

# Neural Style Transfer

LEON A. GATYS, ALEXANDER S. ECKER, MATTHIAS BETHGE

Werner Chao

Felipe Perez

Xiyang Chen

March 14<sup>th</sup>, 2019

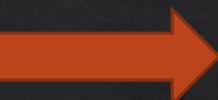
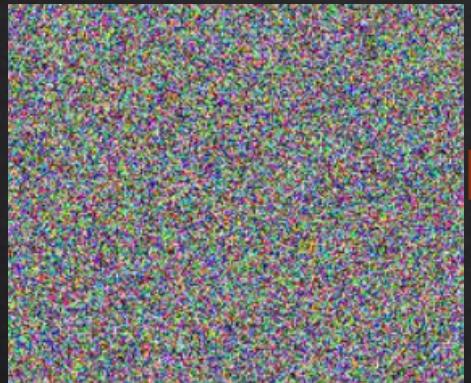
# Gatys et al 2015



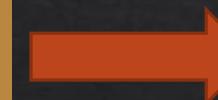
# Motivation

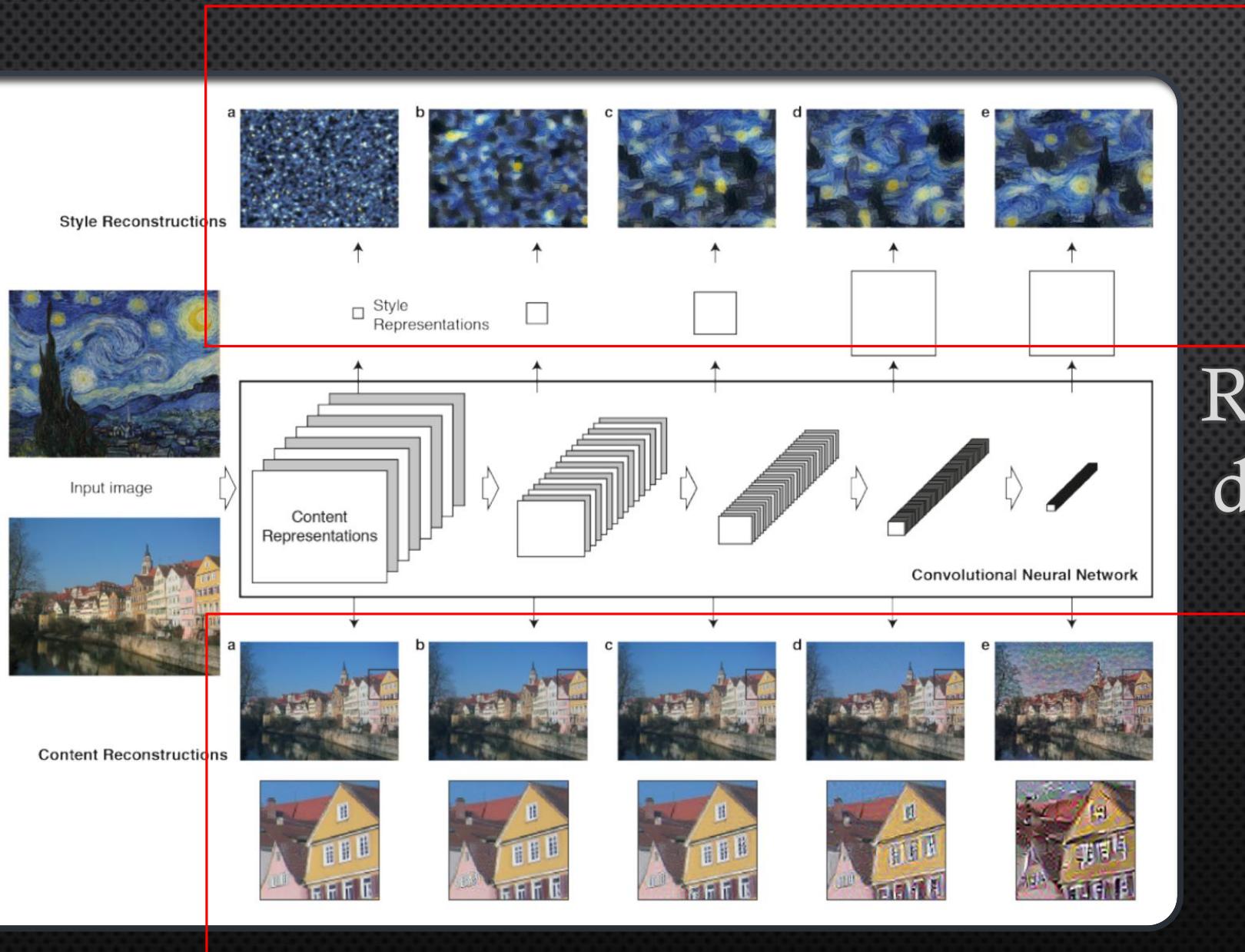
- ❖ Humans are great at creative things, such as art
- ❖ We don't understand well why humans are able to do this
- ❖ No artificial intelligence system exist that can create art parallel to human
- ❖ Can we have an algorithmic understanding how humans create art?

