changes due only to noise, addition of subsets of:

- Completely static scenes
- foreground moving but not background

## 4.1. Training dataset

**Subpixel motion generation:** managing subpixel motion, it depends on demosaicking algorithm, so there is an approach with:

- > **Reconstructing images in the continuous domain** using bicubic interpolation, before applying translation or resizing
- > **Generating images at a higher resolution** (motions appear larger) and then downsampling (to the desired size), with a Gaussian filter to reduce aliasing
- > **Adding uniform quantization noise** before final quantization to ensure minor intensity changes are preserved: each pixel has a chance of rounding up proportional to its rounding residual

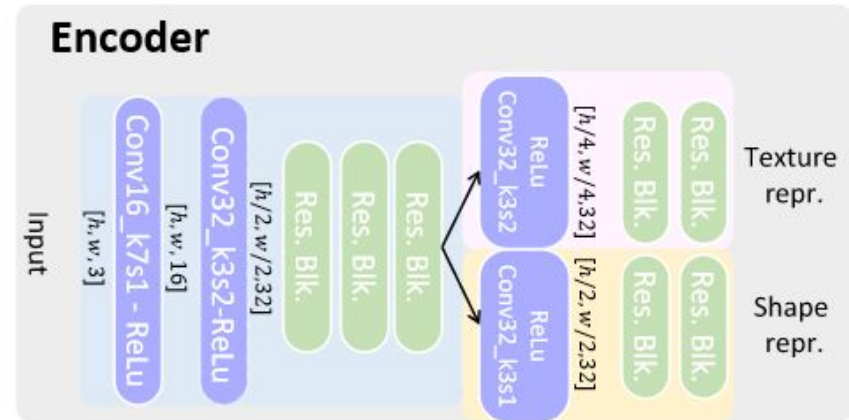
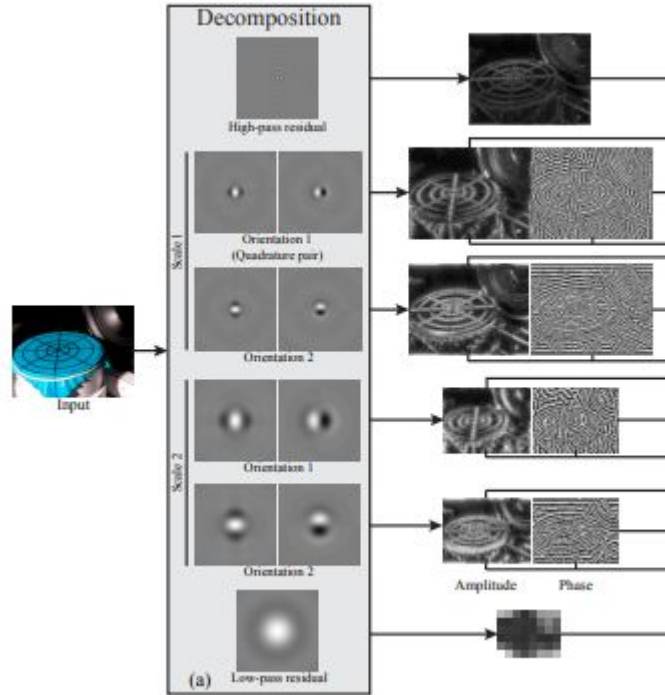
## 4.2. Comparison of Two Approaches

	Phase-based	Learning-based
<b>First step</b>	Decomposing with a Complex Steerable Pyramid	Encoder Ge
<b>Second step</b>	Filtering and Phase Removal Applying Motion Magnification	Manipulator Gm
<b>Third step</b>	Reincorporating High and Low Pass Components	Decoder Gd
<b>Target</b>	Phase of the steerable pyramid	Shape representation Ma
	Amplitude of the steerable pyramid	Texture representation Mb

