

# Style transfer models

Victor Kitov

[v.v.kitov@yandex.ru](mailto:v.v.kitov@yandex.ru)



# Table of Contents

- 1 Neural style transfer introduction
- 2 Offline style transfer
- 3 Online ST (Johnson)
- 4 Online ST (Ulyanov)

# Neural style transfer

- Input: content image, style image.
- Style transfer - application of artistic style from style image to content image.
- Easy to do with convolutional neural networks.

# Example



# Applications

- **Enhance social communication**
  - adding personality, emotions
- **Extended drawing tools**
  - 2D drawings for painters, designers, architects.
- **Apply special effects** to movies, computer games (interactive style, depending on situation)
- **Cheap cartoon creation** from filmed scenes with actors.
- **Transfer learning**, e.g:
  - photos->paintings, learn to identify people on paintings
  - day photos->night photos, learn image segmentation at night

# Table of Contents

1 Neural style transfer introduction

2 Offline style transfer

3 Online ST (Johnson)

4 Online ST (Ulyanov)

## Seminal work<sup>1</sup>

- Consider convolutional network (VGG).
- Content is reconstructed from activations at layer  $l_c$ .
- Style is reconstructed from correlations of its activations at layers  $l \in L_s$

---

<sup>1</sup>Gatys et al (2015). A Neural Algorithm of Artistic Style.

# Content image



# Style application



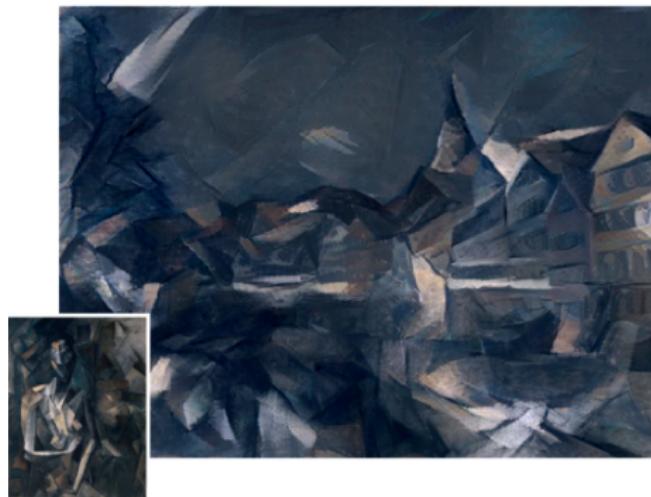
# Style application



# Style application



## Style application



# Notation

- Images:
  - $x_c$  - content image;
  - $x_s$  - style image
  - $x$  - stylized image (to be found)
- Consider convolutional layer  $l$ :
  - feature map has  $M_l = H_l \times W_l$  spatial dimensionality
  - $N_l$  filters
- $\Phi_j^l(x) \in \mathbb{R}^{M_l}$  - spatially vectorized output of  $i$ -th filter on  $j$ -th position on layer  $l$ .
  - individual element is  $\Phi_{ij}^l(x)$ ,  $i = \overline{1, N_l}$ ,  $j = \overline{1, M_l}$ .
  - all feature maps:  $\Phi^l(x) \in \mathbb{R}^{N_l \times M_l}$
- $\|A\|_F^2 = \sum_{i,j} a_{i,j}^2$  - squared Frobenius matrix norm.

## Gram matrix

- Define Gram matrices  $\in \mathbb{R}^{N_I \times N_I}$ :

$$G_{ij}(x) = \frac{1}{M_I} \sum_{k=1}^{M_I} \Phi_{ik}(x) \Phi_{jk}(x) = \frac{1}{M_I} \Phi^I(x) (\Phi^I(x))^T$$

- They capture uncentered correlations<sup>2</sup> between channels (feature maps)
  - spatial information is aggregated

---

<sup>2</sup>Centering gives same results but is more computationally expensive.