# Overall process



LOSS  
Function





# Reconstructions from different CNN layers In VGG Network

# Content and style in cnn are separable

$$L_{total}(\vec{c}, \vec{s}, \vec{x}) = \alpha L_{content}(\vec{c}, \vec{x}) + \beta L_{style}(\vec{s}, \vec{x})$$

- ◆  $\vec{c}$  = content
- ◆  $\vec{s}$  = style
- ◆  $\vec{x}$  = white noise
- ◆  $\alpha/\beta = 1e - 3 \text{ or } 1e - 4$

# Important notes

- ❖ No model learning: frozen VGG network
- ❖ image reconstruction only
- ❖ Target:
  - ❖ activation from cnn layer(s) in pre-trained vgg network
- ❖ Only 1 pair of content and style image are used

# Content loss

$$L_{content}(\vec{c}, \vec{x}, l) = \frac{1}{2} \sum_{i,j,k} (F_{i,j,k}^l - C_{i,j,k}^l)^2$$

White noise

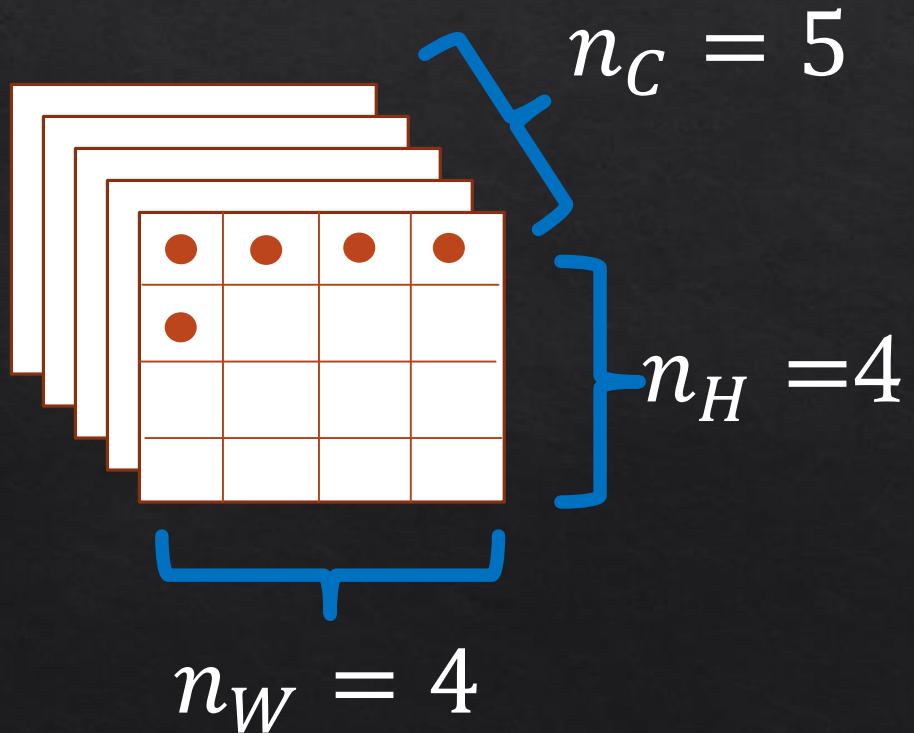


Target content, output  
from a VGG mid-layer

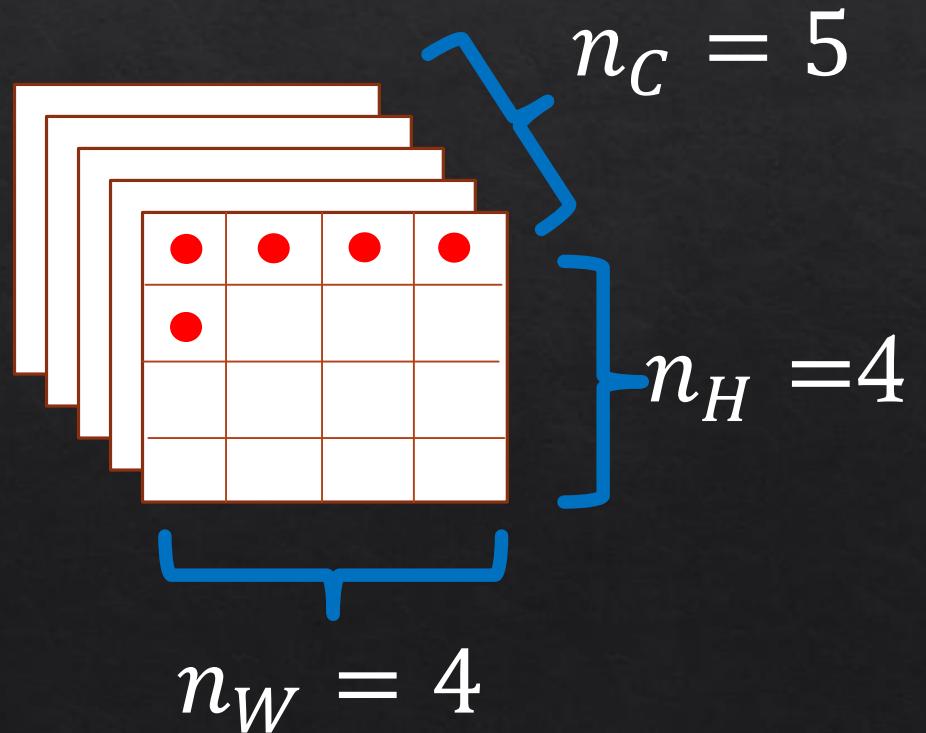
- ❖  $i$ : index of height
- ❖  $j$ : index of width
- ❖  $k$ : index of channel
- ❖  $l$ : certain CNN layer in VGG

