

# CVPR 2016

## Image Style Transfer Using Convolutional Neural Networks

**Leon A. Gatys**

Centre for Integrative Neuroscience, University of Tübingen, Germany

Bernstein Center for Computational Neuroscience, Tübingen, Germany

Graduate School of Neural Information Processing, University of Tübingen, Germany

**Alexander S. Ecker**

Centre for Integrative Neuroscience, University of Tübingen, Germany

Bernstein Center for Computational Neuroscience, Tübingen, Germany

Max Planck Institute for Biological Cybernetics, Tübingen, Germany

Baylor College of Medicine, Houston, TX, USA

**Matthias Bethge**

Centre for Integrative Neuroscience, University of Tübingen, Germany

Bernstein Center for Computational Neuroscience, Tübingen, Germany

Max Planck Institute for Biological Cybernetics, Tübingen, Germany

# Introduction: What is style transfer?

Merging **content** from one image with the artistic **style** of another



Image Source: [Tuebingen Neckarfront with beautiful old houses.](#)



Image Source: [The Starry Night by Vincent van Gogh, 1889](#)



## Problem Definition: Why style transfer is challenging?

- Content and style are deeply intertwined. Separation of content and style of a natural image is extremely difficult.
- Traditional methods rely on low-level features of the target image to inform the texture transfer (synthesise photorealistic natural textures by **resampling** the pixels of a given source texture)

## Solution: *A Neural Algorithm of Artistic Style*

A **texture transfer algorithm** that constrains a texture synthesis method by feature representations from state-of-the-art **Convolutional Neural Networks**.

# Background: Hierarchical Features Through CNNs

CNNs learn hierarchical features through:

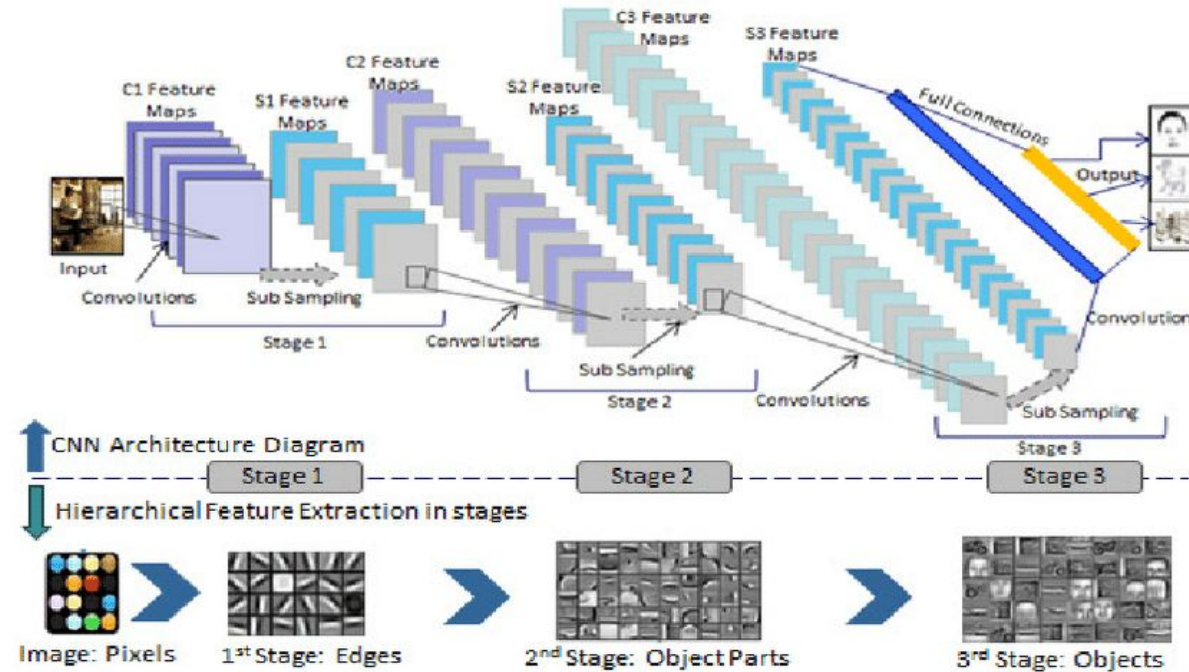


Image Source: [Learning hierarchy of visual features in CNN architecture](#)

- Lower layers: Basic pixel-based attributes (edges, colors, gradients etc.)
- Higher layer: High level details (objects & shapes, scenes & context, faces & expressions etc.)