# Results and Evaluations

## 4.2. Comparison of Two Approaches

40

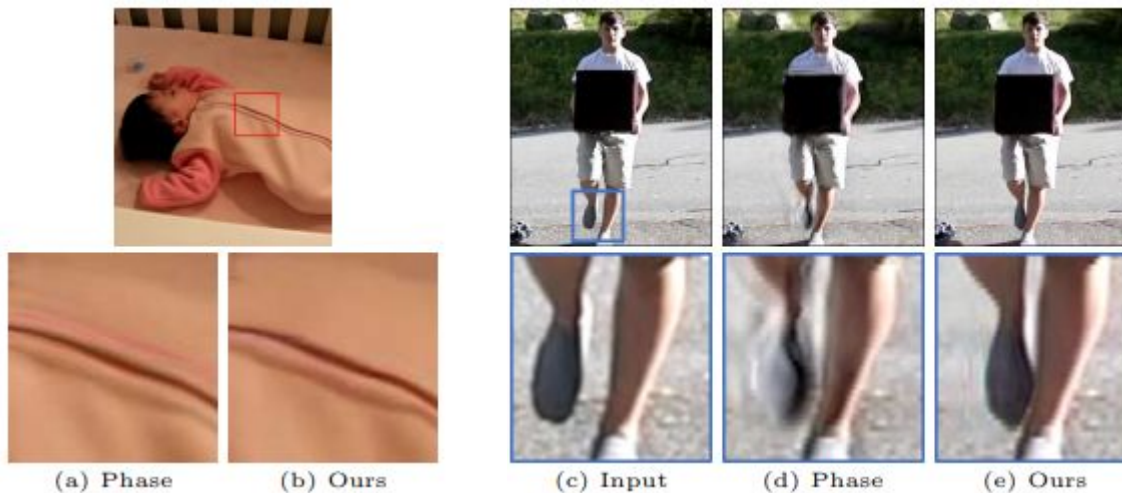




# Results and Evaluations

## 4.2. Comparison of Two Approaches

41

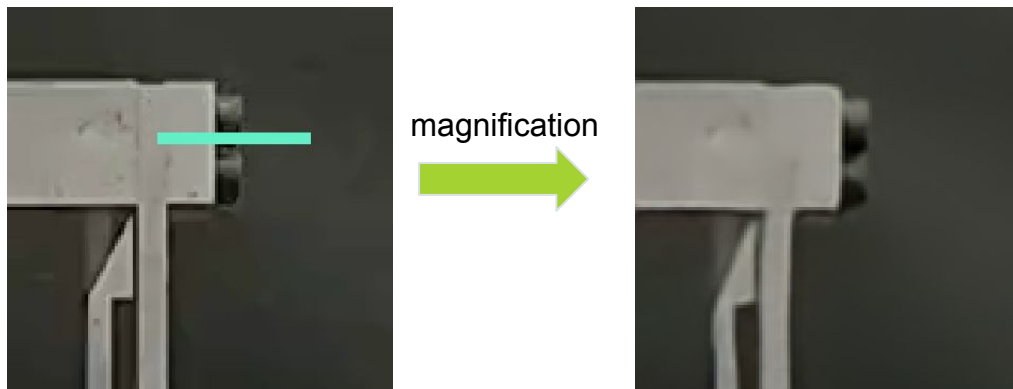


⇒ More ringing artifacts and blurring in the phase-based method than the learning based one.

⇒ Learning-based representation is trained **end-to-end** using example motion data. In contrast, the phase-based method relies on manually designed **multi-scale representations**, which struggle to accurately handle strong edges.

## 4.3. Comparison with state of the art

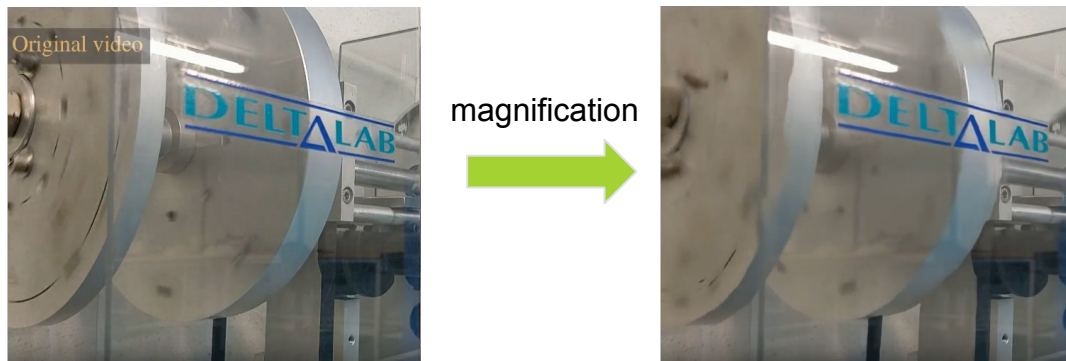
- Could not find benchmarks for motion magnification.
- However the Swin Transformer motion magnification model seems to be the most used.
- Different architecture, also trained on synthetiques training data