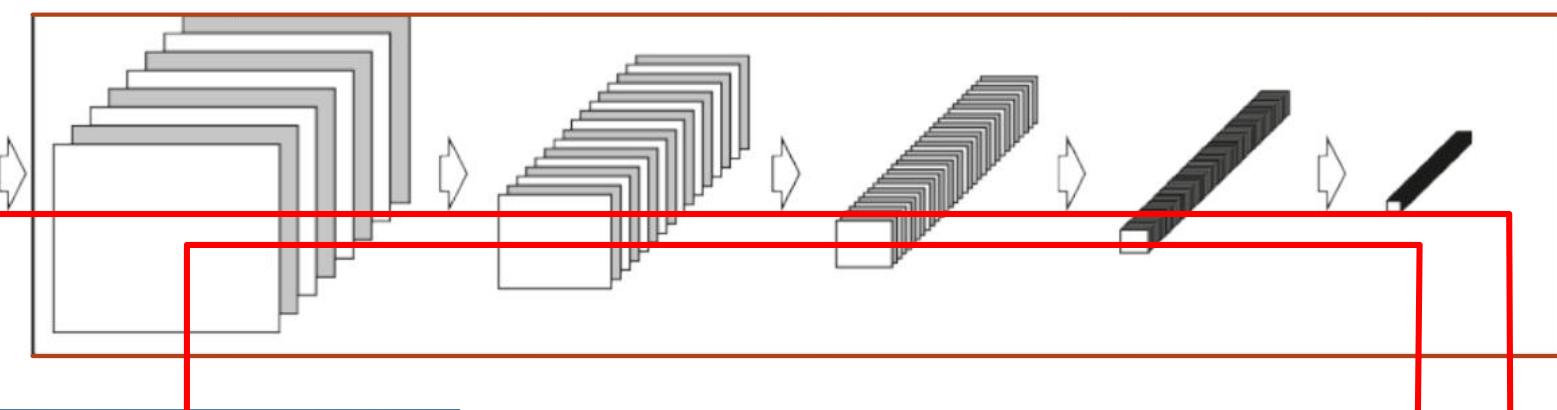
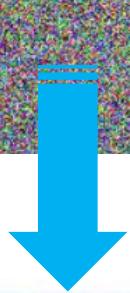
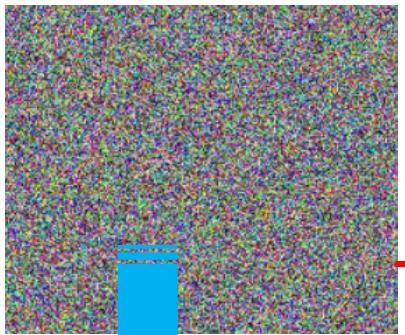
# Content loss – output of 1 cnn layer in vgg