- The key insight to this paper: CNNs can encode texture (style) and object structure (content) separately.



# Content Representation

- Content representations of images are extracted from higher layers of a pre-trained CNN (VGG-19 with 16 convolutional layers and 5 pooling layers).
- Content informations are encoded within higher layers of the CNN as **feature maps**.
- To visualize these content information encoded at different layers, gradient descent can be performed on a white noise image, using the following loss function.

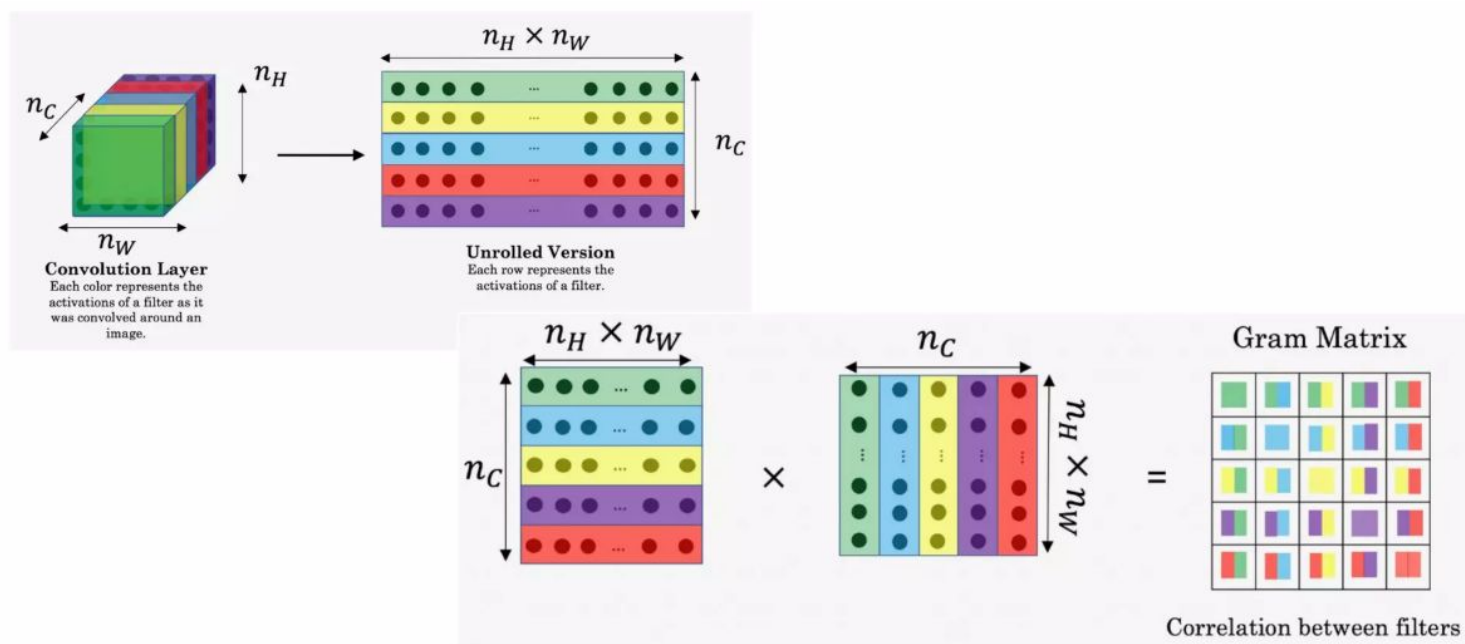
$$\mathcal{L}_{\text{content}}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2$$

$F_{ij}^l$ : feature representation (of the  $i^{\text{th}}$  filter at  $j^{\text{th}}$  position) of the generated image in the  $l^{\text{th}}$  layer.

$P_{ij}^l$ : feature representation (of the  $i^{\text{th}}$  filter at  $j^{\text{th}}$  position) of the original content image in the  $l^{\text{th}}$  layer.