- Edge conservation
- Less artifacts

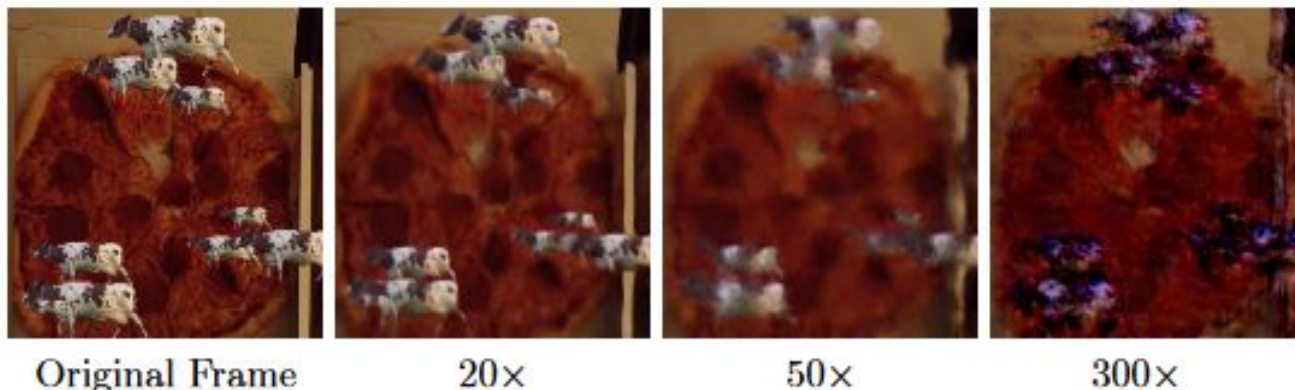
## 4.3. Comparison with state of the art

- Could not find benchmarks for motion magnification.
- However the Swin Transformer motion magnification model seems to be the most used.
- Different architecture, also trained on synthetiques training data s



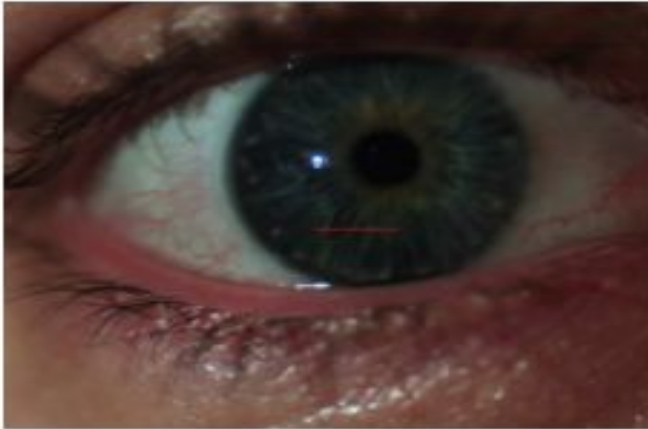
- Edge conservation
- Less artifacts

### Performance degradation with high $\alpha$



**Issue:** Blurring and color artifacts with high magnification factors.

### Blindness to very small motions



Input



Ours with temporal filter



Phase-based [24]



### Blindness to very small motions

Ours with temporal filter



Phase-based [24]

