

CSE 152: Computer Vision

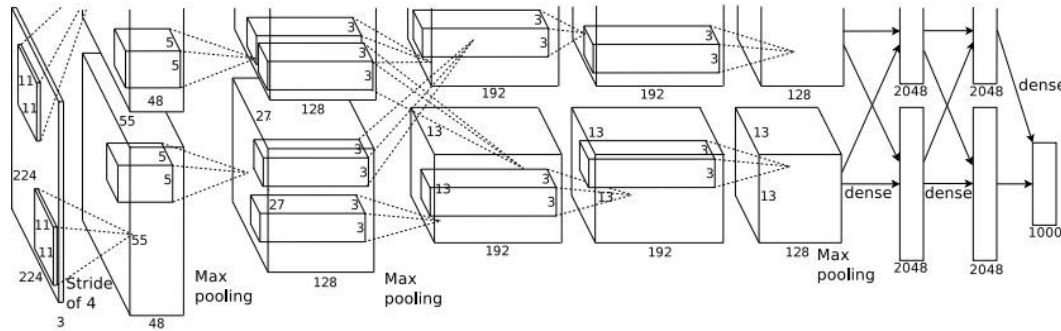
Hao Su

Lecture 11: Visualizing Networks, Fun Applications



What's going on inside ConvNets?

[This image is CC0 public domain](#)



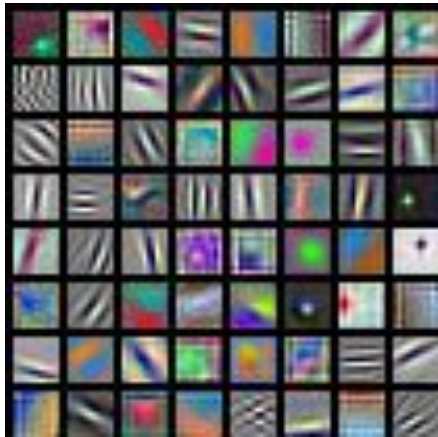
Class Scores:
1000 numbers

Input Image:
3 x 224 x 224



What are the intermediate features
looking for?

First Layer: Visualize Filters



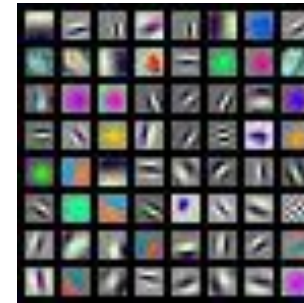
AlexNet:
64 x 3 x 11 x
11



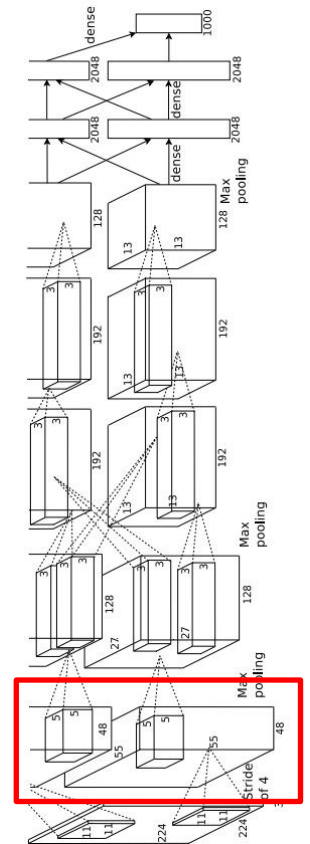
ResNet-
18: 64 x 3
x 7 x 7



ResNet-1
01: 64 x 3
x 7 x 7



DenseNet-
121: 64 x
3 x 7 x 7



Visualize the filters/kernels (raw weights)

We can visualize filters at higher layers, but not that interesting

(these are taken from ConvNetJS CIFAR-10 demo)

Weights:


layer 1 weights

16 x 3 x 7 x 7

Weights:


layer 2 weights

20 x 16 x 7 x 7

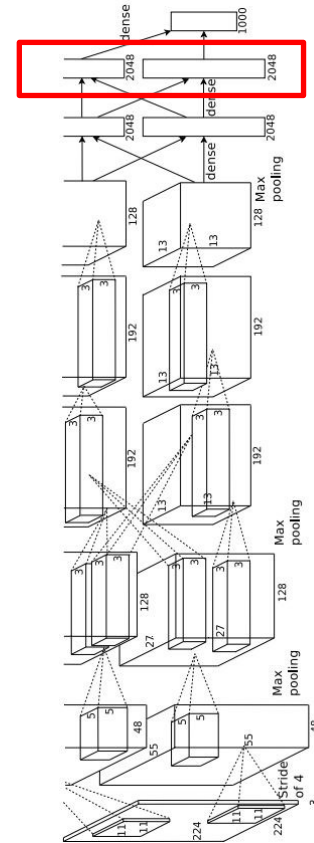
Weights:


layer 3 weights

20 x 20 x 7 x 7

Last Layer: Nearest Neighbors

4096-dim vector



Test image L2 Nearest neighbors in feature space



Recall: Nearest neighbors
in pixel space

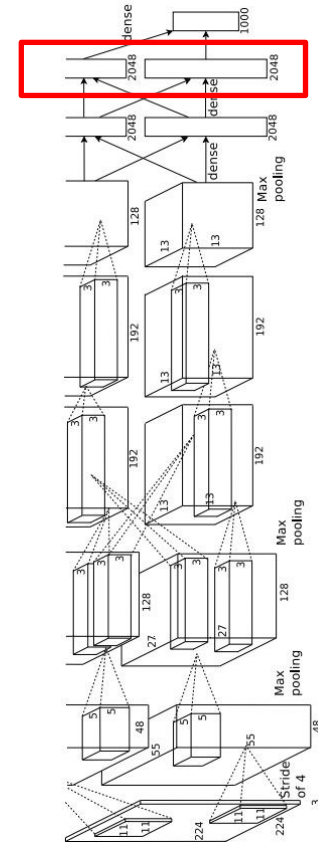
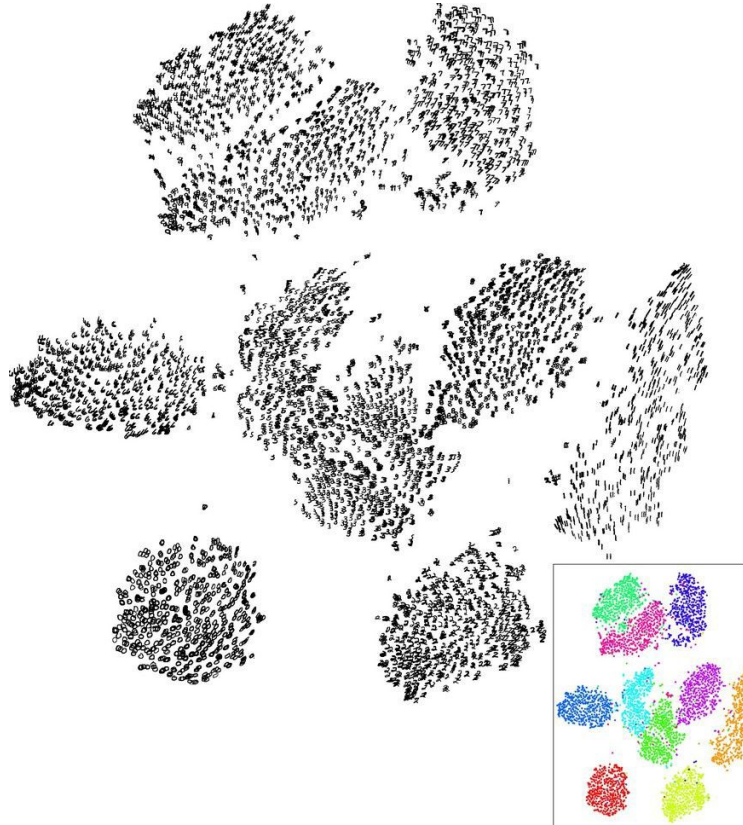


Last Layer: Dimensionality Reduction

Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

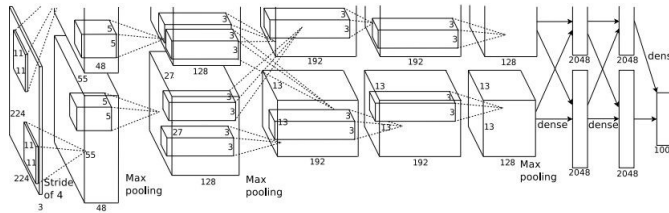
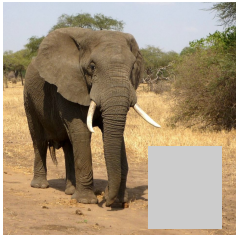
Simple algorithm: Principal Component Analysis (PCA)

More complex: **t-SNE**

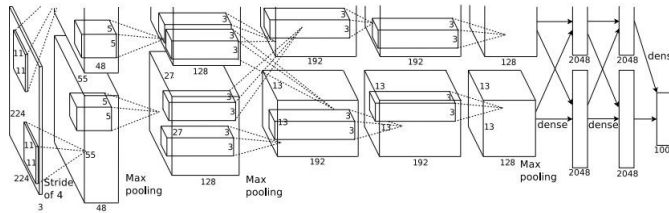
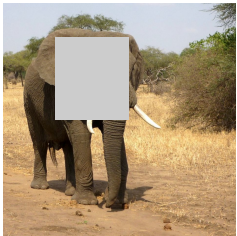


Which pixels matter: Saliency vs Occlusion

Mask part of the image before feeding to CNN, check how much predicted probabilities change



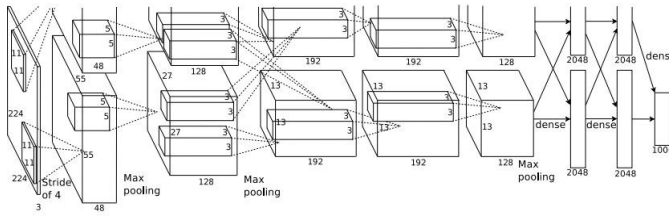
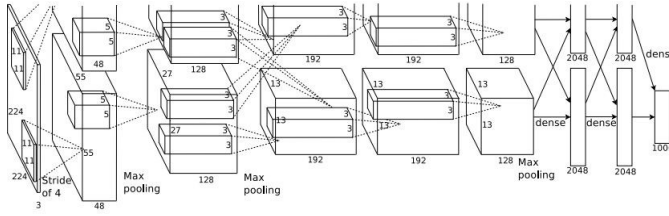
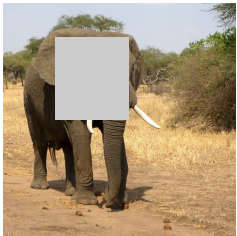
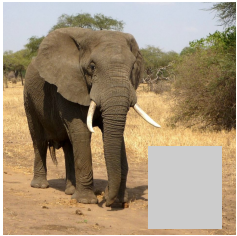
$P(\text{elephant}) = 0.95$



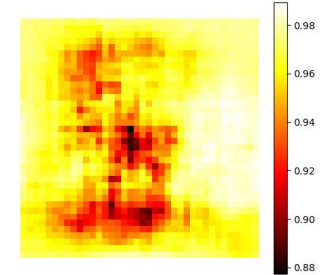
$P(\text{elephant}) = 0.75$

Which pixels matter: Saliency vs Occlusion

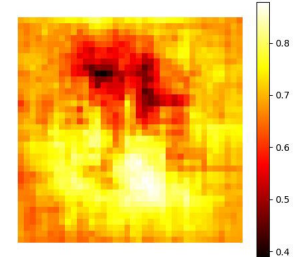
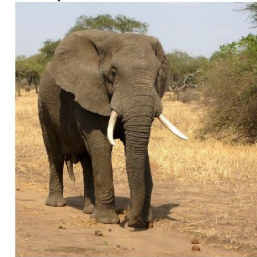
Mask part of the image before feeding to CNN, check how much predicted probabilities change