# Style Representation

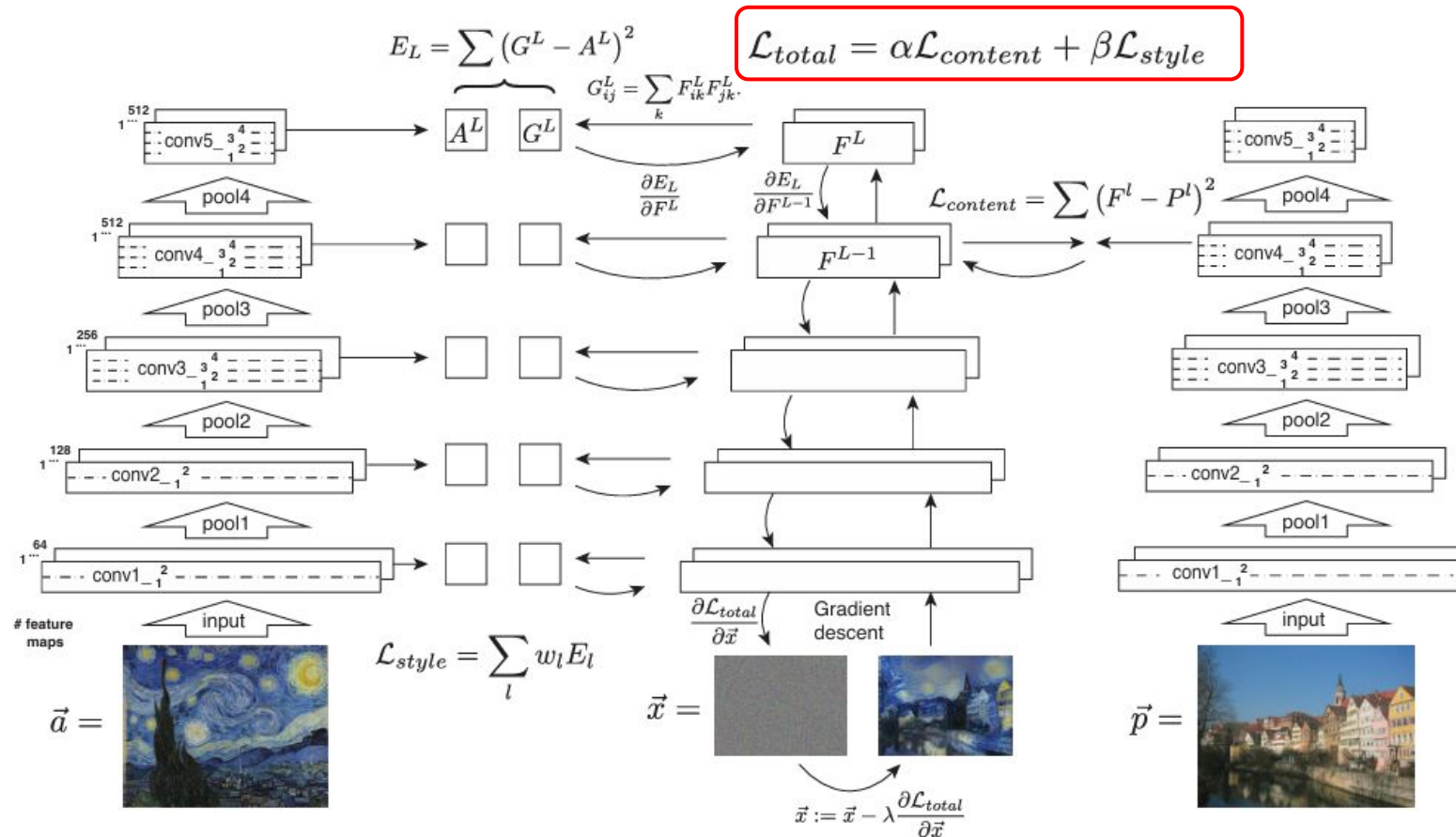
- A **feature space** that captures texture information was designed to obtain style representation of an image.
- The feature space consists with the correlations between different filters of a layer.
- Feature correlations are stored in a Gram matrix.



Multiple layers contribute to capturing stationary, multi-scale texture patterns.

# Style Transfer Algorithm

- Synthesise a new image that simultaneously matches content from one image and the style of an artistic image.



