# ✧ Image Style Transfer ✧

Gatys et. al. 2016

## 1. Introduction

Texture analysis → Efros & Freeman (correspondence map)
→ Hertzman et.al. (image analogies)
→ Ashikhmin (transferring the high-frequency
→ Lee et.al. improved texture info +preserve
+ edge orientation info.        coarse scale))

remarkable results but limitation = only low-level image features
of the target img to inform the
texture transfer.

◆ a style transfer algo. should be able to extract the semantic img content from the target img. (e.g. objects & scenery.) and then inform a texture transfer procedure to render the semantic content of the target img. in the style of the source img

➡ to seperate content from style, extract high-level semantic info. by using Convolutional Neural Networks (CNNs).

introduced: A Neural Algo. of Artistic Style.
✚ combines a parametric texture model based on CNNs with a method to invert their img. ~~represtat~~ representations.

## 2. Deep Image Representations

• basis → VGG network which is for object recognition & localization.

VGG-19 } 16 convolutional & 5 pooling layers
total 19 layers = 16 conv + 3 fully connected

**◉** network is normalized by scaling the weights
→ the mean activation of each conv. filter over imgs. and positions is equal to one.
**✸** re-scaling can be done for the VGG net. without changing its output; because it contains only rectifying linear activation (ReLU) funcs. and no normalization or pooling over feature maps.

- not used any of fully connected layers.
- for img. synthesis avg. pooling is slightly better than max. pooling.

## 2.1 Content Representation

- each layer in network defines a non-linear filter bank.
  ✚ complexity of layer increases with the position of layer.
  ✚ hence, input img. $\vec{x}$ is encoded in each layer of CNN by filter responses.

$N_l$ → distinct layers
$M_l^l$ → size of each feature map (height x width)
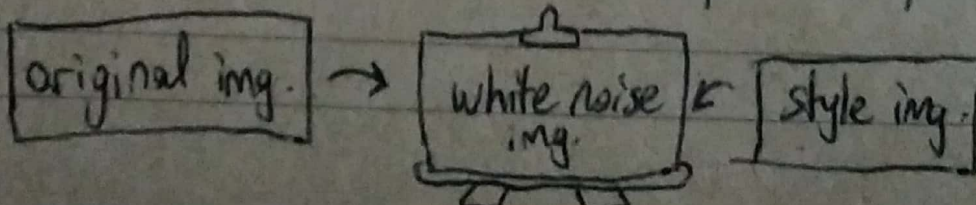$F^l$ → respons' in layer $l$
↓
matrix $\quad F^l \in R^{N_l \times M_l}$, $\quad F_{ij}^l$ → activation of the $i^{th}$ filter at position $j$ in layer $l$.

☼ perform gradient descent on a white noise img. to find another img. that matches the features responses of original img.

♻ original img. → white noise img. ← style img.

let $\vec{p}$ & $\vec{x}$ be original img., $P^{\ell}$ & $F^{\ell}$ their respective feature representations in layer $\ell$.

(white-noise)

Then, squared-error loss:

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

Derivative of loss func. wrt. activations in layer $\ell$:

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

$\vec{x}$ can be computed using standard error back-propagation

> low layers reproduce exact pixel values of original img.
> so, higher layers are preferred
> → as content representation.

## 2.2. Style Representation

- for style, a feature space to capture texture info.
  + this feature map can be built on top of the filter responses in any layer of net.
  + consists of correlations between different filter responses → represented by Gram matrix.

$G^{\ell} \in R^{N_{\ell} \times M_{\ell}}$, $G_{ij}^{\ell}$ → inner product between vectorised feature map $i$ and $j$ in $\ell$.

$$G_{ij}^{\ell} = \sum_{k} F_{ij}^{\ell} \cdot F_{jk}^{\ell}$$

- feature correlations of multiple layers → stationary a
  + texture info, not global arrangement (multiscale representation)

◆ visualize info. by style feature spaces :

$\overset{\ast}{\underset{\ast}{\overset{\ast}{\ast}}} \overset{\ast}{\underset{\ast}{\ast}} \overset{\ast}{\longrightarrow}$ by using gradient descent from a white noise img. to minimise the mean-squared distance between entries of Gram matrices from original img. & Gram matrices of img. to be generated.

let $\vec{a}$ & $\vec{x}$ be original img. and generated img.
$A^l$ & $G$ respectively style representation in layer $l$.
Then, contribution of layer $l$ to total loss:

$$E_\ell = \frac{1}{4 N_\ell^2 M_\ell^2} \sum_{ij} \left( G_{ij}^\ell - A_{ij}^\ell \right)$$

Total style loss :

$$\mathcal{L}_{style}(\vec{a}, \vec{x}) = \sum_{l=0}^{L} w_\ell \cdot E_\ell$$

→ weighting factor

Derivative of $E_\ell$ :

$$\frac{\partial E_\ell}{\partial F_{ij}^\ell} = \begin{cases} \frac{1}{N_\ell^2 M_\ell^2} \left( (F^\ell)^T \cdot (G^\ell - A^\ell) \right)_{ij} & \text{if } F_{ij}^\ell > 0 \\ 0 & \text{if } F_{ij}^\ell < 0 \end{cases}$$

## 2.3. Style Transfer

The loss function that is minimized :

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

- $\alpha$ & $\beta$ are weighting factors.
- $\dfrac{\partial h_{total}}{\partial \vec{x}}$ can be used as input for some numerical optimisation strategy.

- here L-BFGS is used ( best for img. synthesis)

- always resized style imgs to same size as the content.

- not regularised synthesis results with img. priors.

## 3. Results

### 3.1. Trade-off between content & style matching

### 3.2. Effect of different layers of the CNN

- Matching the style representations up to higher layers in the network preserves local imgs. structures on increasingly large scale. ⊛ 'conv1_1', 'conv2_1', 'conv3_1', 'conv4_1' and 'conv5_1'

- lower layers → more content from original img.

### 3.3. Initialisation of Gradient Descent

- white noise is used but one can also initialize it with content or style img; but it causes bias.

### 3.4. Photorealistic Style Transfer

# 4. Discussion

- Most limiting factor is resolution of the synthesised img.

- dimensionality of the optimisation problem | grow
  + the number of units in the CNN | linearly
  so, speed depends on resolution — with # of pixels.

$512 \times 512$ pixels } in this → prohibits online & interactive
Nvidia K40 GPU } paper   applications
one hour

- seperation of img. content from style is not well-defined problem.

## Bonus: L-BFGS : Limited-memory BFGS

- an optimization algorithm that approximates BFGS* algorithm (Broyden-Fletcher-Goldfarb-Shanno) using a limited amount of computer memory. family: quasi-Newton methods.
  * an iterative method for solving unconstrained ~~linear equa~~ non-linear optimization. $O(n^2)$

- popular algo. for parameter estimation in ML.
- target problem is to minimize $f(x)$ over constrained values.

Quasi-Newton Method: used to either find local maxima & minima of functions or zeroes. They can be used if the Jacobian or Hessian is unavailable or too expensive.

- L-BFGS aka "the algo. of choice" for fitting log-linear (MaxEnt) and conditional random fields with $l_2$-regularization.

◆ a way of finding a (local) minimum of an objective func., making use of objective func. values and the gradient of the obj. func.