Content Image



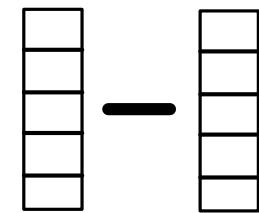
White Noise



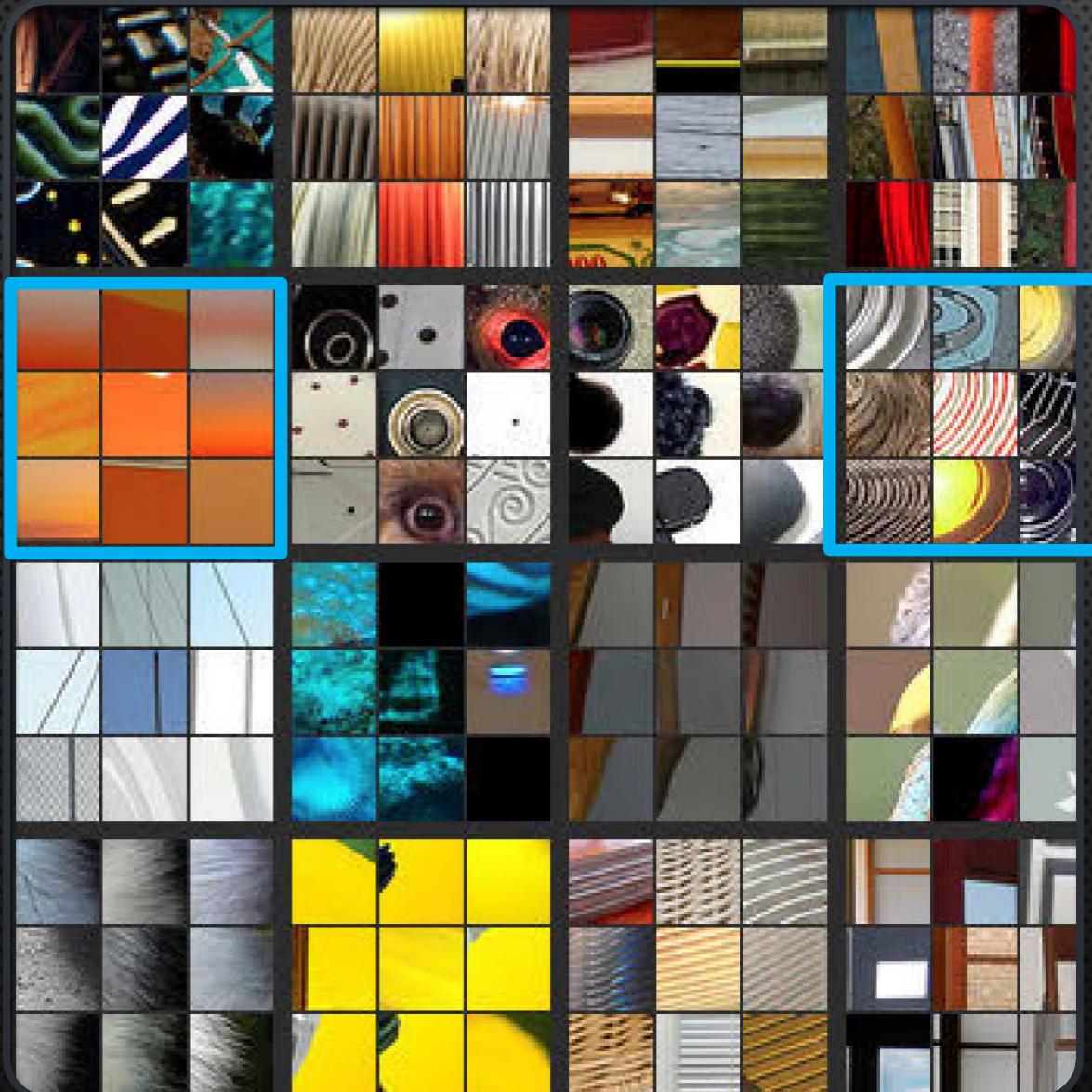


*Minimize*

Target  
content



White  
noise



# Gram matrix - correlation intuition

# Style loss

$$L_{style}(\vec{s}, \vec{x}, l) = \frac{1}{(2n_W n_H n_c)^2} \sum_{k,k'} (G_{k,k'}^l, -S_{k,k'}^l)^2$$

White noise

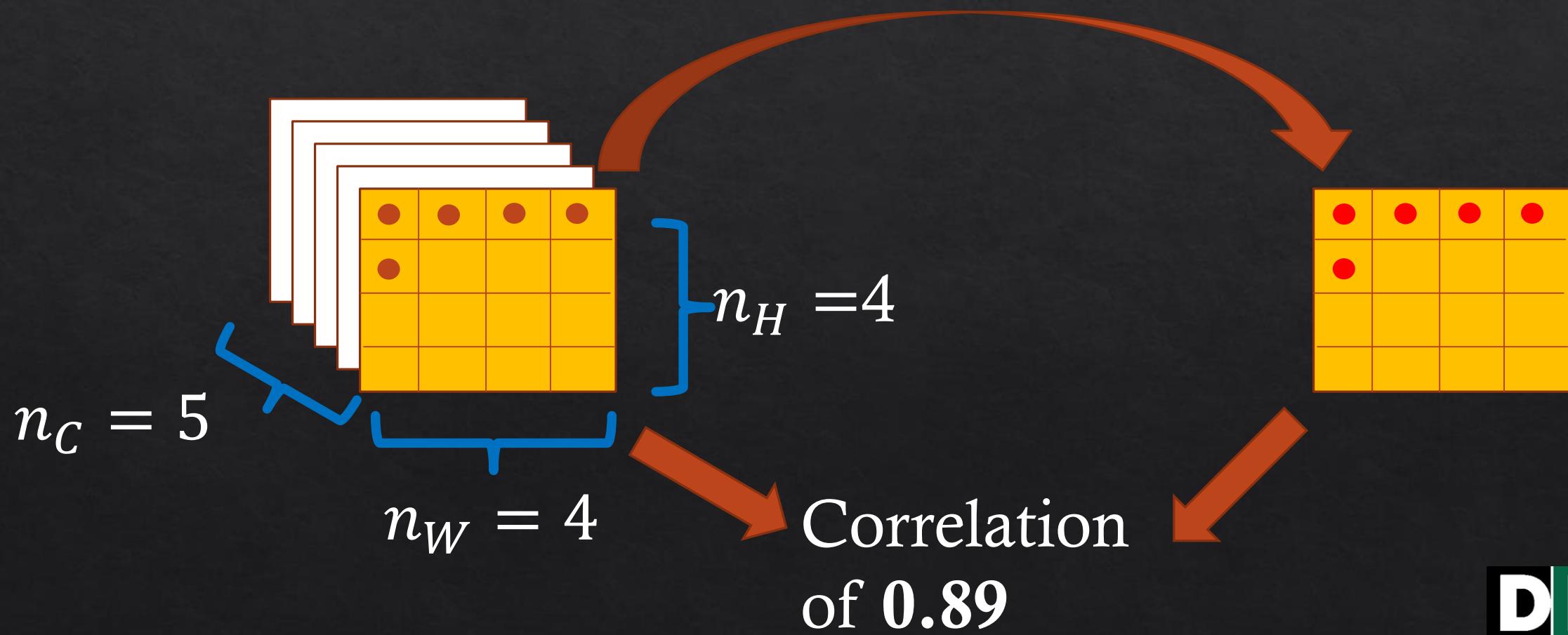
Target style,  
output from a  
VGG CNN layer