# Style Transfer Equation

- The total loss is calculated by:

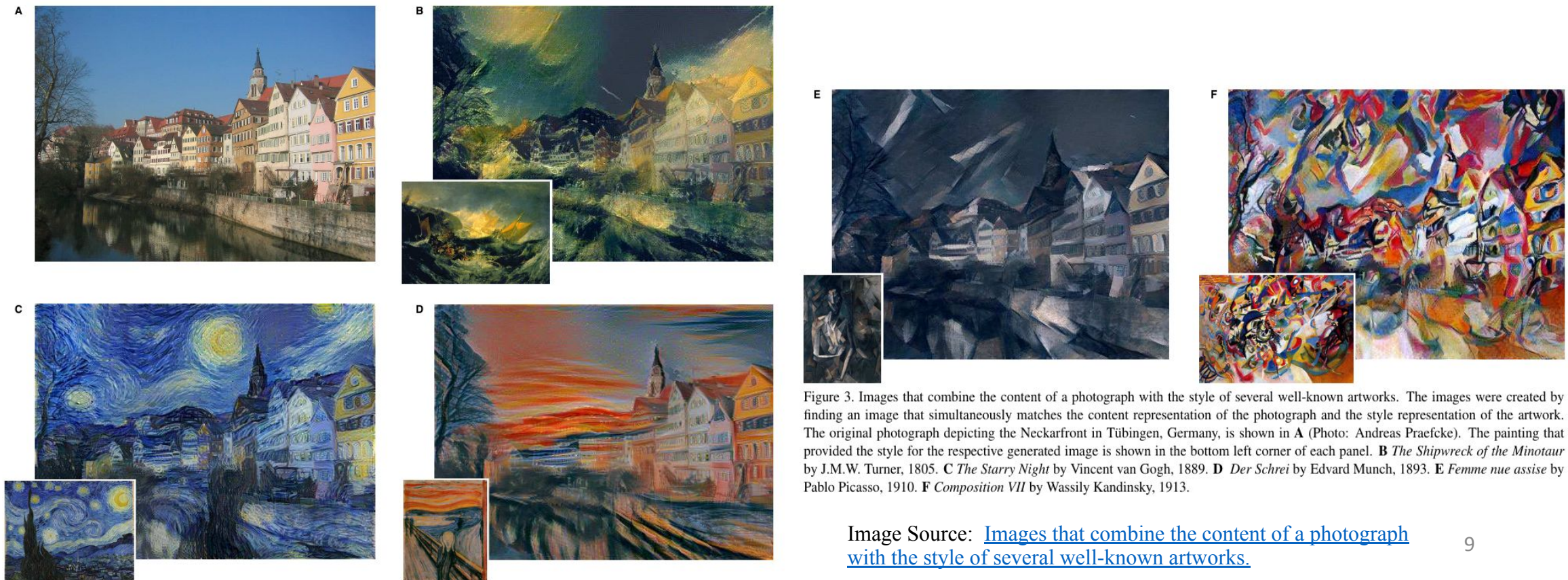
$$\mathcal{L}_{\text{total}}(\vec{p}, \vec{a}, \vec{x}) = \alpha \mathcal{L}_{\text{content}}(\vec{p}, \vec{x}) + \beta \mathcal{L}_{\text{style}}(\vec{a}, \vec{x})$$

- Tuning parameters:
  1. High  $\alpha$ : more content preserved
  2. High  $\beta$ : stronger style influence
- Optimization Method: **L-BFGS** (Limited-memory Broyden-Fletcher-Goldfarb-Shanno)



# Results of Style Transfer

- Key findings:
  - a. Content and style representation of a CNN is **well separable** (in theory).
  - b. We can smoothly regulate the emphasis on either reconstructing the content or the style (by manipulating  $\alpha/\beta$  ratio).





# Effect of Content Layer Selection

- Matching from lower layers → more detailed pixel information retained, but the texture is merely blended over the image (**weaker style transfer**).
- Matching from higher layers → the fine structure of the image is altered such that it agrees with the style of the artwork (**better fusion of content and style**).



# Effect of Style Layers

- Matching the style representation upto higher layers preserves local images structures an increasingly large scale.
- The style loss for a particular layer:

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

$G_{ij}^l$ : style representation (of the  $i^{\text{th}}$  filter at  $j^{\text{th}}$  position) of the generated image in the  $l^{\text{th}}$  layer.

$A_{ij}^l$ : style representation (of the  $i^{\text{th}}$  filter at  $j^{\text{th}}$  position) of the original image in the  $l^{\text{th}}$  layer.