# Losses

- Content loss (high level features should be similar):

$$\mathcal{L}_{content}(x, x_c) = \frac{1}{N_{I_c} M_{I_c}} \left\| \Phi^{I_c}(x) - \Phi^{I_c}(x_c) \right\|_F^2$$

- Style loss for one layer (feature correlations should be similar):

$$\mathcal{L}_{style}^I = \frac{1}{N_I^2} \left\| G^I(x) - G^I(x_s) \right\|_F^2$$

- Total style loss:

$$\mathcal{L}_{style}(x, x_s) = \sum_{I \in L_s} w_I \mathcal{L}_{style}^I$$

- $I_c$ -content layer,  $L_s$ -layers for style reconstruction.

# Losses

- Higher  $l \Rightarrow$  higher content/style abstraction and dimensionality preserved.
- Content: only 1 deep layer  $l_c$  is used (versatile reconstruction of content object)
- Style: both shallow (low level pixel info) and deep layers (general style) are used
- $x$  is found from

$$x = \arg \min_x \{ \mathcal{L}_{content}(x, x_c) + \alpha \mathcal{L}_{style}(x, x_s) \} \quad (1)$$

- Hyperparameter  $\alpha > 0$  measures relative importance of style, compared to content.

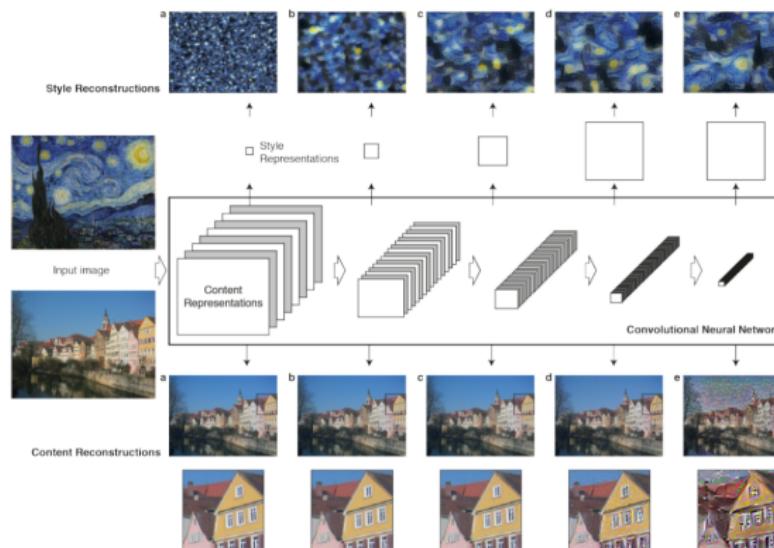
# Detailed style transfer algorithm

INPUT: content image  $x_c$ , style image  $x_s$ .

- ➊ Train classification CNN or take pretrained (e.g. VGG)
- ➋ Using CNN, compute
  - ➌ content image features  $\Phi^{l_c}(x_c)$
  - ➍ compute Gram matrices  $\{G^l(x_s)\}_{l \in L_s}$
- ➎ Initialize new image  $I$  (using content image or randomly)
- ➏ Repeat until convergence:
  - ➐ Pass  $I$  through CNN
  - ➑ Compute loss (1)
  - ➒ Backprop w.r.t  $I$ , holding CNN fixed
  - ➓ Update  $I$
- ➔ Output  $I$  as stylization result

# Image reconstruction from content/style information

Content (lower row) and style reconstruction (upper row) from layers of increasing depth:



Deeper layers contain more abstract information.

# Increasing content importance

$\downarrow \alpha$  ( $\Leftrightarrow \uparrow$ importance of content preservation):



1



2



3



4

## Total variation loss

- In ST we reverse max-pooling=>solution underdetermined.
- Adding total variation loss makes it determined:

$$\mathcal{L}_{tv}(x) = \sum_{i,j} (x_{i+1,j} - x_{i,j})^2 + (x_{i,j+1} - x_{i,j})^2$$

- Total variation loss act as smoothness prior.
- Alternative less common definitions:

$$\sum_{i,j} |x_{i+1,j} - x_{i,j}| + |x_{i,j+1} - x_{i,j}|$$

$$\sum_{i,j} \left( (x_{i+1,j} - x_{i,j})^2 + (x_{i,j+1} - x_{i,j})^2 \right)^\gamma, \quad \gamma > 0$$

