



(Talk)

Photorealistic Image Stylization:

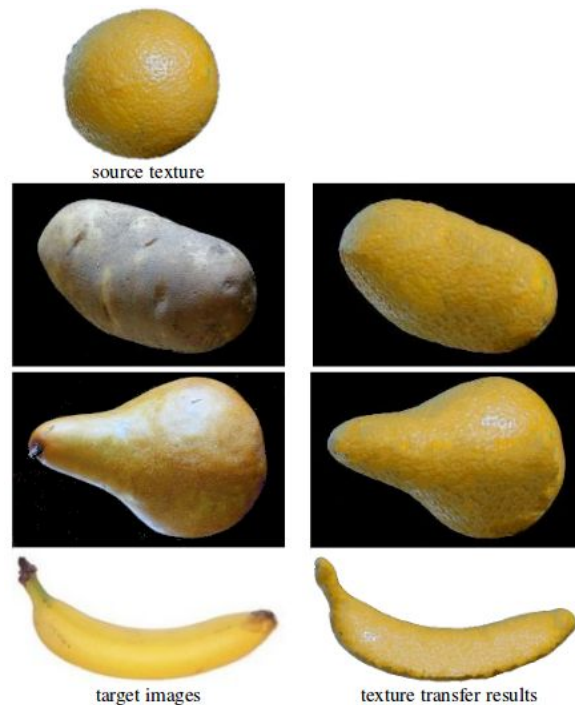
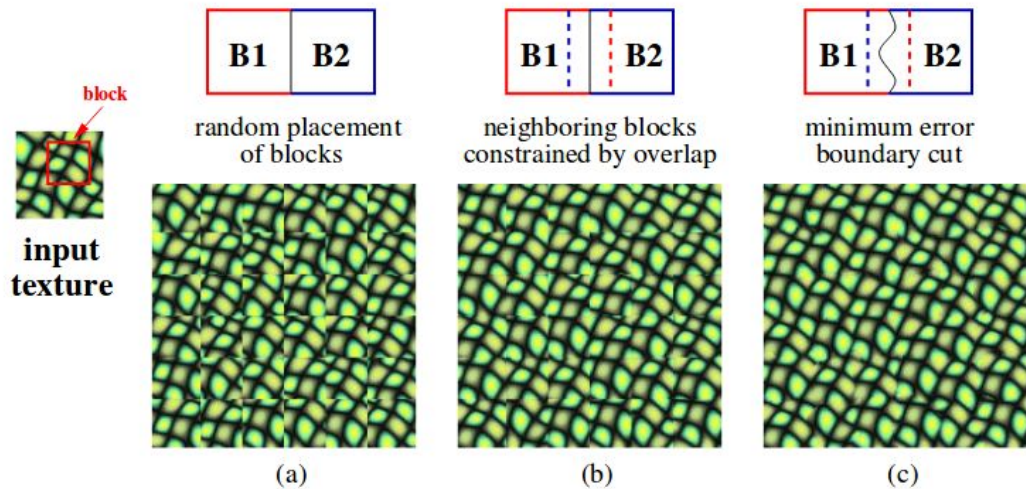
Image Stylization using Convolutional Neural
Networks

