

# Transformation-Based Style Transfer

Victor Kitov

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

# Table of Contents

- 1 Transformation-based ST (Johnson)
- 2 Transformation-based ST (Ulyanov)
- 3 Instance normalization

# Fast style transfer<sup>1</sup>

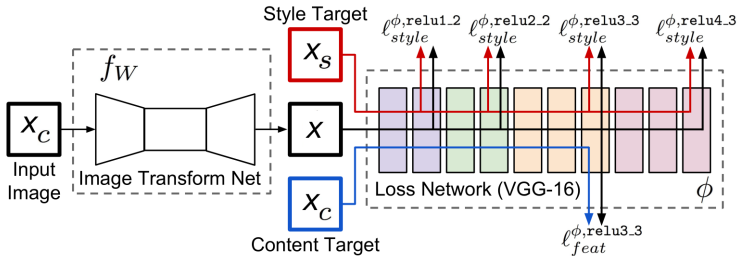
- 2 nets:
  - image transform net (generator)
  - loss net (VGG-16)
- Loss net is held constant
  - it uses pretrained weights from image classification

---

<sup>1</sup>Link to [paper](#) and [technical details](#).

# Train

- Image transform net is trained on
  - single style image
  - different content images
  - loss=offline ST loss



# Test

- To stylize a new image just pass it through the transformer!
  - real-time, 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
$256 \times 256$	3.17	9.52s	15.86s	<b>0.015s</b>	212x	636x	<b>1060x</b>
$512 \times 512$	10.97	32.91s	54.85s	<b>0.05s</b>	205x	615x	<b>1026x</b>
$1024 \times 1024$	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 <sup>2</sup>.
- Upsampling: fractionally strided convolutions
  - later articles: upscaling & convolution - better, no checkerboard artifacts

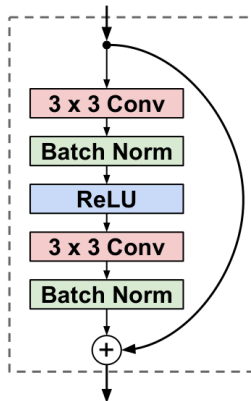
---

<sup>2</sup>May omitting normalization at all be even better?

# Architecture of transformet net

Layer	Activation size
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 Transformation-based ST (Johnson)
- 2 Transformation-based ST (Ulyanov)**
- 3 Instance normalization

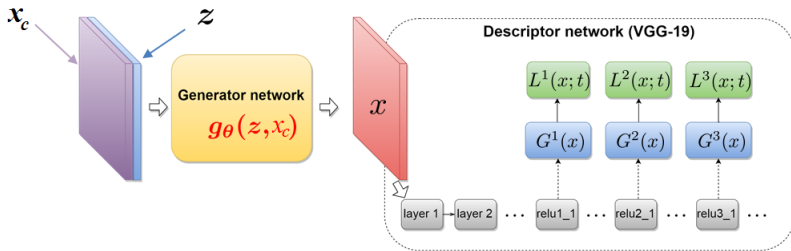
## Idea<sup>3</sup>

- 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 ST loss (Gatys et al.)

---

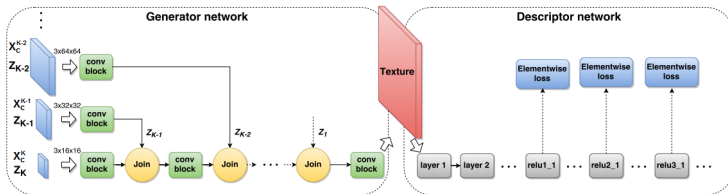
<sup>3</sup>Ulyanov et. al. (2016). Texture Networks: Feed-forward Synthesis of Textures and Stylized Images.