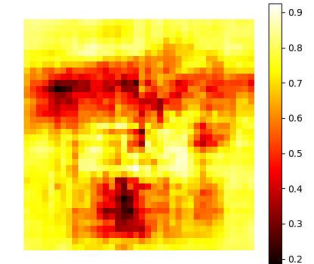
schooner



African elephant, *Loxodonta africana*

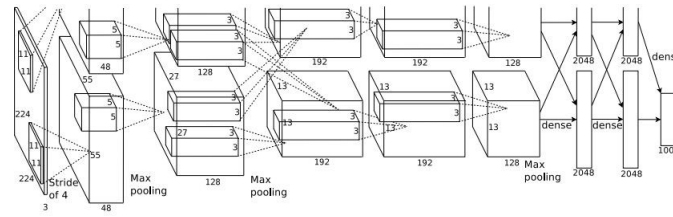


go-kart



Which pixels matter: Saliency via Backprop

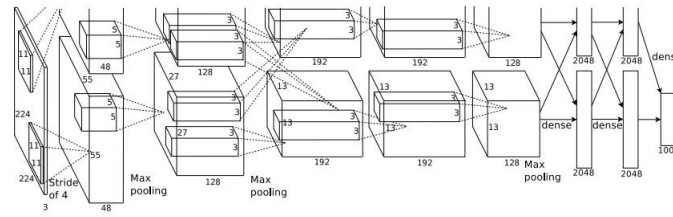
Forward pass: Compute probabilities



Dog

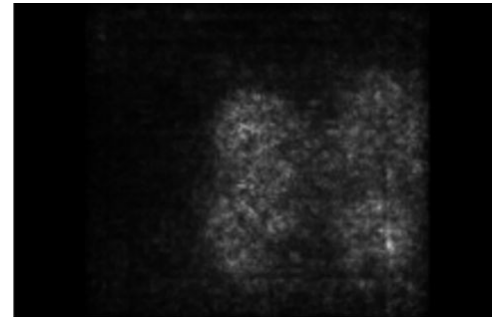
Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



Dog

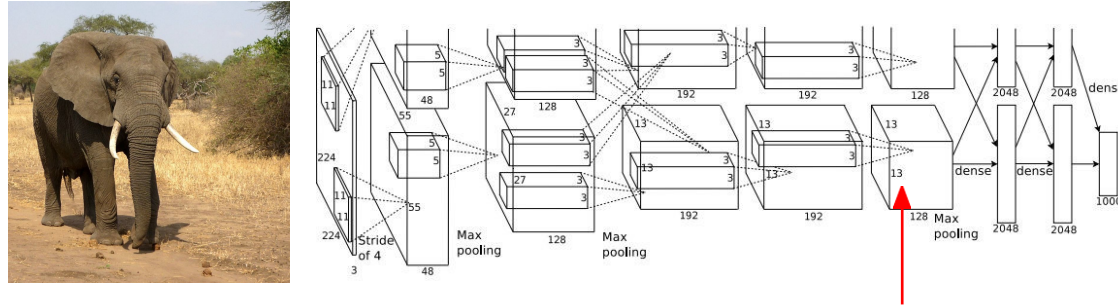
Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels



Saliency Maps



Intermediate Features via (guided) backprop



Pick a single intermediate neuron, e.g. one value in 128 x 13 x 13 conv5 feature map

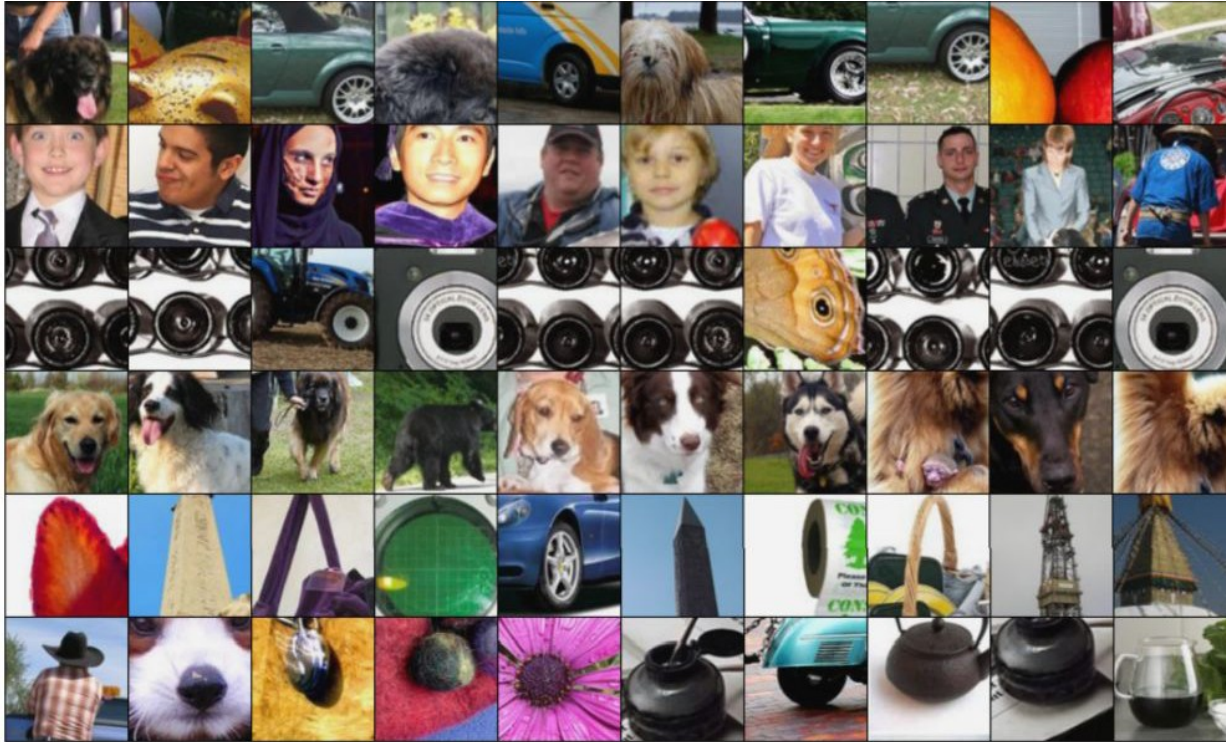
Compute gradient of neuron value with respect to image pixels