# Table of Contents

1 Neural style transfer introduction

2 Offline style transfer

3 Online ST (Johnson)

4 Online ST (Ulyanov)

## Fast style transfer<sup>3</sup>

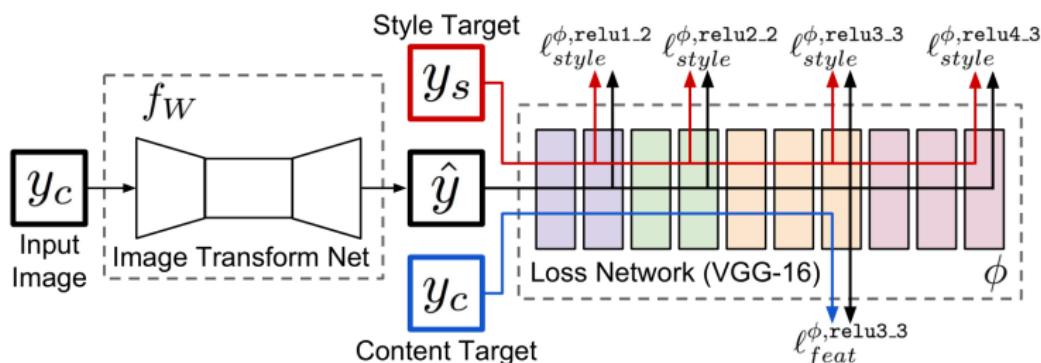
- 2 nets:
  - image transform net (generator)
  - loss net (VGG-16)
- Loss net is held constant
  - it uses pretrained weights from image classification

---

<sup>3</sup>Link to [paper](#) and [technical details](#).

# Image transform net

- Image transform net is trained on
  - single style image
  - different content images
  - loss: offline ST loss (1)+total variation loss



# Application

- To stylize a new image just pass it through the transformer net!
  - no optimization required
  - but need separate transformer net for each style
- Speedup of online ST compared to offline ST:

Image Size	Gatys <i>et al</i>			Ours	Speedup		
	100	300	500		100	300	500
256 × 256	3.17	9.52s	15.86s	<b>0.015s</b>	212x	636x	<b>1060x</b>
512 × 512	10.97	32.91s	54.85s	<b>0.05s</b>	205x	615x	<b>1026x</b>
1024 × 1024	42.89	128.66s	214.44s	<b>0.21s</b>	208x	625x	<b>1042x</b>

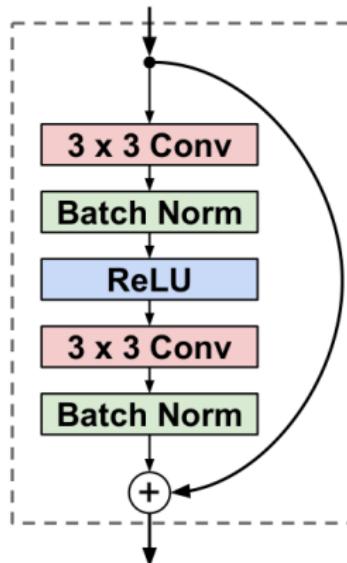
# Architecture of generator

- No pooling layers used
- Downsampling: strided convolutions
- Residual blocks used in the middle
  - good at reproducing identity
  - reasonable: identity gives content image!
- Each convolution-followed by batch normalization
  - later articles: instance normalization better and probably better without normalization at all.
- Upsampling: fractionally strided convolutions
  - in later articles replaced by upsampling with NN & convolution - better