# Proposed architecture

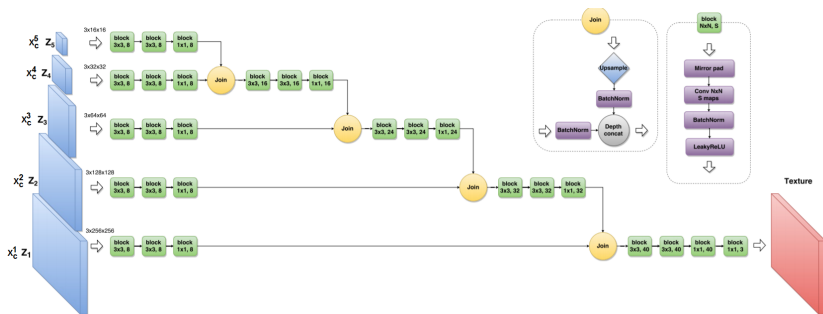


# Architecture

- $x_c^k$  - 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 (??) is calculated.



## Architecture in detail



# Table of Contents

- 1 Transformation-based ST (Johnson)
- 2 Transformation-based ST (Ulyanov)
- 3 Instance normalization

# Instance normalization<sup>4</sup>

- Instance normalization works better than batch normalization in online ST generators.
- Batch normalization: normalize spatial feature within minibatch of objects
- Instance normalization: normalize spatial feature within single object
- Normalizes distribution of features of the content image.
  - easier to map this distribution to given style

---

<sup>4</sup>2017 - Instance normalization - The Missing Ingredient for Fast Stylization  
- Ulyanov.

# Formula

- Let  $a_{icxy}$  be neuron activation for object  $n$  in the minibatch, channel  $c$  and coordinates  $x, y$ .
- Batch normalization ( $B$ -size of minibatch):

$$\tilde{a}_{ncxy} = \gamma \frac{a_{ncxy} - \mu_c}{\sqrt{\sigma_c^2 + \varepsilon}} + \beta$$

$$\mu_c = \frac{1}{BWH} \sum_{n=1}^B \sum_{x=1}^W \sum_{y=1}^H a_{ncxy}, \quad \sigma_c^2 = \frac{1}{BWH} \sum_{n=1}^B \sum_{x=1}^W \sum_{y=1}^H (a_{ncxy} - \mu_c)^2$$

- Instance normalization (proposed approach)

$$\tilde{a}_{ncxy} = \gamma \frac{a_{ncxy} - \mu_{nc}}{\sqrt{\sigma_{nc}^2 + \varepsilon}} + \beta$$

$$\mu_{nc} = \frac{1}{WH} \sum_{x=1}^W \sum_{y=1}^H a_{ncxy}, \quad \sigma_{nc}^2 = \frac{1}{WH} \sum_{x=1}^W \sum_{y=1}^H (a_{ncxy} - \mu_{nc})^2$$



# Comments

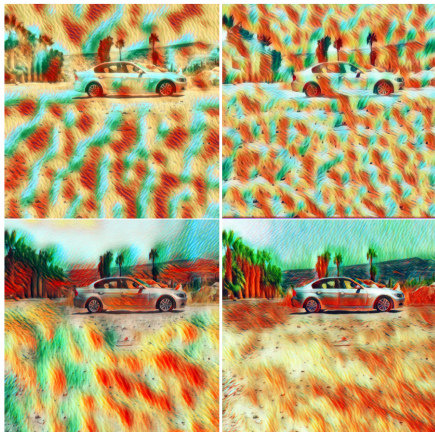
- $\gamma, \beta$  - learned parameters.

At test time:

- batch normalization: substitute train sample mean & std.deviation instead of  $\mu_c, \sigma_c^2$ .
- instance normalization: recalculate  $\mu_{nc}, \sigma_{nc}^2$  for test object.

## Results with instance normalization are better

Results for generator of Ulyanov (2016) (left) and Johnson (2016) (right). Batch normalization (upper row) and instance normalization (lower row).



Results are robust to contrast/brightness of content



(a) Content image.



(b) Stylized image.