$$G_{k,k'}^l = \sum_{i,j} a_{i,j,k}^l a_{i,j,k'}^l$$

- ◊ Gram matrix: inner product between channels

# Style loss – gram matrix

Style Image



Style loss –  
gram matrix

$$G_{k,k'}^l = \sum_{i,j} a_{i,j,k}^l a_{i,j,k'}^l$$

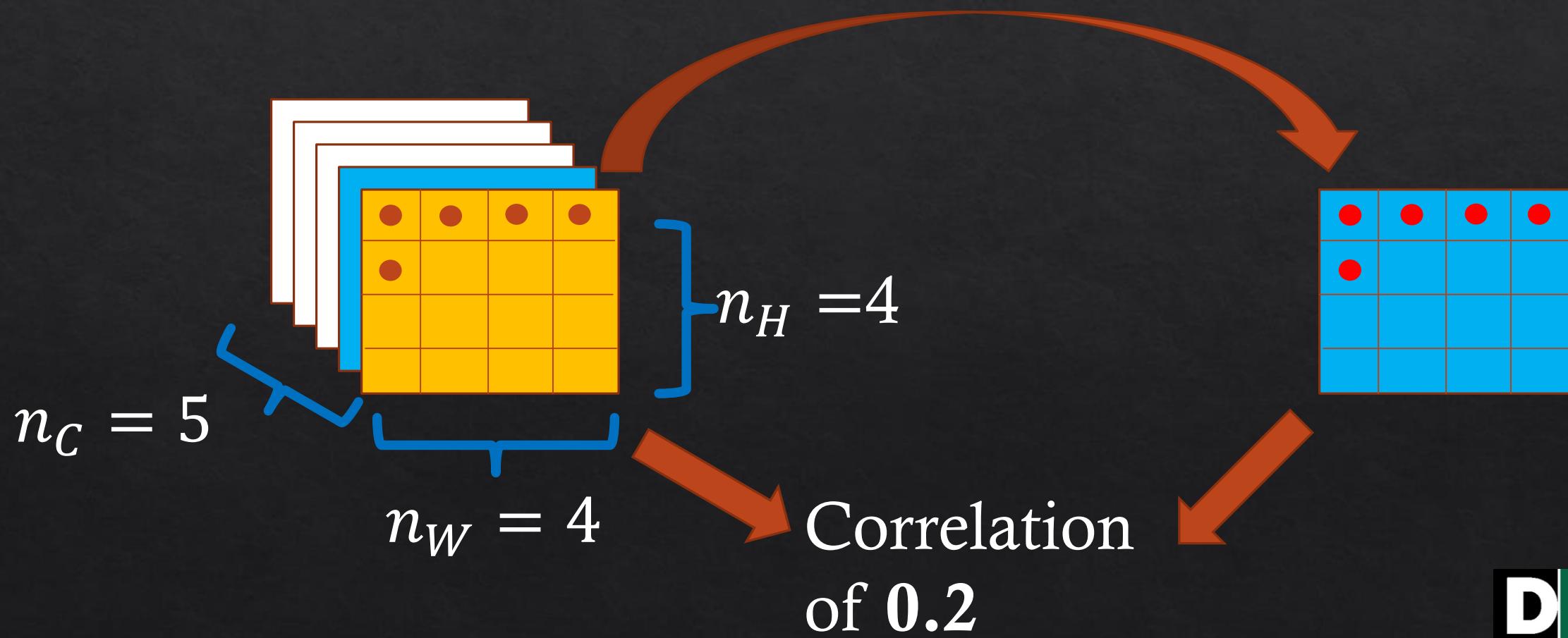
$n_C$

$n_C = 5$

0.89				

# Style loss – gram matrix

Style Image



# Style loss – gram matrix

$$G_{k,k'}^l = \sum_{i,j} a_{i,j,k}^l a_{i,j,k'}^l$$

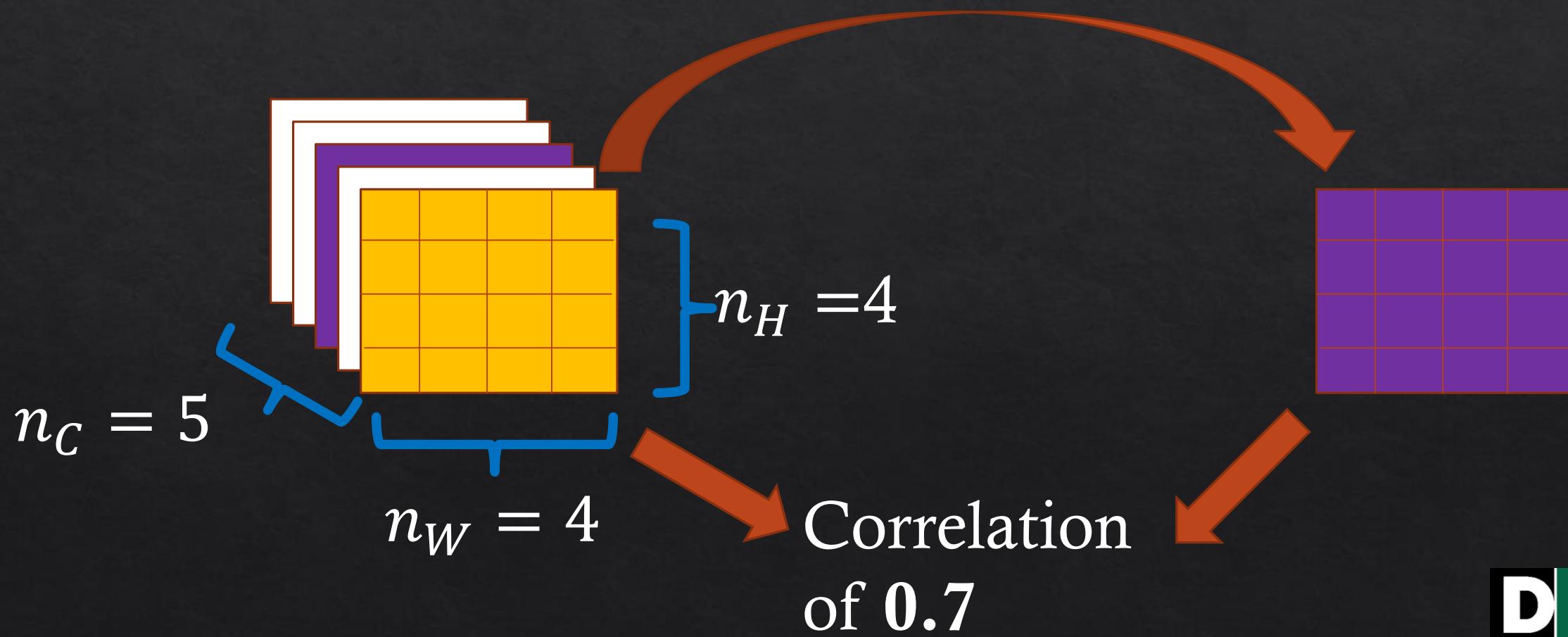
$n_C$

$n_C$

0.89	0.2			

# Style loss – gram matrix

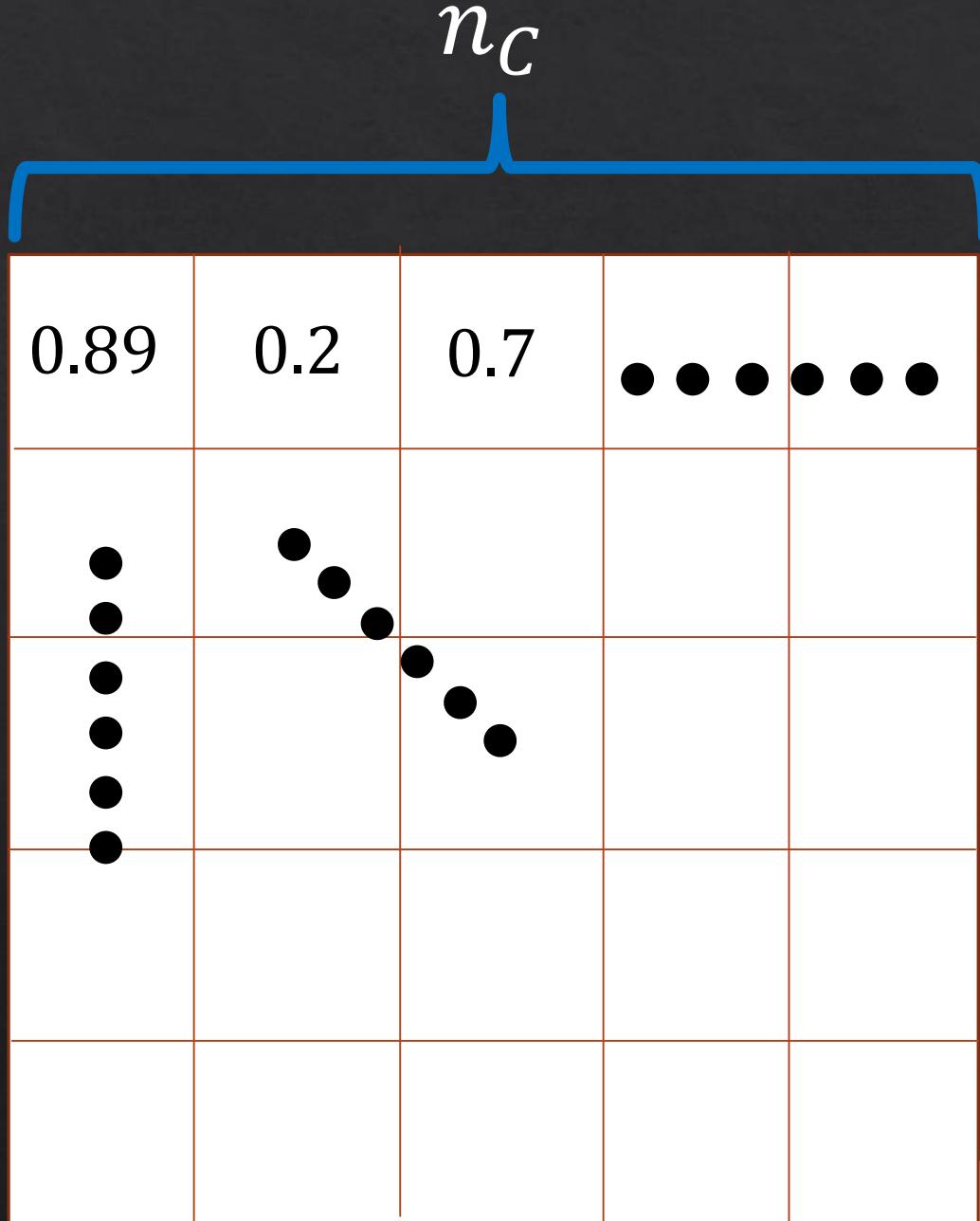
Style Image



# Style loss – gram matrix

$$G_{k,k'}^l = \sum_{i,j} a_{i,j,k}^l a_{i,j,k'}^l$$