- The loss from each individual layer will contribute to the total style loss as follows:

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

# Trade-off Between Content & Style

- Adjusting the  $\alpha/\beta$  ratio of the following loss equation will affect the results.

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

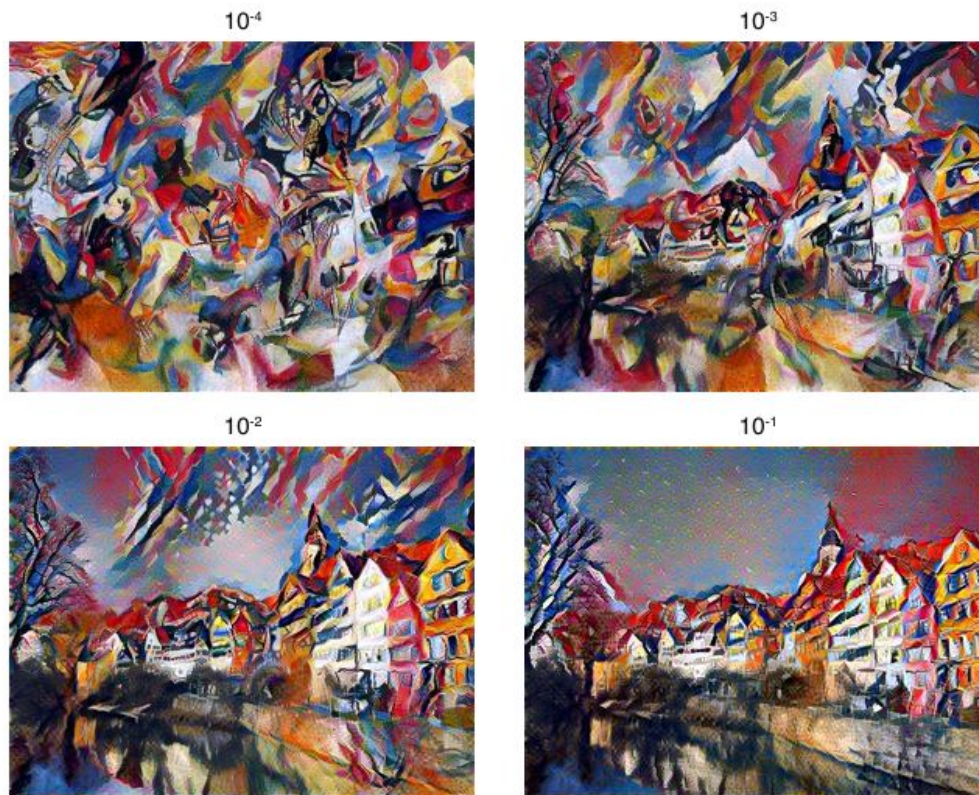


Figure 4. Relative weighting of matching content and style of the respective source images. The ratio  $\alpha/\beta$  between matching the content and matching the style increases from top left to bottom right. A high emphasis on the style effectively produces a texturised version of the style image (top left). A high emphasis on the content produces an image with only little stylisation (bottom right). In practice one can smoothly interpolate between the two extremes.



# Impact of the Initialization

- We can initialize the style transfer algorithm either on **white noise** or some **deterministic content**.
  - white noise → diverse outputs
  - content image → more stable, deterministic results

