



**IMT Atlantique**

Bretagne-Pays de la Loire

École Mines-Télécom

# Style Transfer by Relaxed Optimal Transport and Self-Similarity

Group 5:

Houda GHALLAB

Renzo MORALES

Carlos ARGUILAR

Mohamed Salim ARIFA

Ghaith MAKHLOUF

# SUMMARY

1. Introduction
2. Fundamentals : Optimal Transport & Self Similarity
3. Methodology
4. Implementation and region control
5. Experiments and comparison to the related work
6. Conclusion



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

# INTRODUCTION



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

# CHAPITRE 1 : INTRODUCTION

4

## 1.1 What is style transfer?



Content



Style



Combination

# CHAPITRE 1 : INTRODUCTION

5

## 1.1 What is style transfer?



## 1.2 Intuition about Style and Content

**Style:** a distribution over features extracted by a deep neural network



**Style:** a distribution over features extracted by a deep neural network



## 1.2 Intuition about Style and Content

### Content:

Self similarity: objects often have repeating patterns or elements within themselves





### Content:

Human visual system: Relative appearance and surroundings



### Content:

**Self similarity:** objects often have repeating patterns or elements within themselves

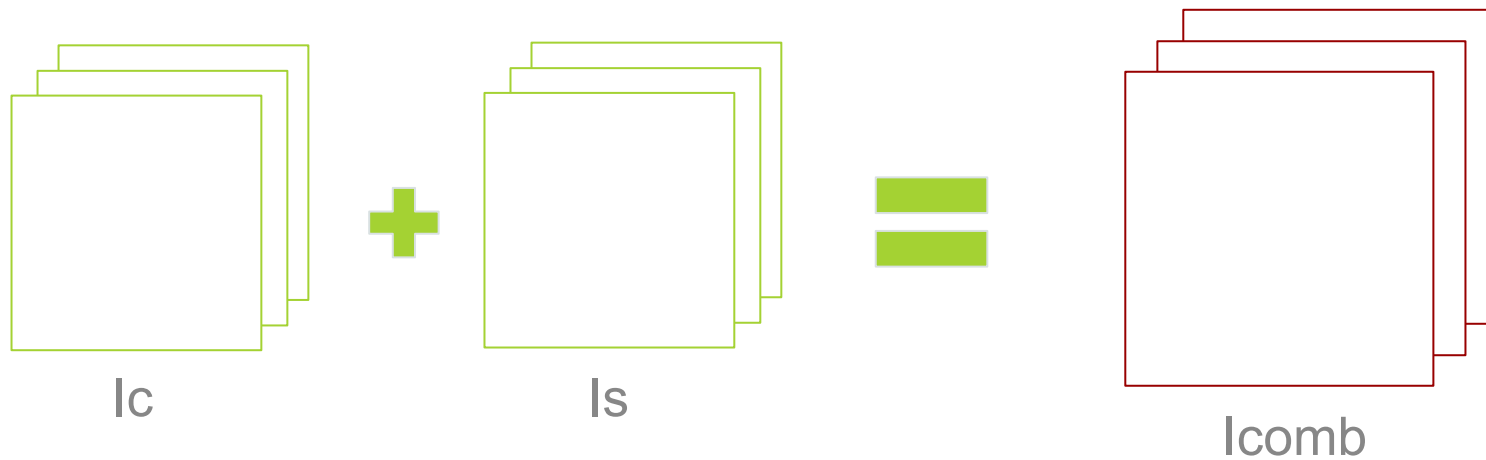
**Human visual system:**  
Relative appearance and surroundings



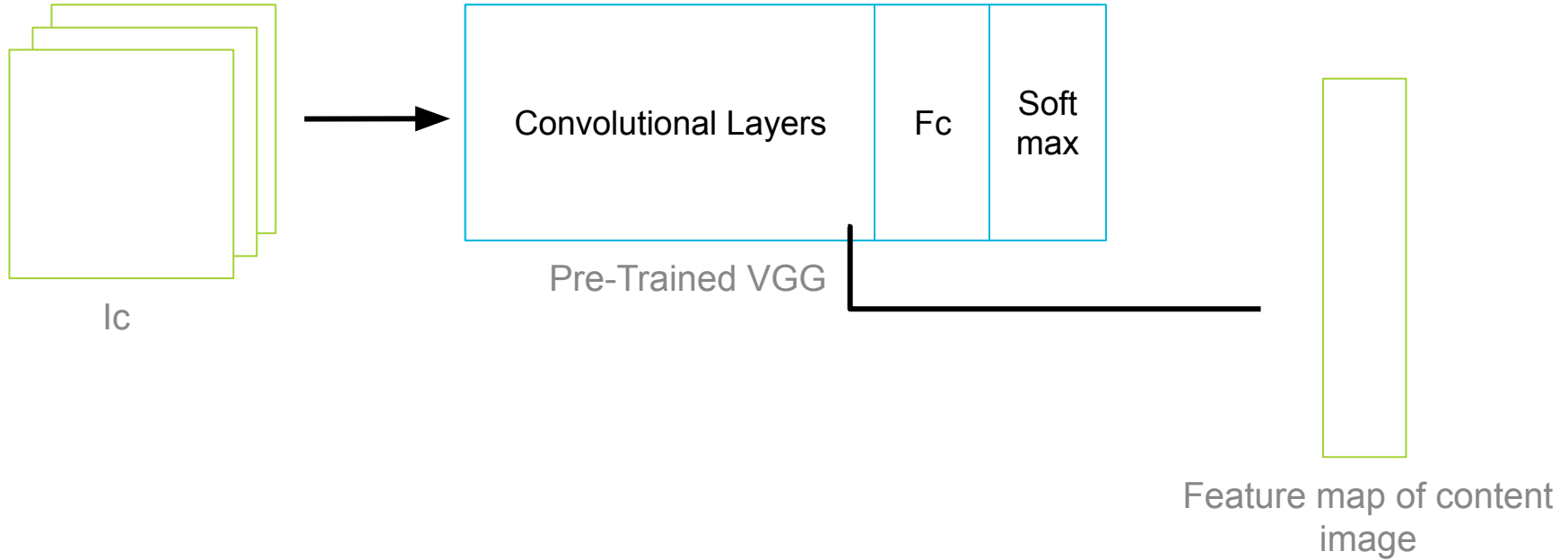
Preserve **semantics** and **spatial** layout of the content image

Content is not focused on the exact pixel values of the image, but rather on the relationships between those pixels.