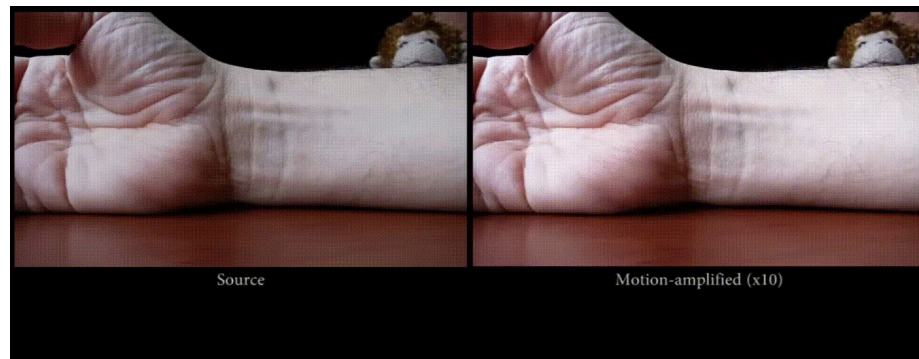


**Issue:** Occasionally magnifying the motion: “*Patchy Magnification*”

→ **Possible solution:** Improving compatibility with temporal filters.

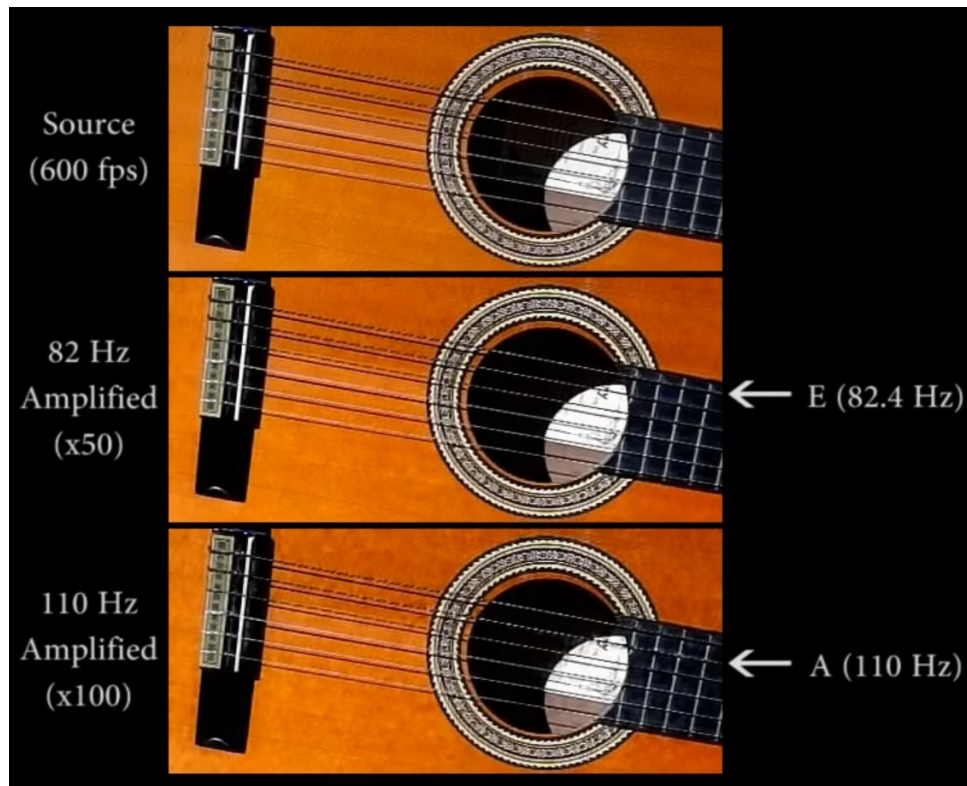
## Key takeaways from this lecture:

- **Phase-based approach:** Access to motion through the local phase of each sub band of fourier decomposition
- **Learning-based approach:** CNNs directly learn motion magnification filters from data.
- Both options present advantages and disadvantages.
- **Applications:**
  - Medical applications: heart rate pulse measurement





- Predicting acoustic wave from the motion





Thank you for your attention



**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

This paper dates from 2018, what are the new approaches ?



The synthetic dataset that they designed contains a ground truth for each sample or is this an unsupervised problem? If these ground truths exist, what are they exactly?



What are some potential real-world applications of this motion magnification method?