Leon A. Gatys, Alexander S. Ecker, Matthias Bethge
CVPR 2016

Presented by: Abhishek Jha

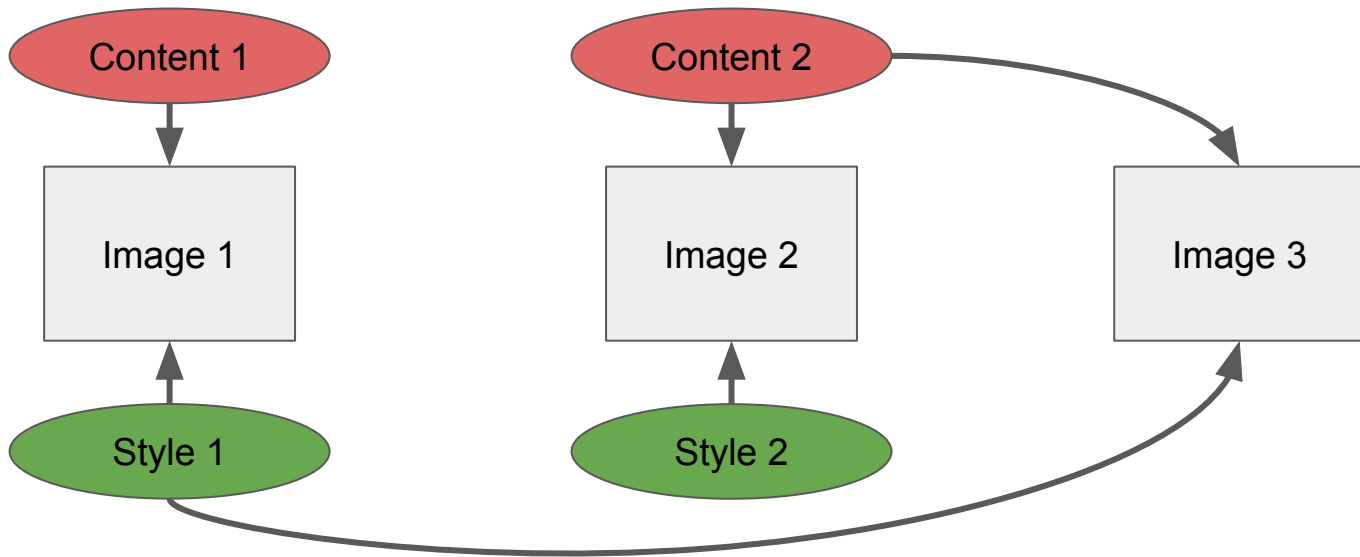
Related Work

Efros and Freeman, 2001.



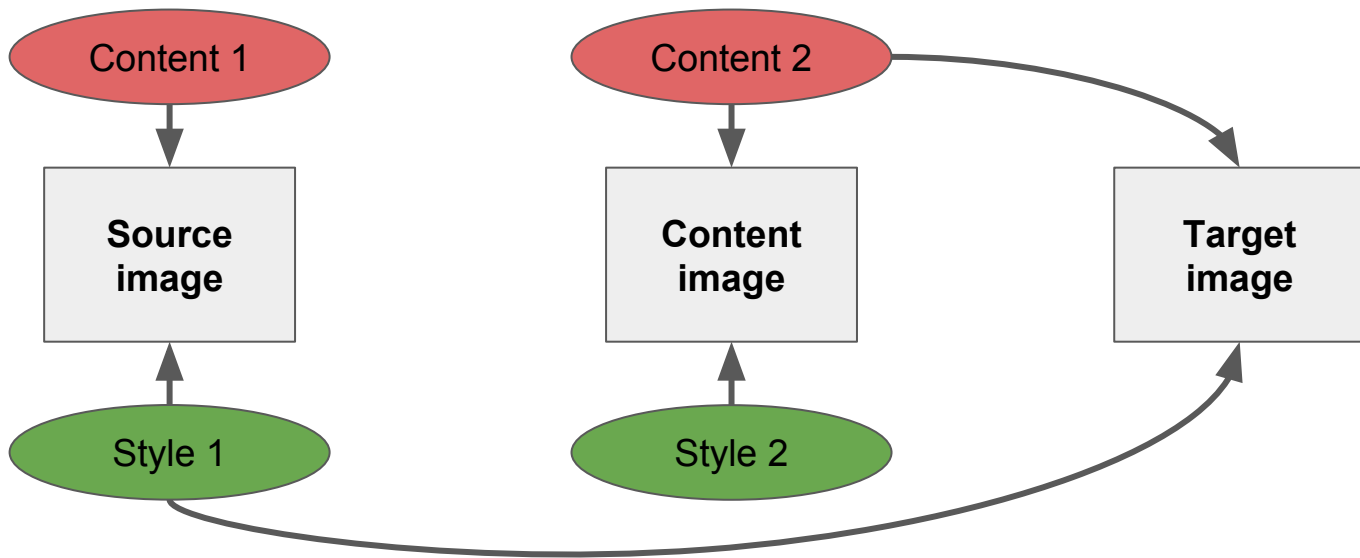
Motivation

- Any image can be factorize into content and style.
- Content provides the semantics of the entities present in the image.
- Style provide the textures corresponding to those entities.