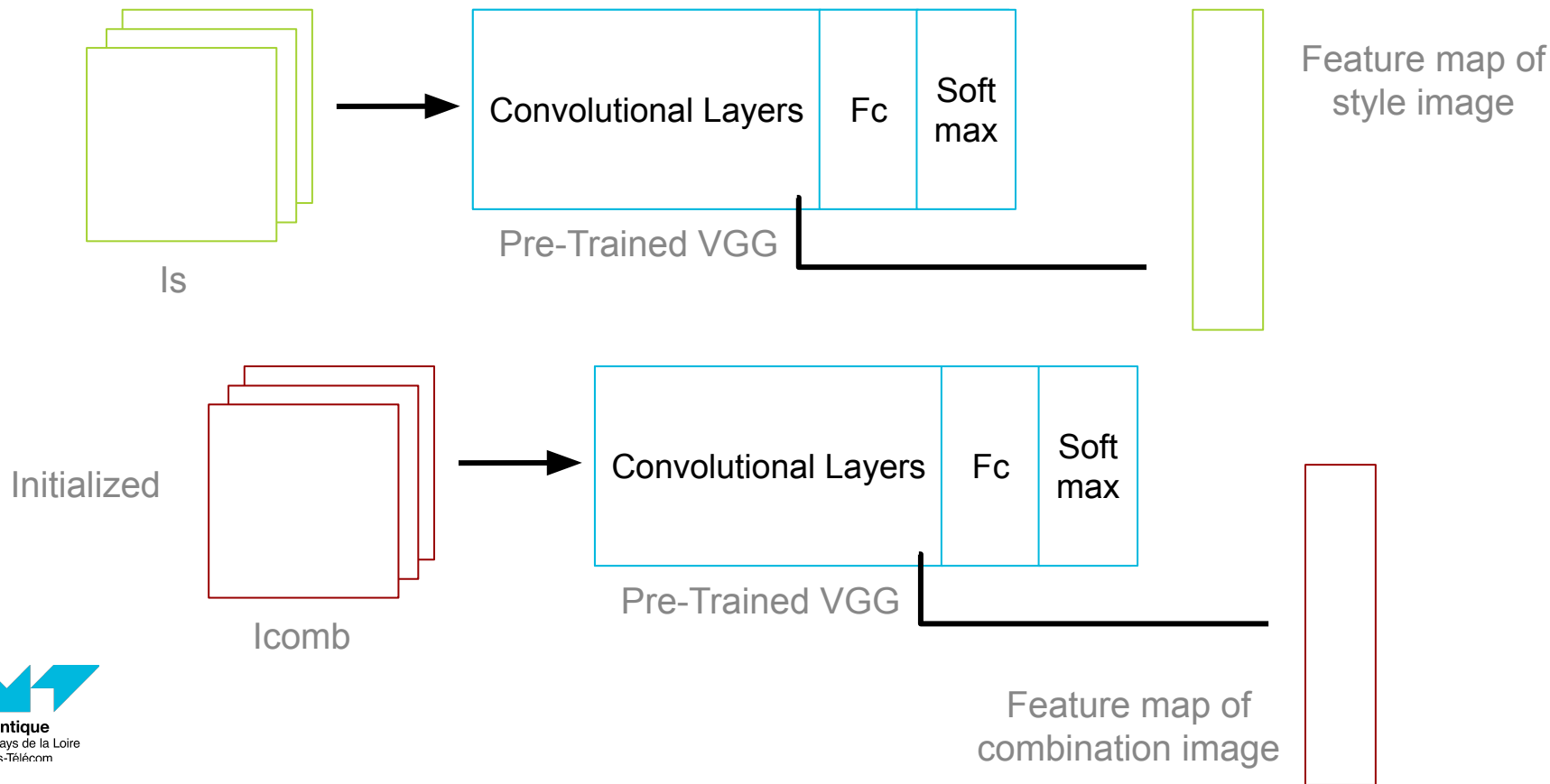
## 1.3 Algorithm overview



## 1.3 Algorithm overview



## 1.3 Algorithm overview



## 1.3 Algorithm overview

Features extraction:  
**feature maps** of  
**content, style,** and  
**combination** images

**Optimization:** We want to measure how far our combination image is from a “Good combination”

**Minimization of a loss function**  
 **$\text{lcomb\_Loss} = \text{Loss}(\text{lc}) + \text{Loss}(\text{ls})$**

Features extraction:  
**feature maps** of  
**content, style,** and  
**combination** images

Gradient descent

```
init lcomb  
for e in range(epochs):  
    calculate lcomb_loss  
    derivative of loss  
    update lcomb accordingly
```

# FUNDAMENTALS



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



- ▶ Let's break down the title of the paper:

### **Style Transfer by Relaxed Optimal Transport and Self-Similarity**

- ▶ Let's break down the title of the paper:

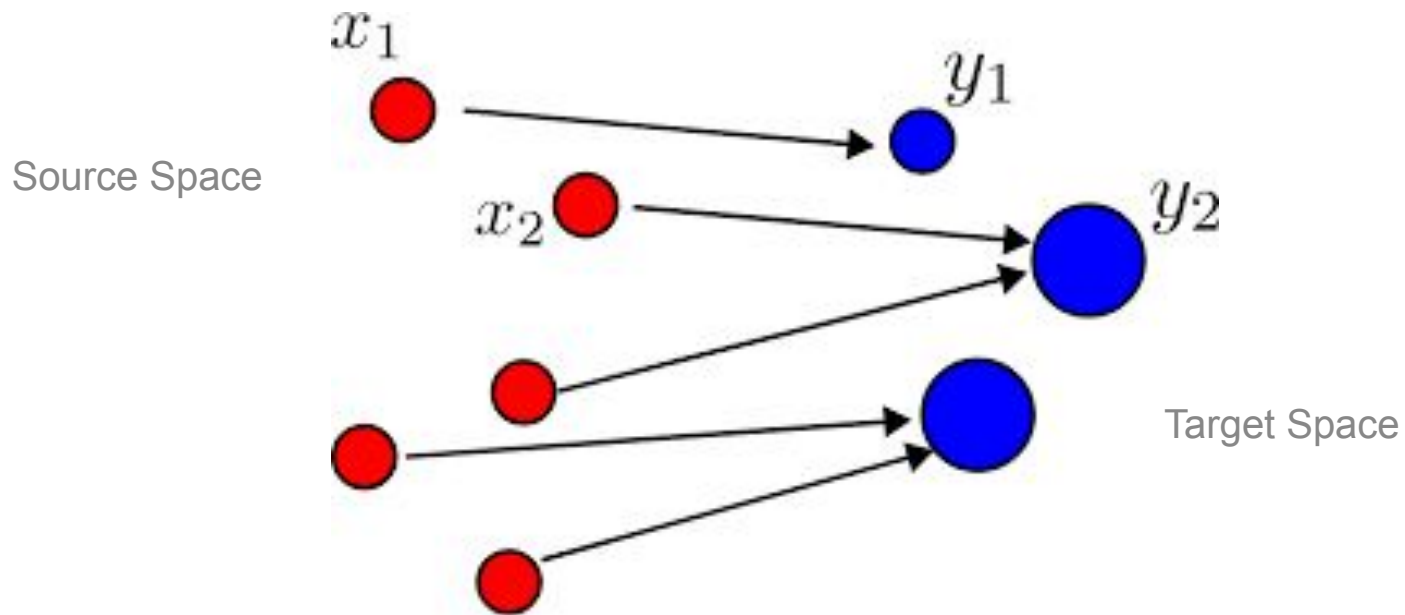
### Style Transfer by Relaxed Optimal Transport and Self-Similarity

### Optimal transport:

**Measures:** initial distribution of "mass" in a space and where we want to move it to

**Cost function:** cost of moving mass between points in the source and target spaces

**Transport map:** how to move the mass from the source distribution to the target distribution while minimizing the total cost



### Optimal transport:

Given two probability measures  $\mu$  (source) and  $\nu$  (target) on a Polish space (a complete, separable metric space)  $X$ , and a cost function  $c: X \times X \rightarrow \mathbb{R}$ , the optimal transport problem seeks a transport map  $T: X \rightarrow X$  that minimizes the total cost of transporting the mass:

don't know if  
this slide is  
necessary

$$\int_X c(x, T(x)) \, d\mu(x)$$

subject to the constraints that  $T$  pushes forward  $\mu$  to  $\nu$  (i.e., for any measurable set  $A$  in  $X$ ,  $\nu(A) = \mu(T^{-1}(A))$ )

## 2.1 Earth mover's distance

### Solution: Earth mover's distance EMD

$$\text{EMD}(A, B) = \min_{T \geq 0} \sum_{ij} T_{ij} C_{ij} \quad (2)$$

$$s.t. \sum_j T_{ij} = 1/m \quad (3)$$

$$\sum_i T_{ij} = 1/n \quad (4)$$

$$C_{ij} = D_{\cos}(A_i, B_j) = 1 - \frac{A_i \cdot B_j}{\|A_i\| \|B_j\|}$$

A and B: Sets of vectors

T: Transport matrix

C: Cost matrix

n: len(A)

m: len(B)

$$O(\max(m, n)^3)$$

Exact match between features

### Approximation: Relaxed earth mover's distance REMD

$$R_A(A, B) = \min_{T \geq 0} \sum_{ij} T_{ij} C_{ij} \quad s.t. \quad \sum_j T_{ij} = 1/m \quad (5)$$

$$R_B(A, B) = \min_{T \geq 0} \sum_{ij} T_{ij} C_{ij} \quad s.t. \quad \sum_i T_{ij} = 1/n \quad (6)$$

A and B: Sets of vectors

T: Transport matrix

C: Cost matrix

n: len(A)

m: len(B)

### Approximation: Relaxed earth mover's distance REMD

$$\ell_r = REMD(A, B) = \max(R_A(A, B), R_B(A, B)) \quad (7)$$

This is equivalent to:

$$\ell_r = \max \left( \frac{1}{n} \sum_i \min_j C_{ij}, \frac{1}{m} \sum_j \min_i C_{ij} \right) \quad (8)$$

A and B: Sets of vectors

T: Transport matrix

C: Cost matrix

n: len(A)

m: len(B)

→ More imperfections in feature matching, more natural looking style transfer



## Ideas to explain

- ▶ what is self similarity
- ▶ the mathematical basis and proofs used
- ▶ Link between theory and
- ▶ In id tortor sit amet augue.

#here the need for a self similarity descriptor to preserve the original content is described

### Trying to define Self-Similarity



### Trying to define Self-Similarity



**= Same “content”**  
**Same internal self similarity**

**≠ different colours**  
**different edges**  
**different texture**



another example of  
content preservation:  
Pareidolia

#trying to define what  
“content” is, importantly to  
say that the structure cannot  
be perceived logically but  
should be analysed within  
the context of an image

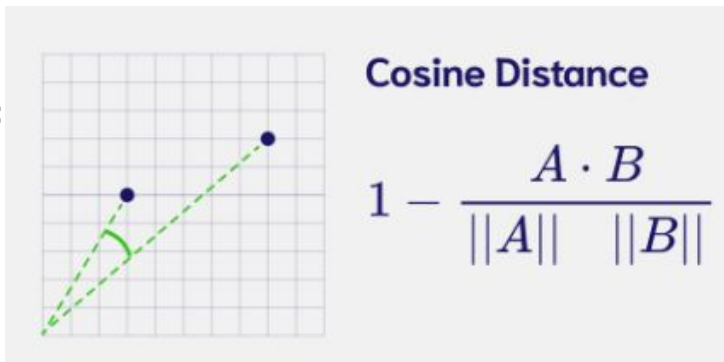
## 2.2 Self Similarity

**Solution:**

Self similarity descriptor=

**doesn't care about:**

- luminance
- texture
- scale
- Colour



## 2.2 Self Similarity



Self similarity Map (1- cosine\_distance)

## 2.2 Self Similarity

$$D_{ij} = \begin{matrix} & \begin{matrix} a & & b & c \end{matrix} \\ \begin{matrix} a \\ b \\ c \end{matrix} & \begin{bmatrix} 0 & \dots & 0.3 & 0.84 \\ \vdots & 0 & \vdots & \vdots \\ 0.3 & \dots & 0 & \vdots \\ 0.84 & \dots & \dots & 0 \end{bmatrix} \end{matrix}$$

**Self similarity Matrix (of cosine distances)**

# METHODOLOGY



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

How are optimal transport and style similarity concepts used in our work ?