Intermediate features via (guided) backprop



Maximally activating patches
(Each row is a different neuron)

Intermediate features via (guided) backprop



Maximally activating patches
(Each row is a different neuron)

Fooling Images / Adversarial Examples

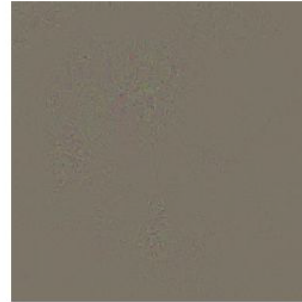
African elephant



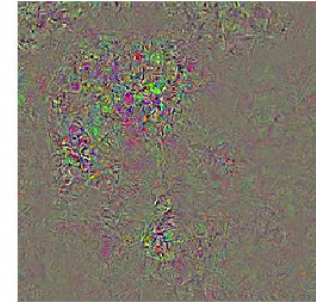
koala



Difference



10x Difference



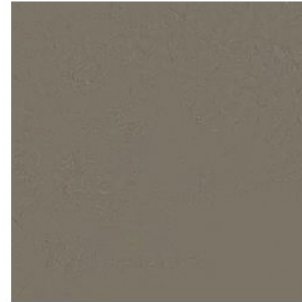
schooner



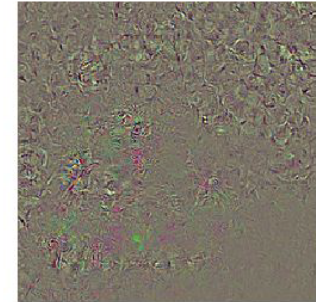
iPod



Difference



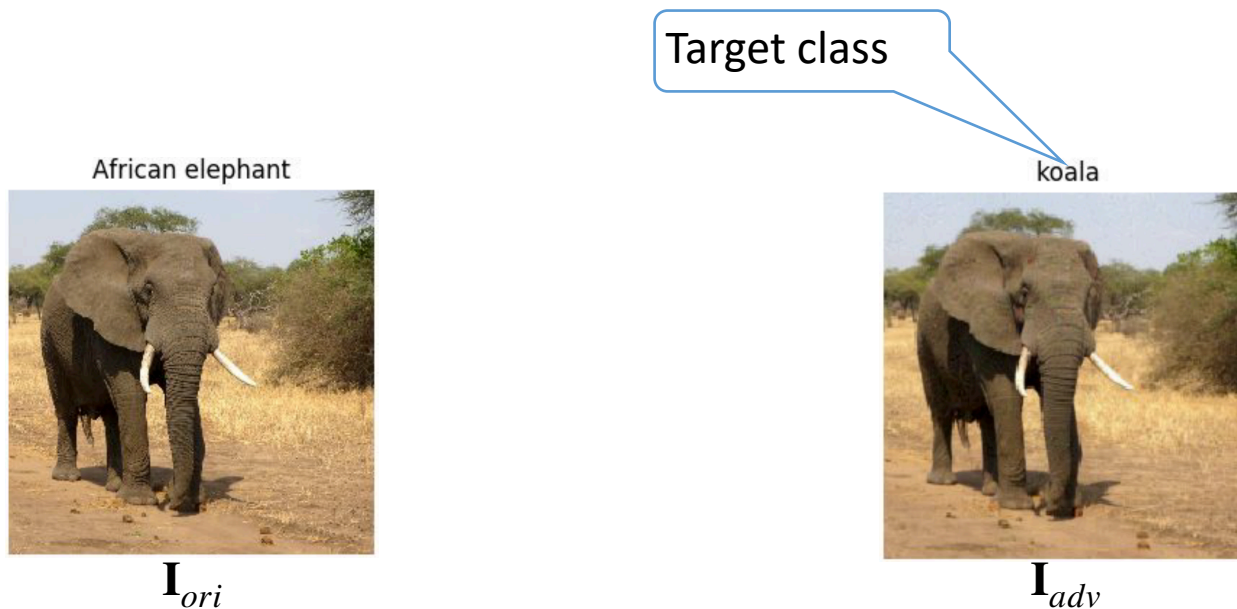
10x Difference



Fooling Images / Adversarial Examples

- (1) Start from an arbitrary image
- (2) Pick an arbitrary class
- (3) Modify the image to maximize the class
- (4) Repeat until network is fooled

Optimization Formulation



$Score_c(\mathbf{I}; \theta)$: the confidence score of an image belonging to class c , using a network of parameters θ

Attack: Modify the image \mathbf{I} to increase $Score_{target\ class}(\mathbf{I}; \theta)$

$$\begin{array}{ll} \underset{\mathbf{I}_{adv}}{\text{maximize}} & Score_{target\ class}(\mathbf{I}_{adv}) \\ \text{subject to} & \|\mathbf{I}_{adv} - \mathbf{I}_{ori}\| \leq \epsilon \end{array}$$

Gradient-based Attack

Fast Gradient Sign Method:

$$\mathbf{I}^{adv} = \mathbf{I} + \epsilon \text{sign}(\nabla_{\mathbf{I}} \text{Score}_{target\ class}(\mathbf{I}))$$