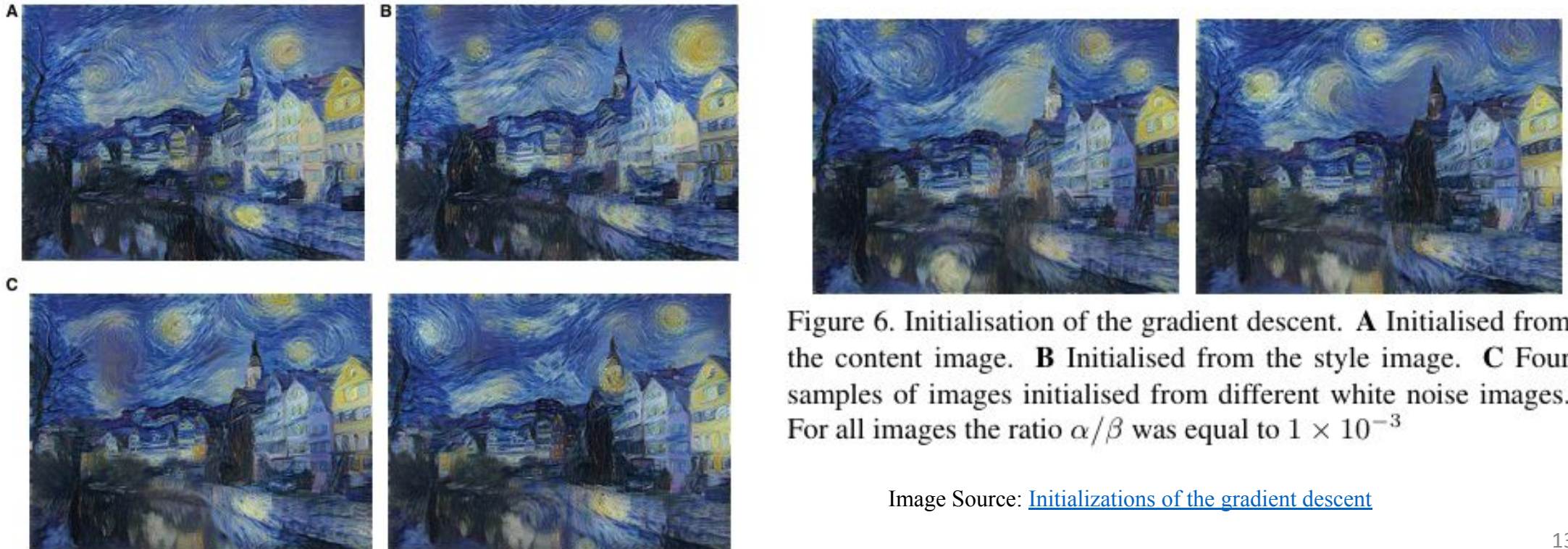
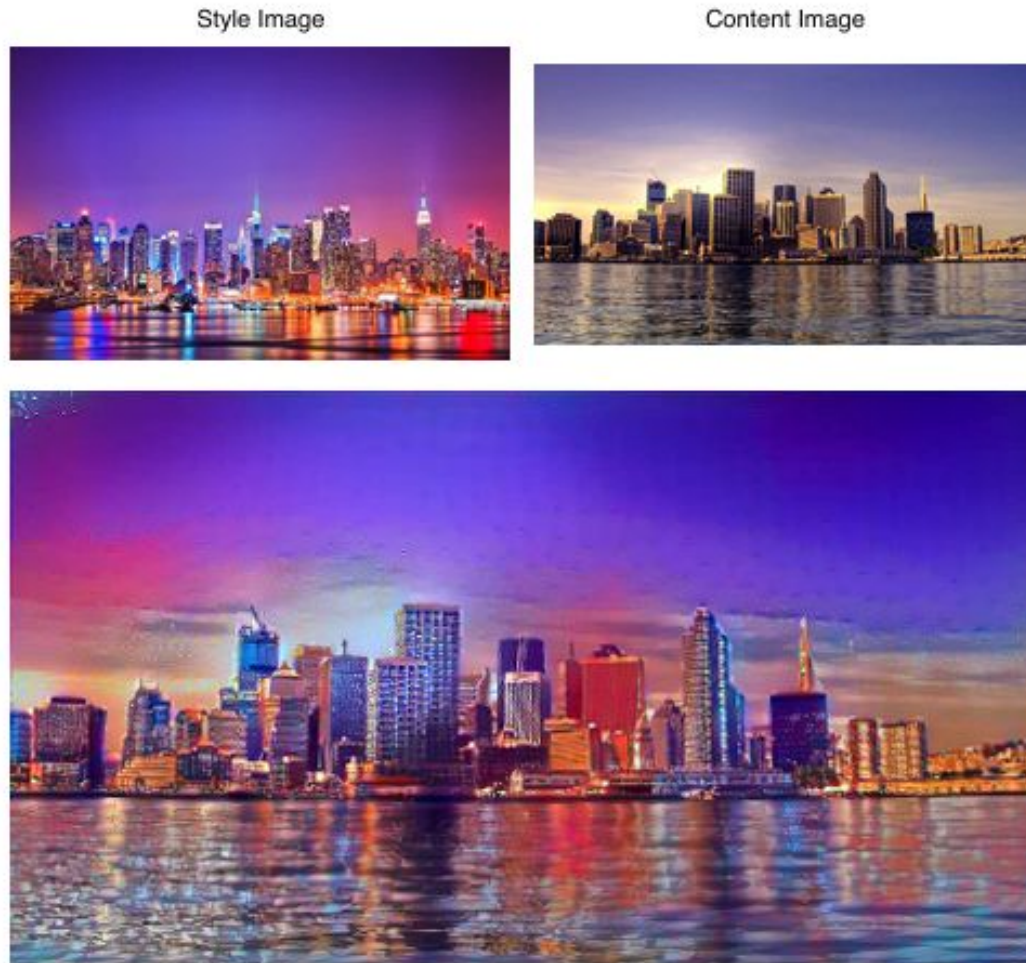