Loss Function :

$$L(X, I_C, I_S) = \frac{\alpha \ell_C + \ell_m + \ell_r + \frac{1}{\alpha} \ell_p}{2 + \alpha + \frac{1}{\alpha}}$$

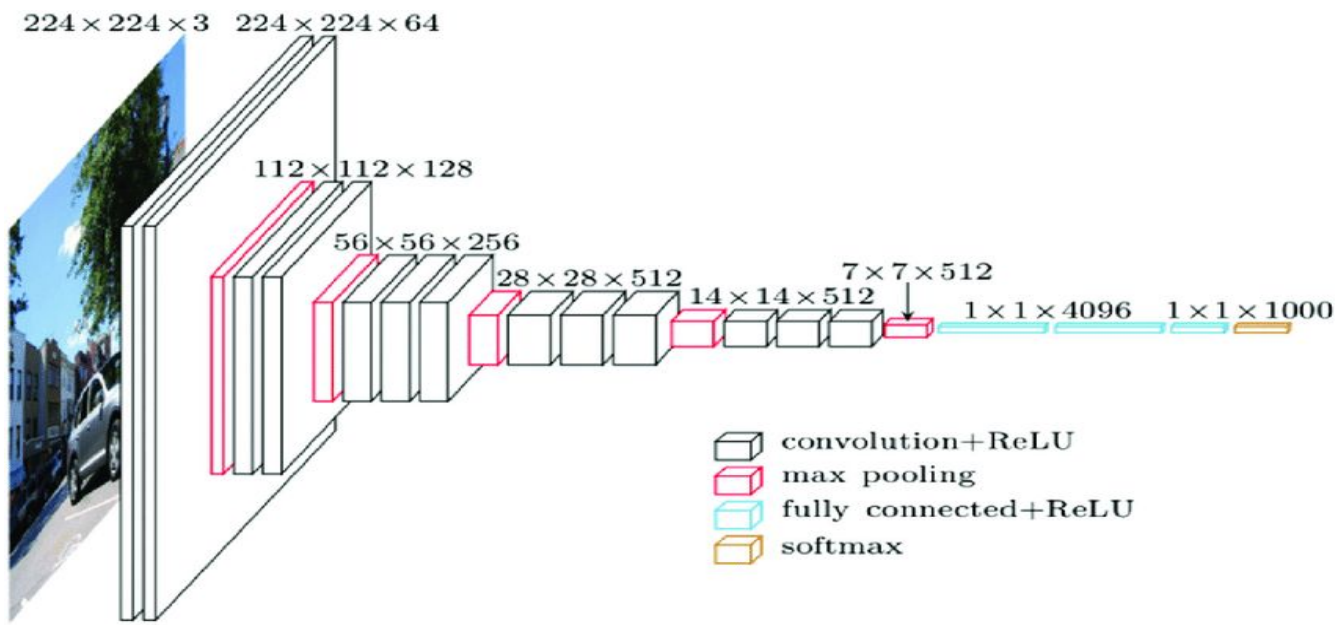


$$L(X, I_C, I_S) = \frac{\overbrace{\alpha l_C}^{\text{Content Loss}} + \overbrace{l_m + l_r + \frac{1}{\alpha} l_p}^{\text{Style Loss}}}{2 + \alpha + \frac{1}{\alpha}}$$

$\alpha$ : Hyperparameter: relative importance of content preservation to stylization

Both our style and content loss terms rely upon extracting a rich feature representation from an arbitrary spatial location.