African elephant



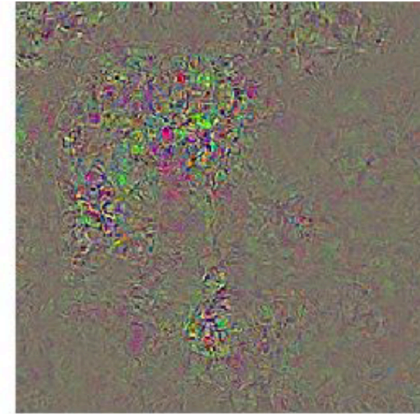
koala



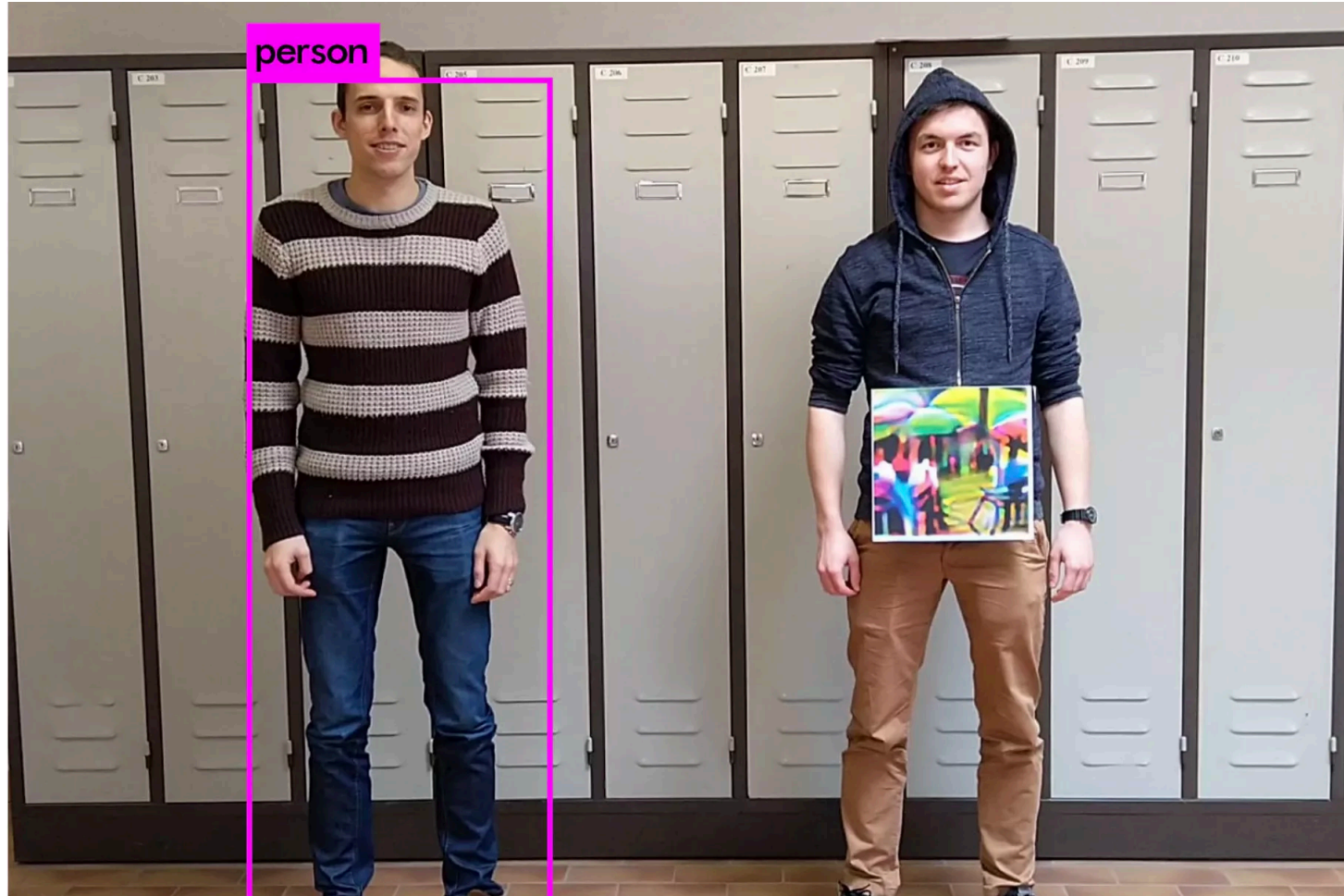
Difference



10x Difference



Patch-based Attack (Spatially Localized)



Dangerous!



(a)



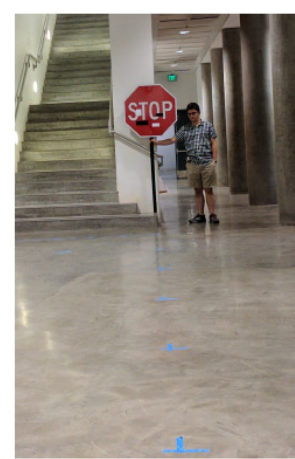
(b)



(c)



(d)



(e)

Neural Style Transfer

Content Image



This image is licensed under [CC-BY3.0](#)

+

Style Image



Starry Night by Van Gogh is in the public domain

Neural Style Transfer

Content Image



[This image](#) is licensed under [CC-BY3.0](#)