Image Source: [Initializations of the gradient descent](#)



# Photorealistic Style Transfer



- Key challenge: It's hard to maintain the realism while changing textures (applying styles).

Figure 7. Photorealistic style transfer. The style is transferred from a photograph showing New York by night onto a picture showing London by day. The image synthesis was initialised from the content image and the ratio  $\alpha/\beta$  was equal to  $1 \times 10^{-2}$

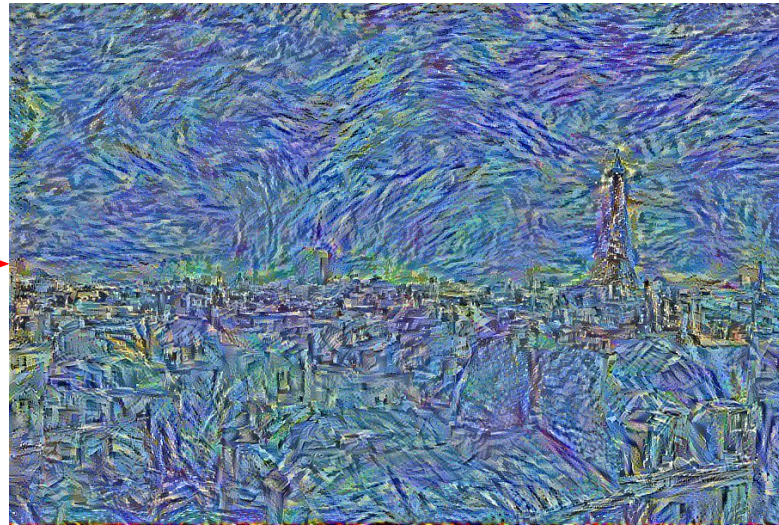


# Limitations

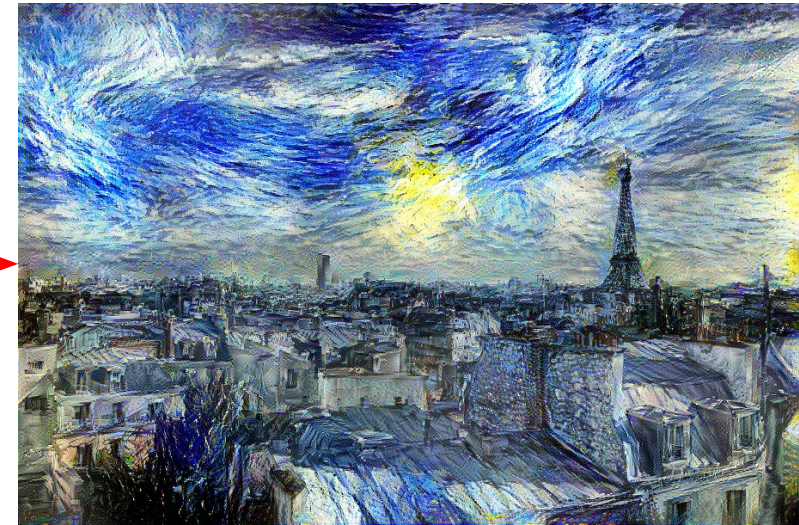
- When the style is too dominant, some content loss can occur.
- It's very difficult to selectively control how much style affects which parts.
- Limitations in computation:
  - a. Optimization is iterating and higher iterations are needed for better visual quality.
  - b. The speed of synthesis process heavily depends on the image resolution.



Original Image



Style Transferred (after 100 iterations)



Style Transferred (after 3000 iterations)

**Let's dive into the code!**

**Thank You**