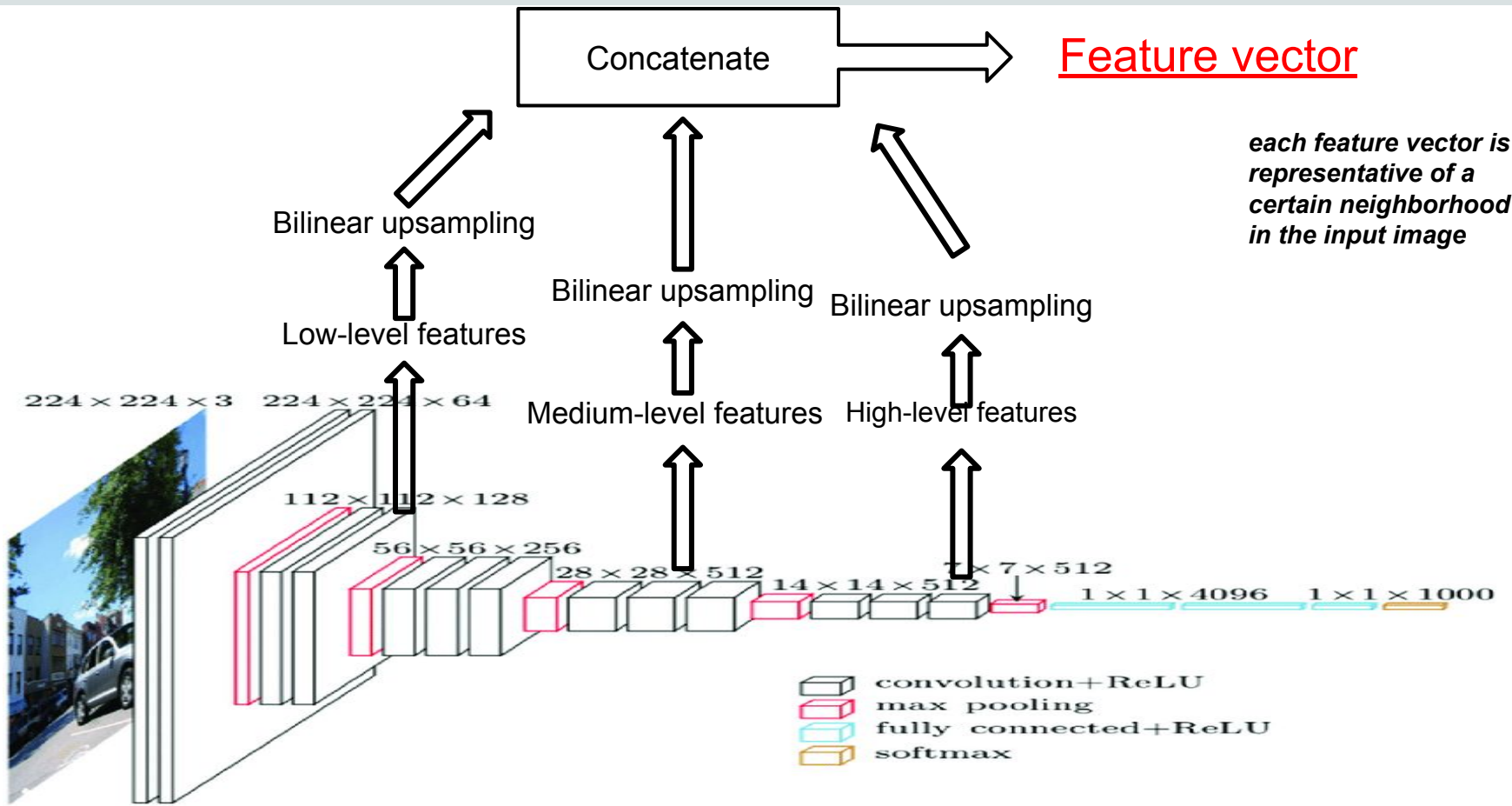
➡ VGG16 trained on ImageNet



# CHAPITRE 3 : Methodology

## VGG-16 for feature extraction

35

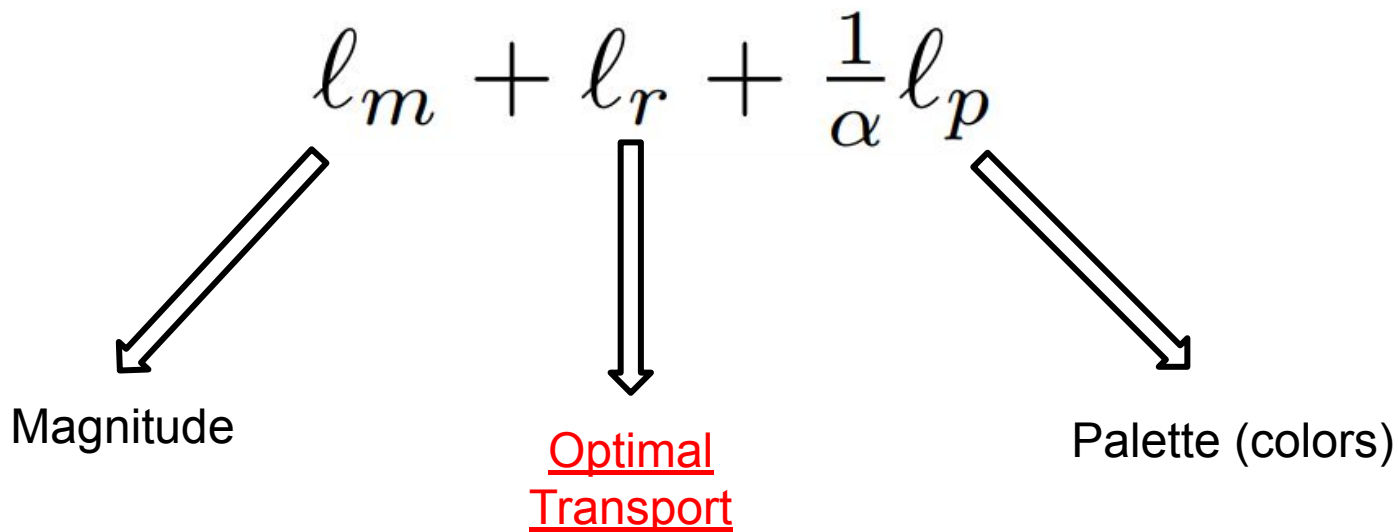


$$\mathcal{L}_{content}(X, C) = \frac{1}{n^2} \sum_{i,j} \left| \frac{D_{ij}^X}{\sum_i D_{ij}^X} - \frac{D_{ij}^{I_C}}{\sum_i D_{ij}^{I_C}} \right|$$

- $D^X$  : The pairwise cosine distance matrix of all feature vectors extracted from  $X(t)$ .
- $D^{I_C}$  : The pairwise cosine distance matrix of all feature vectors extracted from  $I_C$

Normalized cosine  
distance

Self-Similarity



### Optimal Transport Problem

EMD ?  $\implies$  Costly:  $O(\max(m,n)^3)$   $\implies$  R-EMD

$$\ell_r = \max \left( \frac{1}{n} \sum_i \min_j C_{ij}, \frac{1}{m} \sum_j \min_i C_{ij} \right)$$

C: Cost Matrix: How far an element of A (set of feature vectors extracted from X) is from an element of B (set of feature vectors extracted from Is)

$\ell_r$  : Good transfer of the structural forms  
of the source image to the target. BUT..

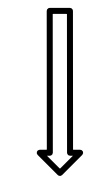
$\ell_p$  : Color matching loss (palette)  $\longrightarrow$  R-EMD between X  
and Is

Magnitude of the  
feature vectors ignored  
by cosine distance  $\longrightarrow$  Visual Artifacts  
(under/over saturation)

$$\ell_m = \frac{1}{d} \|\mu_A - \mu_B\|_1 + \frac{1}{d^2} \|\Sigma_A - \Sigma_B\|_1$$

- $\mu_A$  ( $\mu_B$ ) and  $\Sigma_A$  ( $\Sigma_B$ ) are the mean and covariance of the feature vectors  
in set A (in set B).