Style Image



[Starry Night](#) by Van Gogh is in the public domain

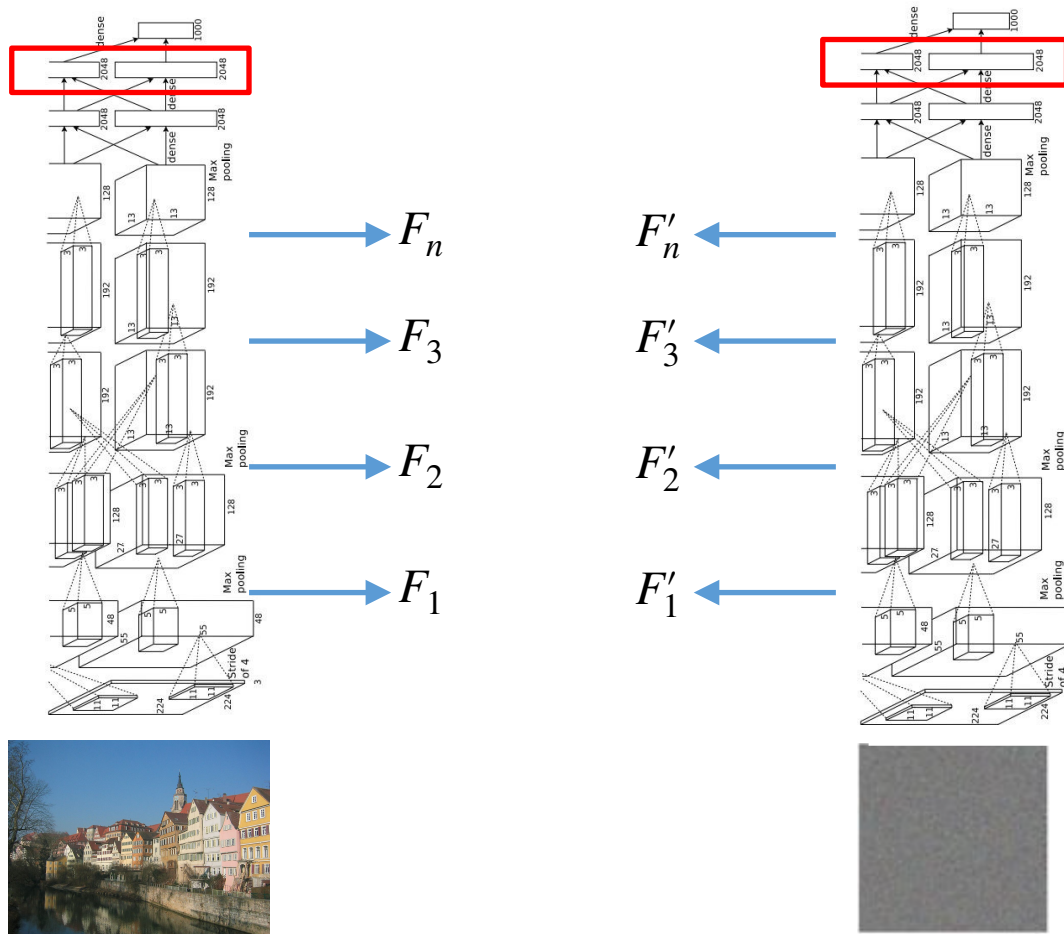
Style Transfer!



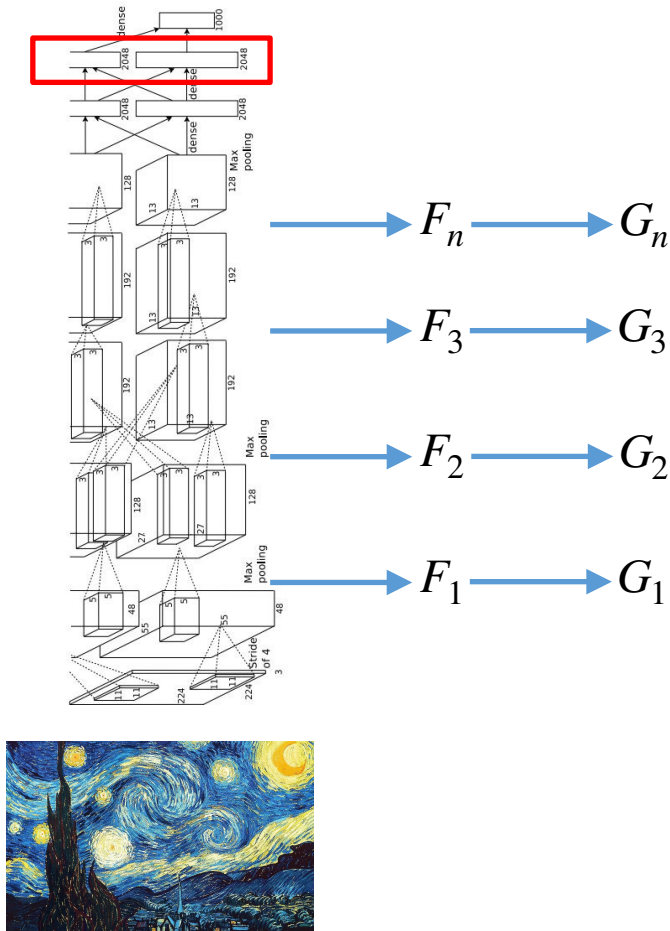
[This image](#) copyright Justin Johnson, 2015. Reproduced with permission.

Content Discrepancy Loss

$$\mathcal{L}_{content}(\mathbf{I}) = \|\mathbf{F}(\mathbf{I}_{ref}) - \mathbf{F}(\mathbf{I})\|^2$$

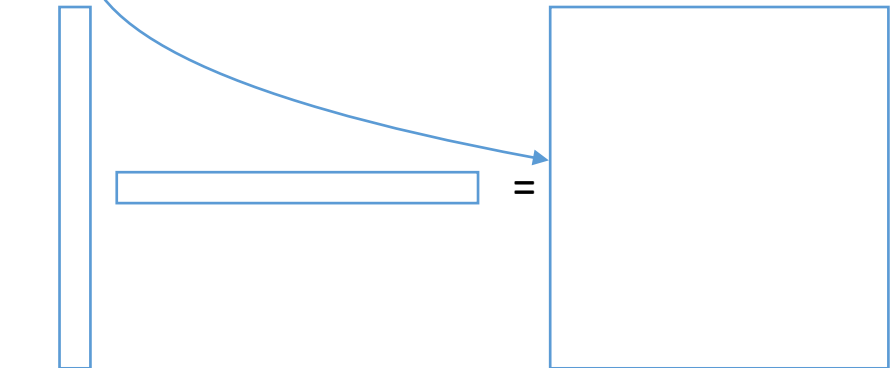


Style Discrepancy Loss

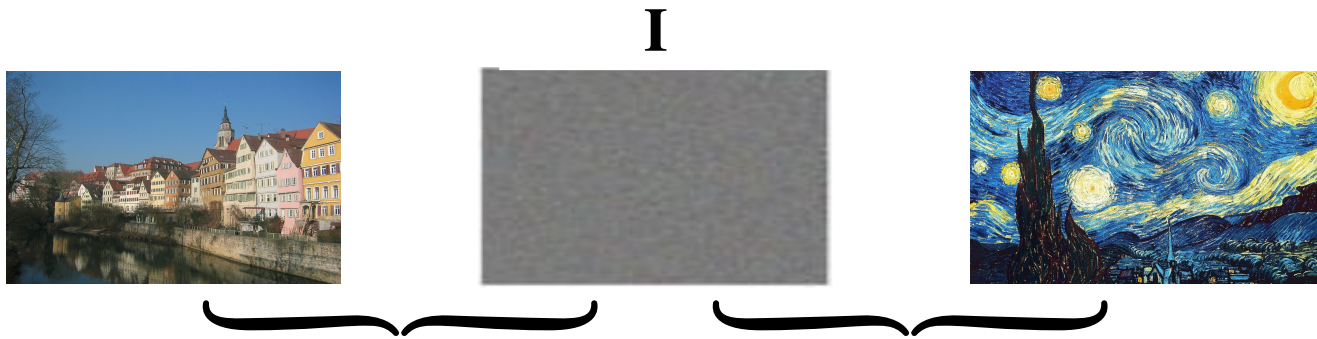


F : the centered feature tensor at some layer
(as in PCA)
G : covariance matrix, a.k.a Gram matrix
(as in PCA)