$n_C$

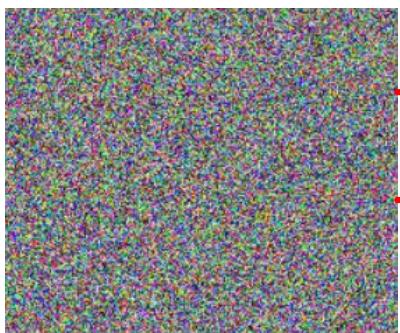


# Style loss

$$L_{style}(\vec{s}, \vec{x}, l) = \frac{1}{(2n_W n_H n_c)^2} \sum_{k,k',l} (G_{k,k'}^l - S_{k,k'}^l)^2$$

**White Noise  
Gram Matrix**   
**Target Style  
Gram Matrix** 

- ❖ Gram matrix: inner product between channels



$$\text{Min} \left( \begin{array}{c|c} G & G \\ \hline - & \end{array} \right) \quad \text{Min} \left( \begin{array}{c|c} G & G \\ \hline - & \end{array} \right) \quad \text{Min} \left( \begin{array}{c|c} G & G \\ \hline - & \end{array} \right) \quad \text{Min} \left( \begin{array}{c|c} G & G \\ \hline - & \end{array} \right) \quad \text{Min} \left( \begin{array}{c|c} G & G \\ \hline - & \end{array} \right)$$