Palette shifting is at odds with  
content preservation



$$\frac{1}{\alpha} \ell_p$$

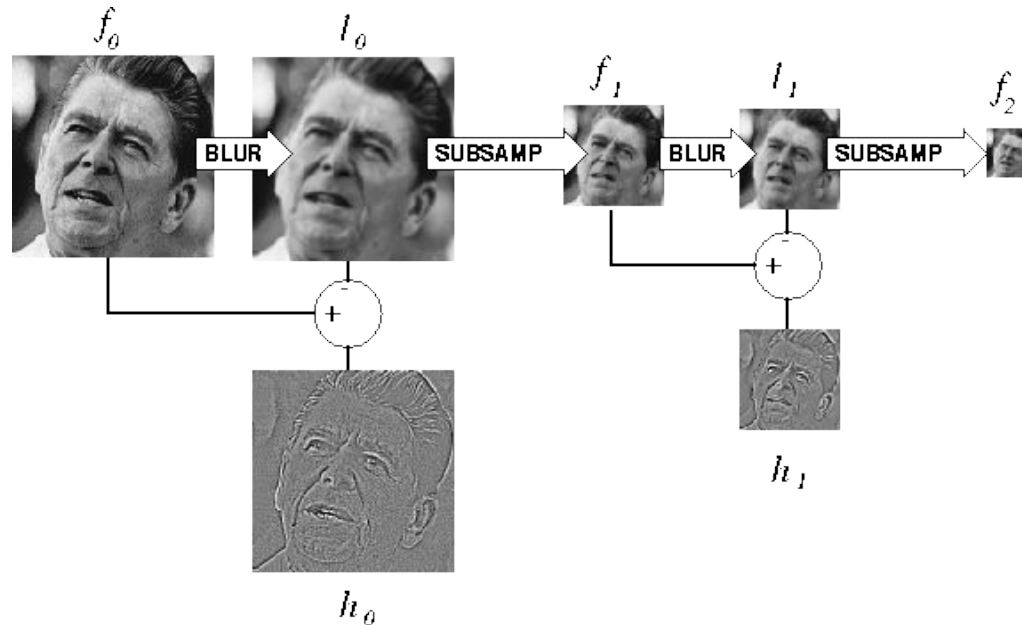
# IMPLEMENTATION AND REGION CONTROL



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



Laplacian Pyramid: Image representation consisting of band-pass images and low-frequency residual image:



# IMPLEMENTATION AND USER CONTROL

42

Step by Step

Goal:



Content Image



Style Image



Stylized Image

### STEP 0: INITIALIZATION



### STEP 1: LAPLACIAN PYRAMID DE COMPOSITION



### STEP 2: LOSS CALCULATION

 $L$ 

Content Image  
(downscaled)

,



Style Image  
(downscaled)