$F = F.flatten().reshape(-1, 1)$
 $G = F @ F.T$

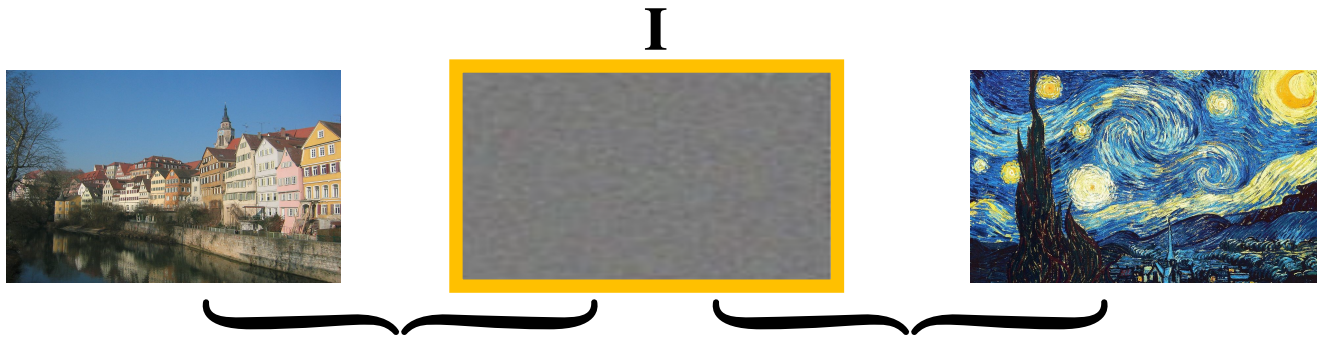


Total Loss



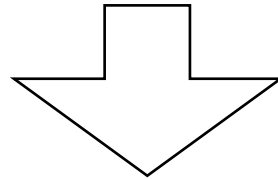
$$\mathcal{L}_{total}(\mathbf{I}) = \alpha \mathcal{L}_{content}(\mathbf{I}) + \beta \mathcal{L}_{style}(\mathbf{I})$$