*Minimize*

Target  
content

White  
noise

Varying layers  
and  $\alpha/\beta$

$$L_{total} = \alpha L_{content} + \beta L_{style}$$

Deeper  
layers for  
style

Increasing  $\alpha/\beta$  ratio

$10^{-5}$

A Conv1\_1



$10^{-4}$



$10^{-3}$



$10^{-2}$



B Conv2\_1



C Conv3\_1



D Conv4\_1



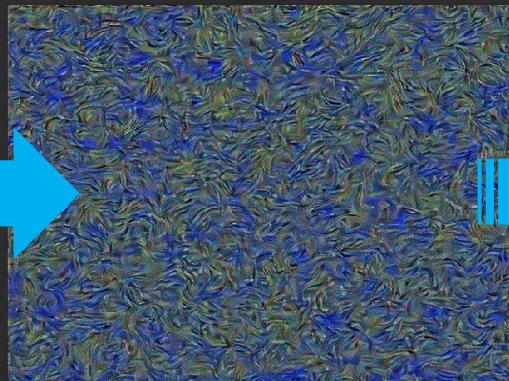
E Conv5\_1



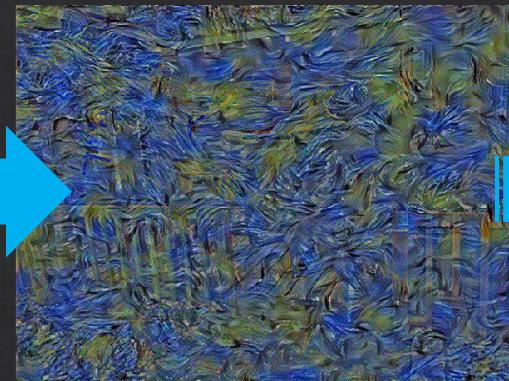
# Optimizer - LBFGS

- ❖ Tracks the second derivative, hessian matrix
- ❖ Terrible if we are doing model learning, millions of weights
- ❖ But, works well for image approximation