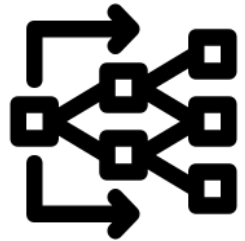
,



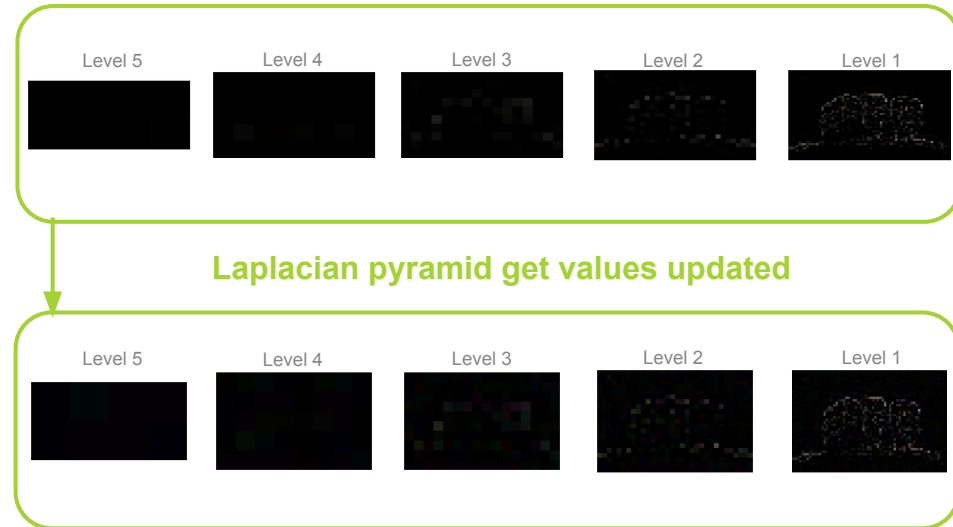
Stylized Image

### STEP 2: BACKPROPAGATE AND UPDATE PYRAMID VALUES

$$L \left( \begin{matrix} \dots \\ \dots \\ \dots \end{matrix} \right)$$



Backpropagation



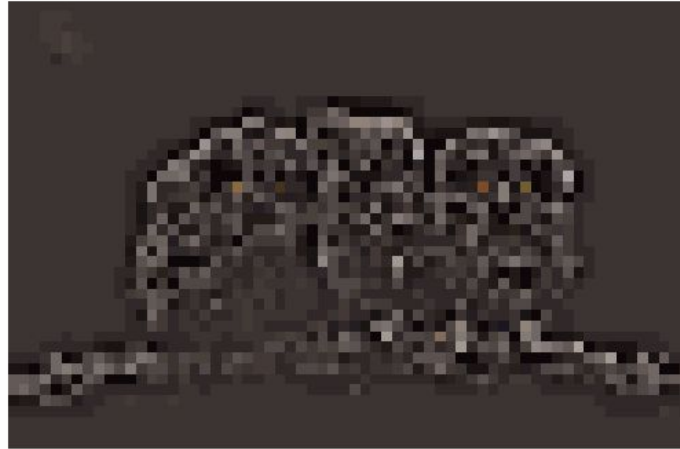
### STEP 3: RECONSTRUCT IMAGE FROM UPDATED PYRAMID

Laplacian pyramid with updated values



Repeat step 1 to 3 (250 iterations) using the updated stylized image

Scale: 64 px, Iteration: 0



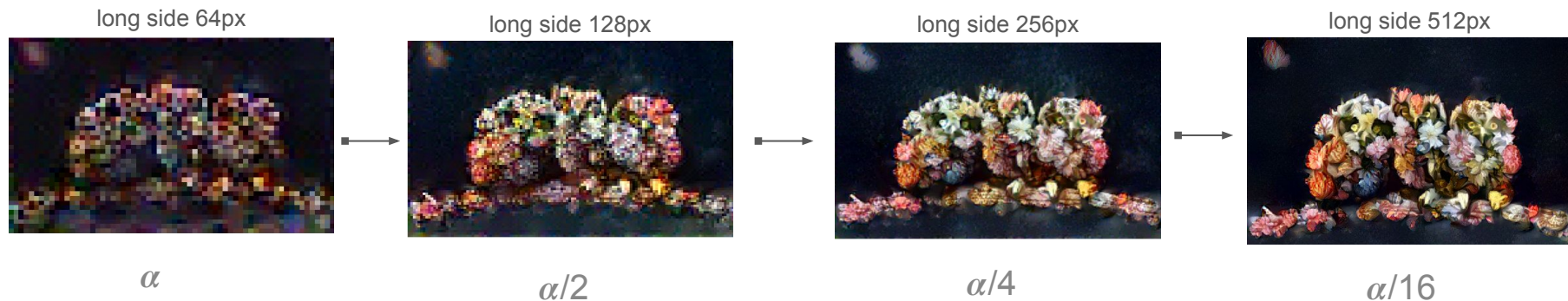


# IMPLEMENTATION AND USER CONTROL

49

## Multiple scale calculation

The iterative process described previously is made for 4 different scales using the (upscaled) output of the previous scale as input, halving the content weight ( $\alpha$ ) for the next scale:



## Multiple scale calculation

The iterative process described previously is made for 4 different scales using the (upscaled) output of the previous scale as input, halving the content weight ( $\alpha$ ) for the next scale

Scale: 64 px, Iteration: 0



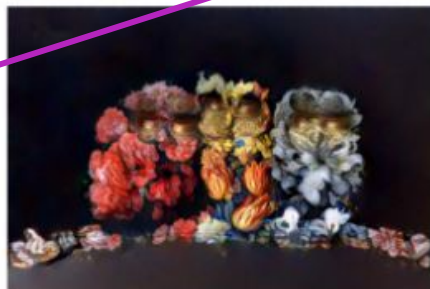
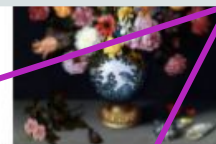
## User control over style

Mask specific areas to have the same style:



## User control over style

Mask specific areas to have the same style:



Control is enforced by making the pairs of points in the same region have higher weight in the loss calculation

# EXPERIMENTS AND RELATED WORK

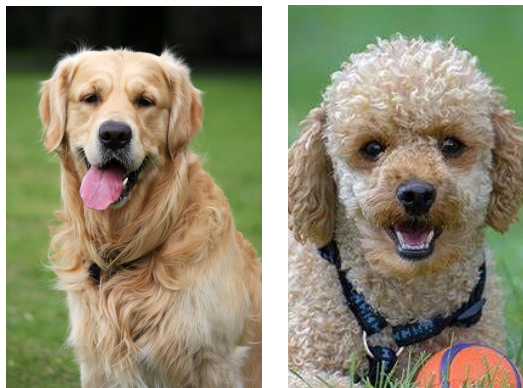


**IMT Atlantique**  
Bretagne-Pays de la Loire  
École Mines-Télécom

## 4.1 Large-Scale human evaluation

Regimes :

**Paired**



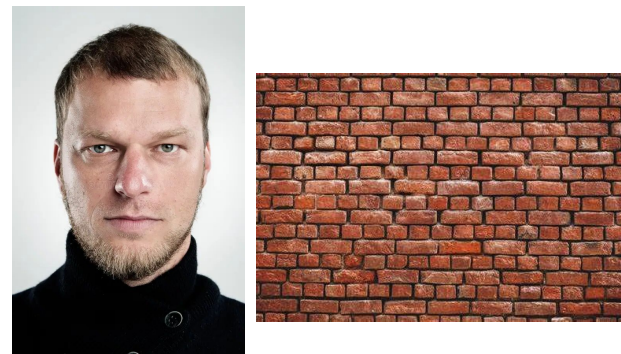
Content image = style image

**Unpaired**



Content image  $\neq$  style image

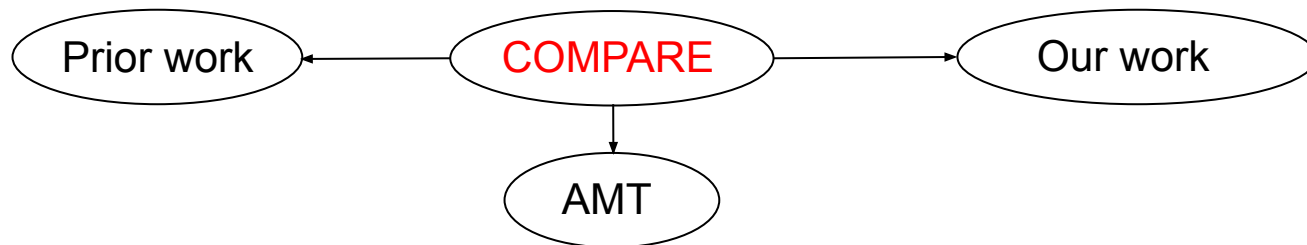
**Texture**



Content: Face photography  
Style: Homogeneous texture

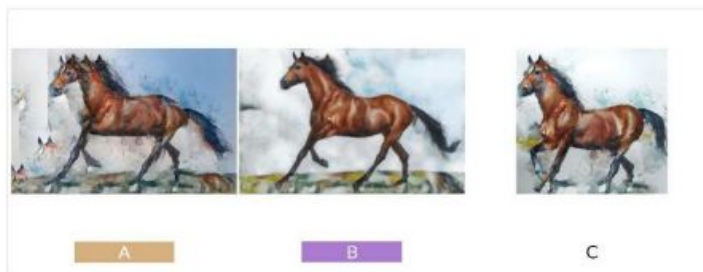
\*30 style/content pairings (total of 90)

## 4.1 Large-Scale human evaluation



An example of worker's interfaces:

Evaluate whether **Image A** or **Image B** has more similar **style** to Image C