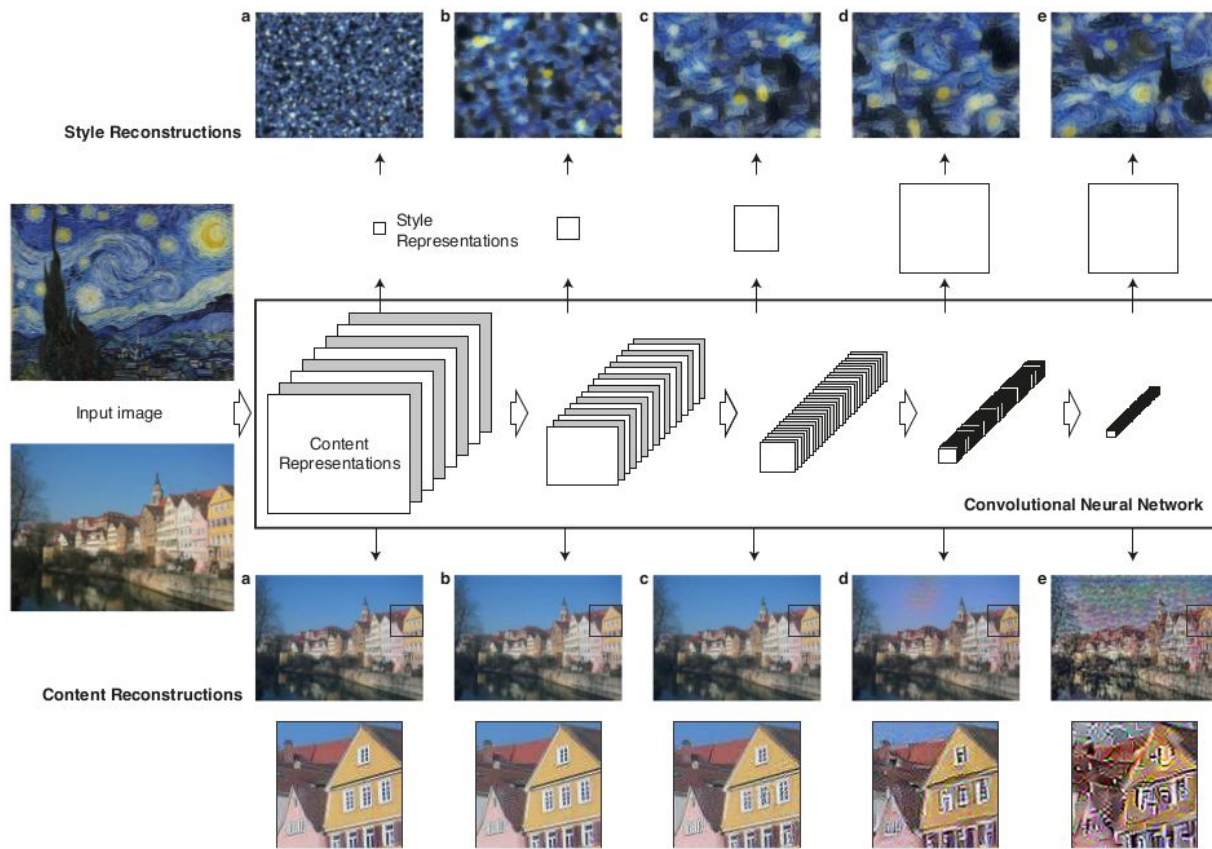


Motivation

- Any image can be factorize into content and style.
- Content provides the semantics of the entities present in the image.
- Style provide the textures corresponding to those entities.



Feature selection



Depth from which the features are taken

Pixel level detail of the reconstructed image

Content Representation

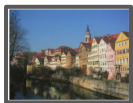
For any image x ; the filter response in layer ℓ : $F^\ell \in \mathcal{R}^{N_\ell \times M_\ell}$

Where N_ℓ is number of filters in layer ℓ and $M_\ell = \text{HxW}$ of the feature map.

Content Loss:

If F_{ij}^ℓ is the activation of i -th filter in the position j in the layer ℓ

Let $p \rightarrow$ original image



Let $x \rightarrow$ generated image



S.t. P^ℓ and F^ℓ are their respective feature response in layer ℓ .

$$\mathcal{L}_{\text{content}}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2. \quad (1)$$

$$\frac{\partial \mathcal{L}_{\text{content}}}{\partial F_{ij}^l} = \begin{cases} (F^l - P^l)_{ij} & \text{if } F_{ij}^l > 0 \\ 0 & \text{if } F_{ij}^l < 0, \end{cases} \quad (2)$$

Style Representation

For any image x ; the Gram matrix for layer l : $G^l \in \mathcal{R}^{N_l \times N_l}$ is the inner product between the vectorised feature maps i and j in layer l .

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

Style Loss:

If F_{ij}^l is the activation of i -th filter in the position j in the layer l .

Let $a \rightarrow$ original image



Let $x \rightarrow$ generated image



S.t. A^l and G^l are their respective feature response in layer l .

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2 \quad (4)$$

$$\mathcal{L}_{\text{style}}(\vec{a}, \vec{x}) = \sum_{l=0}^L w_l E_l, \quad (5)$$

Style Representation

For any image x ; the Gram matrix for layer l : $G^l \in \mathcal{R}^{N_l \times N_l}$ is the inner product between the vectorised feature maps i and j in layer l .