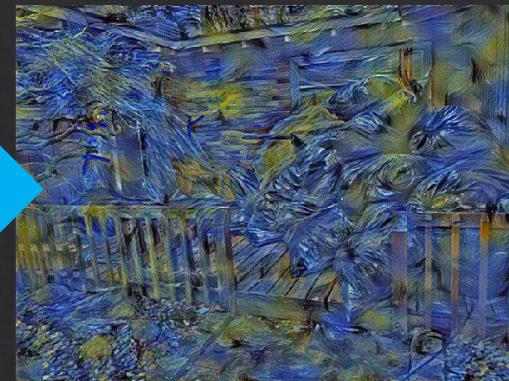
Epoch 100



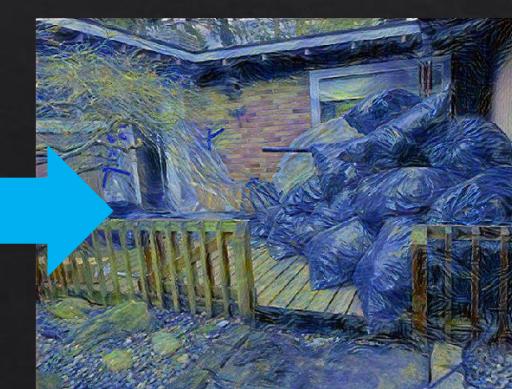
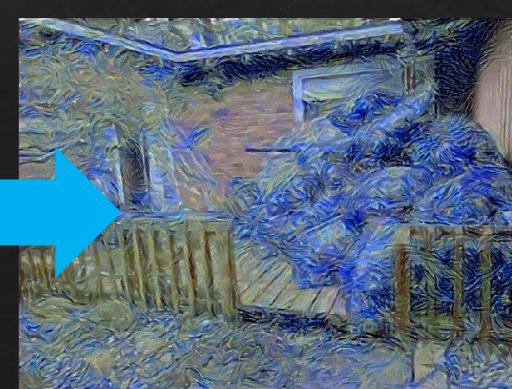
Epoch 300



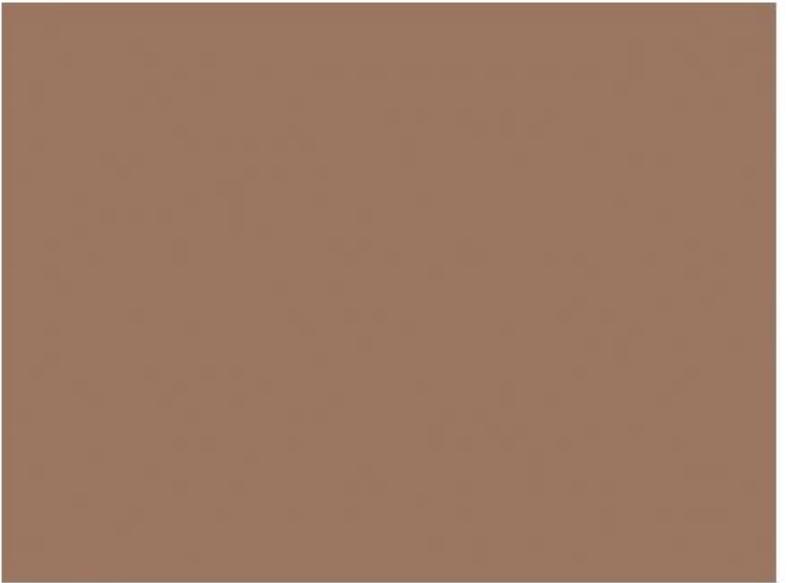
Epoch 500



100 200 300 400







**D L R**  
workshops

# Summary

- ❖ Content and style are separable mathematically
- ❖ The first time that any arbitrary artistic style can be transferred
- ❖ In the past is all filter based/hand crafted techniques
- ❖ Does not need (content, style) pair ground truth

# BREAK

# Discussion

- ❖ Can we separate content/style for data other than images/videos, such as voice, transactions?
- ❖ Demystifying Neural Style Transfer: can we reformulate Gram Matrix to have more flexibility?

# Appendix

# Code

- ❖ Official: <https://github.com/leongatys/PytorchNeuralStyleTransfer/blob/master/NeuralStyleTransfer.ipynb>
- ❖ Major Unofficial: [https://pytorch.org/tutorials/advanced/neural\\_style\\_tutorial.html](https://pytorch.org/tutorials/advanced/neural_style_tutorial.html)
- ❖ Paper: <https://arxiv.org/abs/1508.06576>

# Style loss – gram matrix

- Diagonal: which channels are most active
- Off-diagonal: which channel pairs co-occur