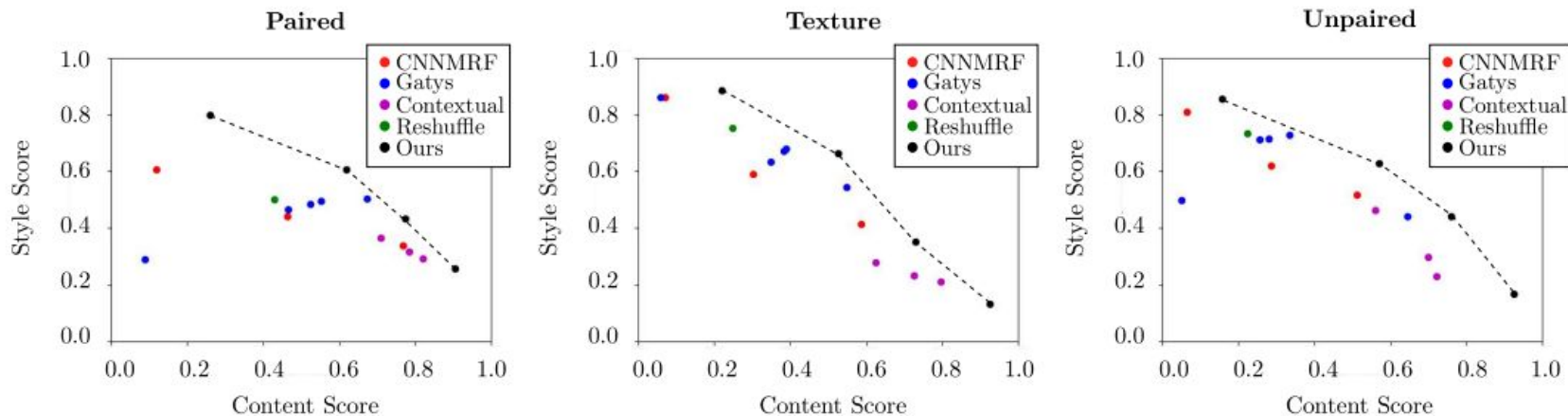
The style of Image C is most similar to...

- a) **Image A**
- b) **Image B**
- c) **Equally** A and B
- d) **Neither** A or B

## 4.1 Large-Scale human evaluation

We test 3 sets of hyper-parameters: Content weight (high and low stylization)

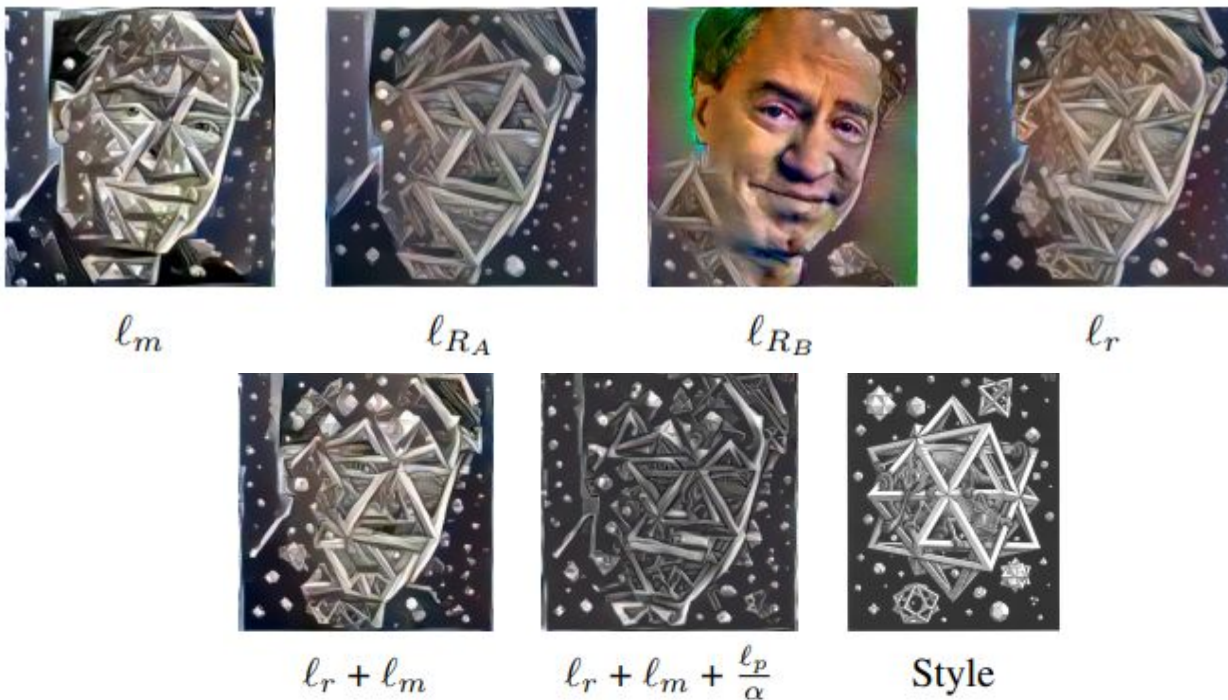
The style score is defined analogously :





## 4.2. Ablation study

Effect of different terms of our style loss



## 4.3. Relaxed EMD Approximation Quality

$$\frac{REMD(A,B)}{EMD(A,B)} \leq 1$$

## 4.4. Timing results

Image size	64	128	256	512	1024
Ours	20	38	60	95	154
Gatys	8	10	14	33	116
CNNMRF	3	8	27	117	X
Contextual	13	40	189	277	X
Reshuffle	-	-	-	69*	-

# CONCLUSIONS AND FUTURE WORK



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

## Conclusions

- Our algorithm demonstrates superior performance in style transfer compared to prior methods.
- We emphasize the importance of style-similarity losses for enhancing stylization quality.
- The simplicity and effectiveness of our earth movers distance approximation highlight its potential in style transfer applications.

## Future work

- Further exploration of more accurate approximations for the earth movers distance.
- Improvement of algorithm speed through the incorporation of feed-forward style transfer methods using our proposed objective function.

# QUESTIONS?



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

**Alex SZAPIRO**

How does this method differ from  
Multimodal Style Transfer via Graph  
Cuts (Zhang et al., 2019) ?



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

**Pierre Monot**

If the user control is also using a  
REMD loss, why not incorporate it in  
the content loss?



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

**Yosr DRIRA**

The paper has mentioned cosine distance as a self similarity descriptor. Are there other descriptors that you could have used ?



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