Style

If F

Let

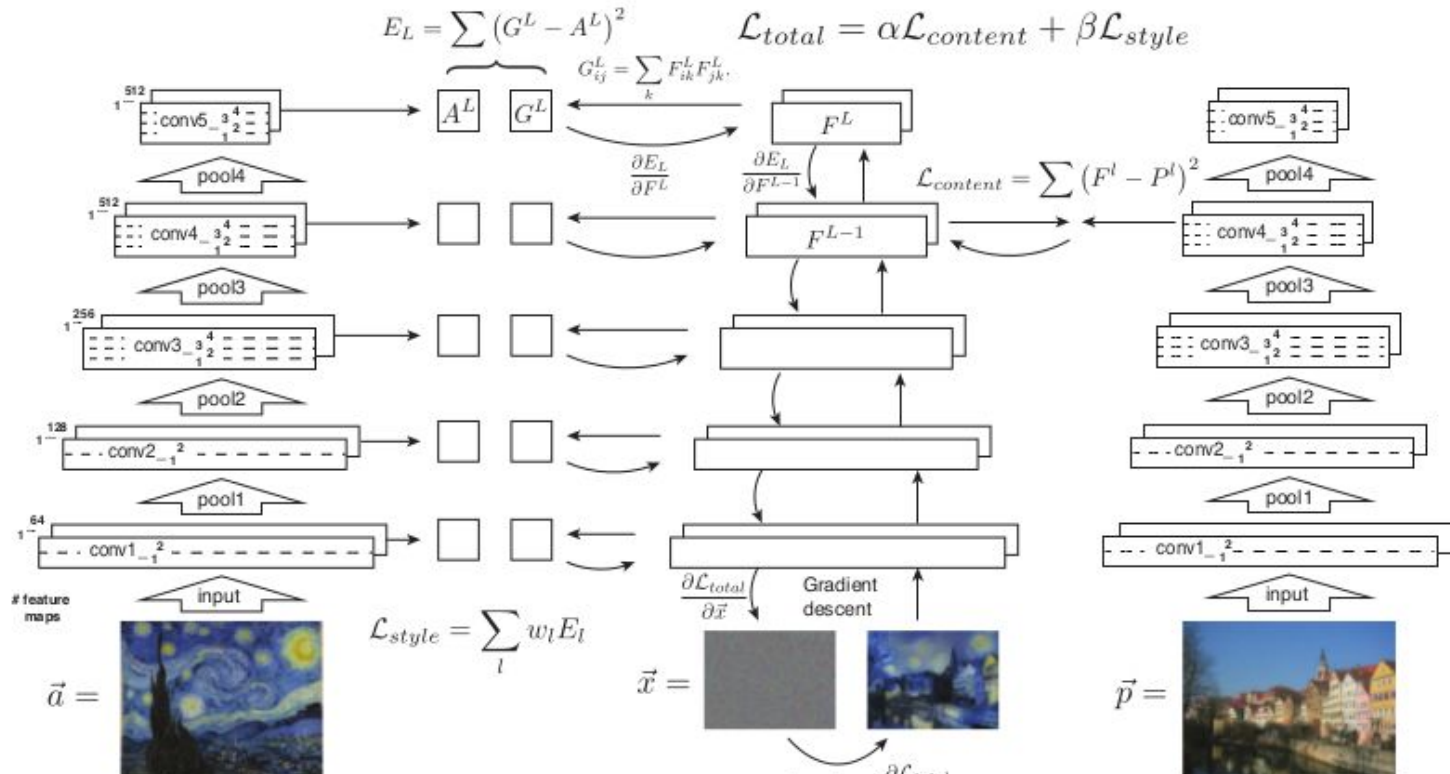
Let

$$\frac{\partial E_l}{\partial F_{ij}^l} = \begin{cases} \frac{1}{N_l^2 M_l^2} ((F^l)^T (G^l - A^l))_{ji} & \text{if } F_{ij}^l > 0 \\ 0 & \text{if } F_{ij}^l < 0. \end{cases} \quad (6)$$

ctive

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2 \quad (4) \quad \mathcal{L}_{\text{style}}(\vec{a}, \vec{x}) = \sum_{l=0}^L w_l E_l, \quad (5)$$

Style Transfer Algorithm



$$\mathcal{L}_{total}(\vec{p}, \vec{a}, \vec{x}) = \alpha \mathcal{L}_{content}(\vec{p}, \vec{x}) + \beta \mathcal{L}_{style}(\vec{a}, \vec{x}) \quad (7)$$

Trade-off between Content and Style: α/β

10^{-4}



10^{-3}



10^{-2}



10^{-1}



Trade-off between Content and Style: α/β

10^{-4}



10^{-3}



10^{-2}



10^{-1}



Selection of layer for content matching

Content Image



Conv2_2



Conv4_2



Effect of different initialization

A



Content image

B



Style image

Effect of different initialization (different white noises)



Thanks!