# Architecture of transformet net

<b>Layer</b>	<b>Activation size</b>
Input	$3 \times 256 \times 256$
Reflection Padding ( $40 \times 40$ )	$3 \times 336 \times 336$
$32 \times 9 \times 9$ conv, stride 1	$32 \times 336 \times 336$
$64 \times 3 \times 3$ conv, stride 2	$64 \times 168 \times 168$
$128 \times 3 \times 3$ conv, stride 2	$128 \times 84 \times 84$
Residual block, 128 filters	$128 \times 80 \times 80$
Residual block, 128 filters	$128 \times 76 \times 76$
Residual block, 128 filters	$128 \times 72 \times 72$
Residual block, 128 filters	$128 \times 68 \times 68$
Residual block, 128 filters	$128 \times 64 \times 64$
$64 \times 3 \times 3$ conv, stride 1/2	$64 \times 128 \times 128$
$32 \times 3 \times 3$ conv, stride 1/2	$32 \times 256 \times 256$
$3 \times 9 \times 9$ conv, stride 1	$3 \times 256 \times 256$

# Residual block



- No padding is used (causes artifacts at image border)
- To match size, identity is cropped at the end

# Table of Contents

- 1 Neural style transfer introduction
- 2 Offline style transfer
- 3 Online ST (Johnson)
- 4 Online ST (Ulyanov)

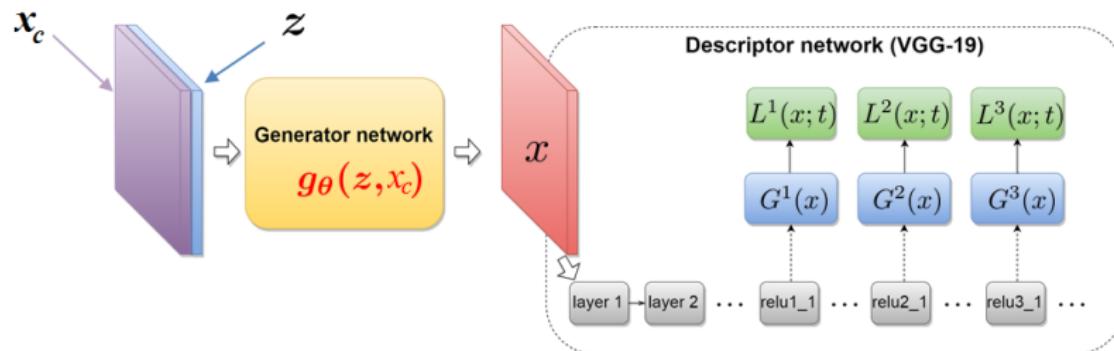
Idea<sup>4</sup>

- Train generative network  $g_\theta(x_c, z)$ 
  - $x_c$ : content image
  - $z$ : uniform random noise
  - $\theta$ : trained parameters (weights)
- By passing different random noise  $z$  we hope to generate many versions of stylized  $x_c$ .
- Use standard offline ST loss (1).

---

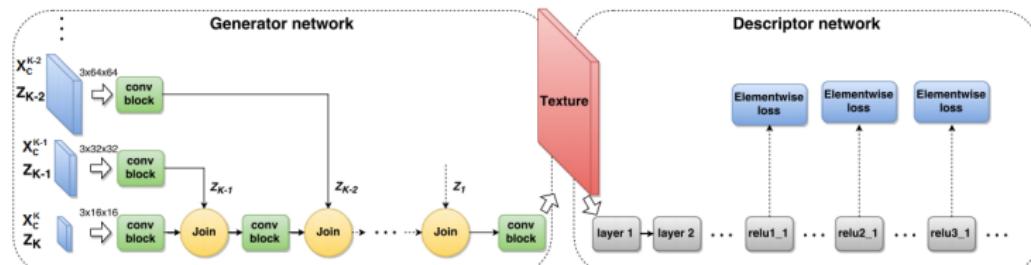
<sup>4</sup>Ulyanov et. al. (2016). Texture Networks: Feed-forward Synthesis of Textures and Stylized Images.

# Proposed architecture



# Architecture

- $x_c^k$  - downsampled content image  $x_c$  by  $2^{k-1}$ ,  $k = \overline{1, K}$ .
- $Z_1, Z_2, Z_3, \dots$  - random noise at different resolutions (abstract levels).
- Only generator network is trained, descriptor network held fixed.
- In descriptor network (first layers of VGG) style transfer loss (1) is calculated.



# Architecture in detail