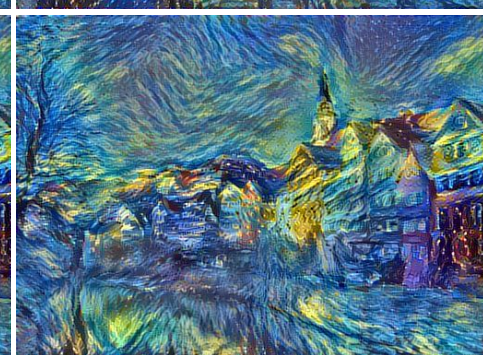
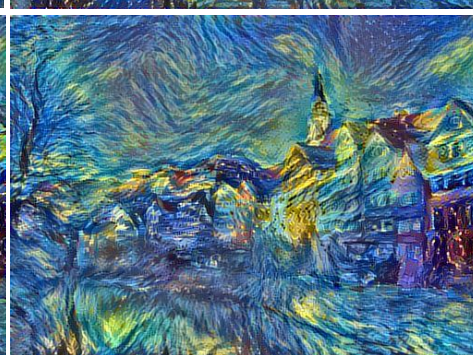
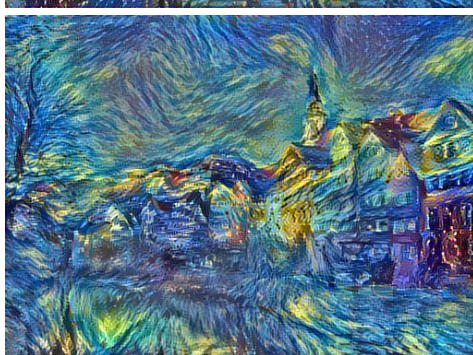
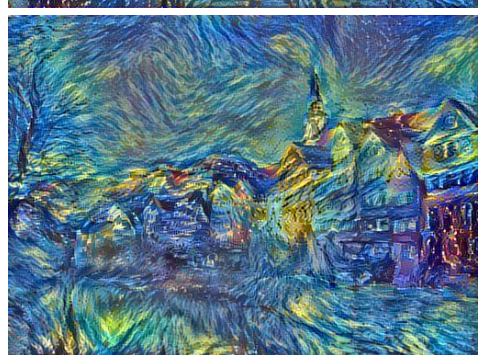
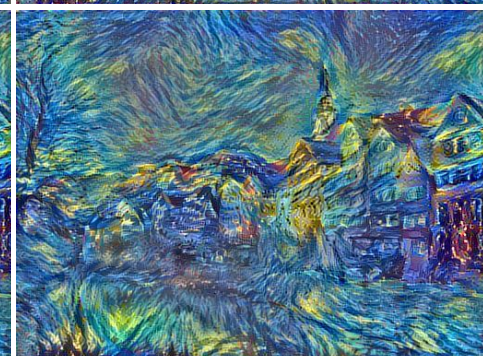
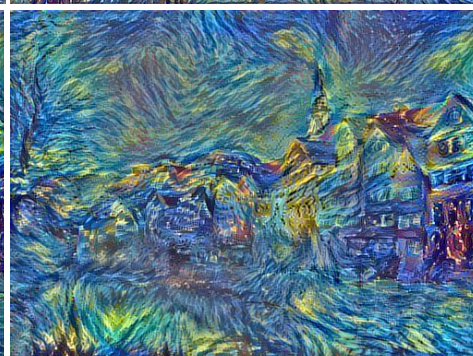
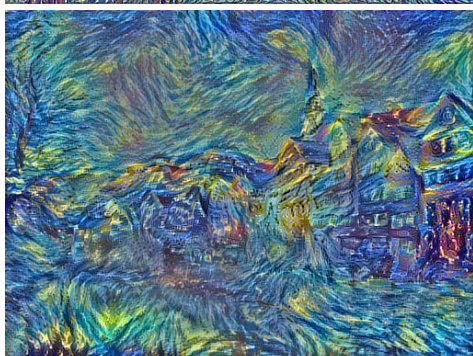
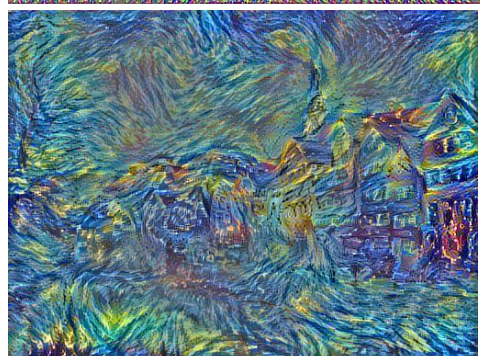
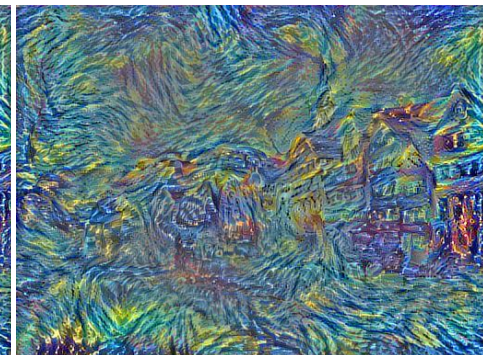
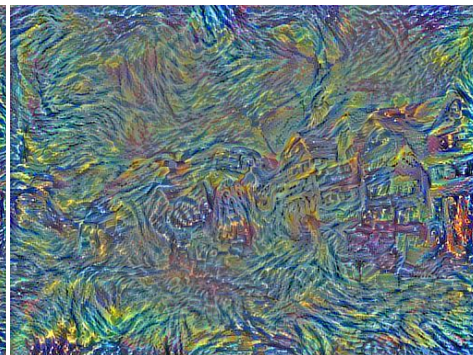
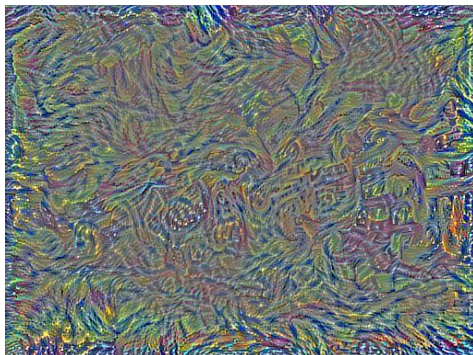
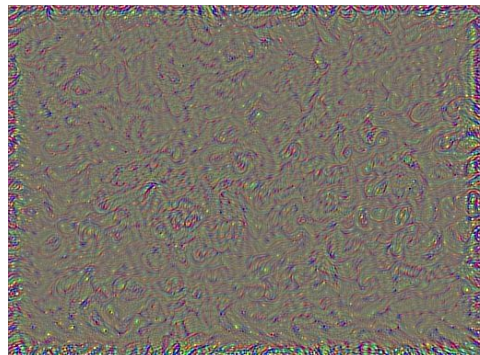
Total Loss



$$\mathcal{L}_{total}(\mathbf{I}) = \alpha \mathcal{L}_{content}(\mathbf{I}) + \beta \mathcal{L}_{style}(\mathbf{I})$$

Minimize total loss





Neural Style Transfer



Example outputs

Neural Style Transfer



More weight to
content loss



More weight to
style loss

Neural Style Transfer

Resizing style image before running style transfer algorithm can transfer different types of features



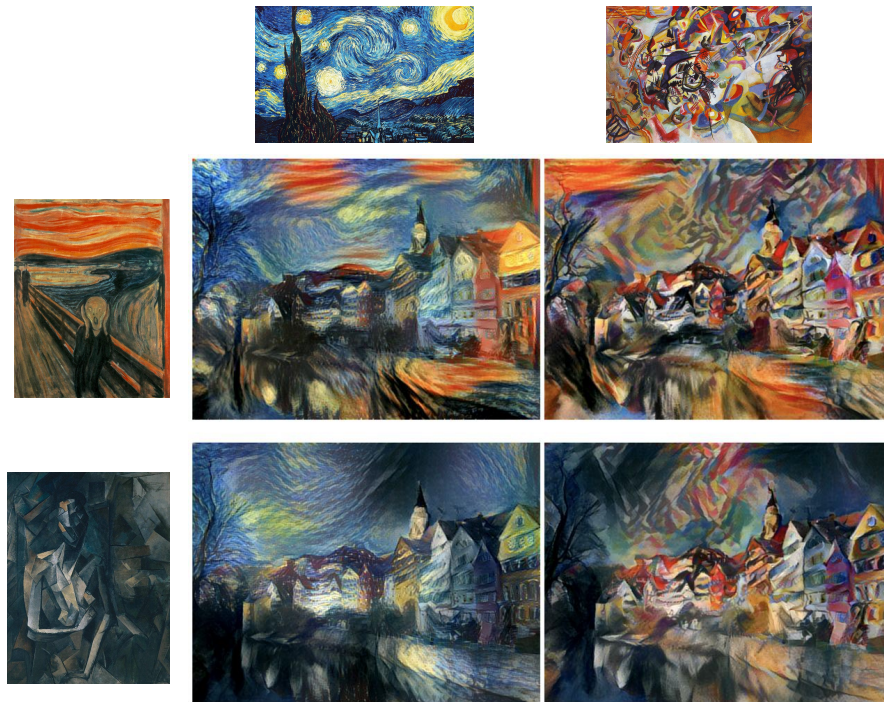
Larger style
image

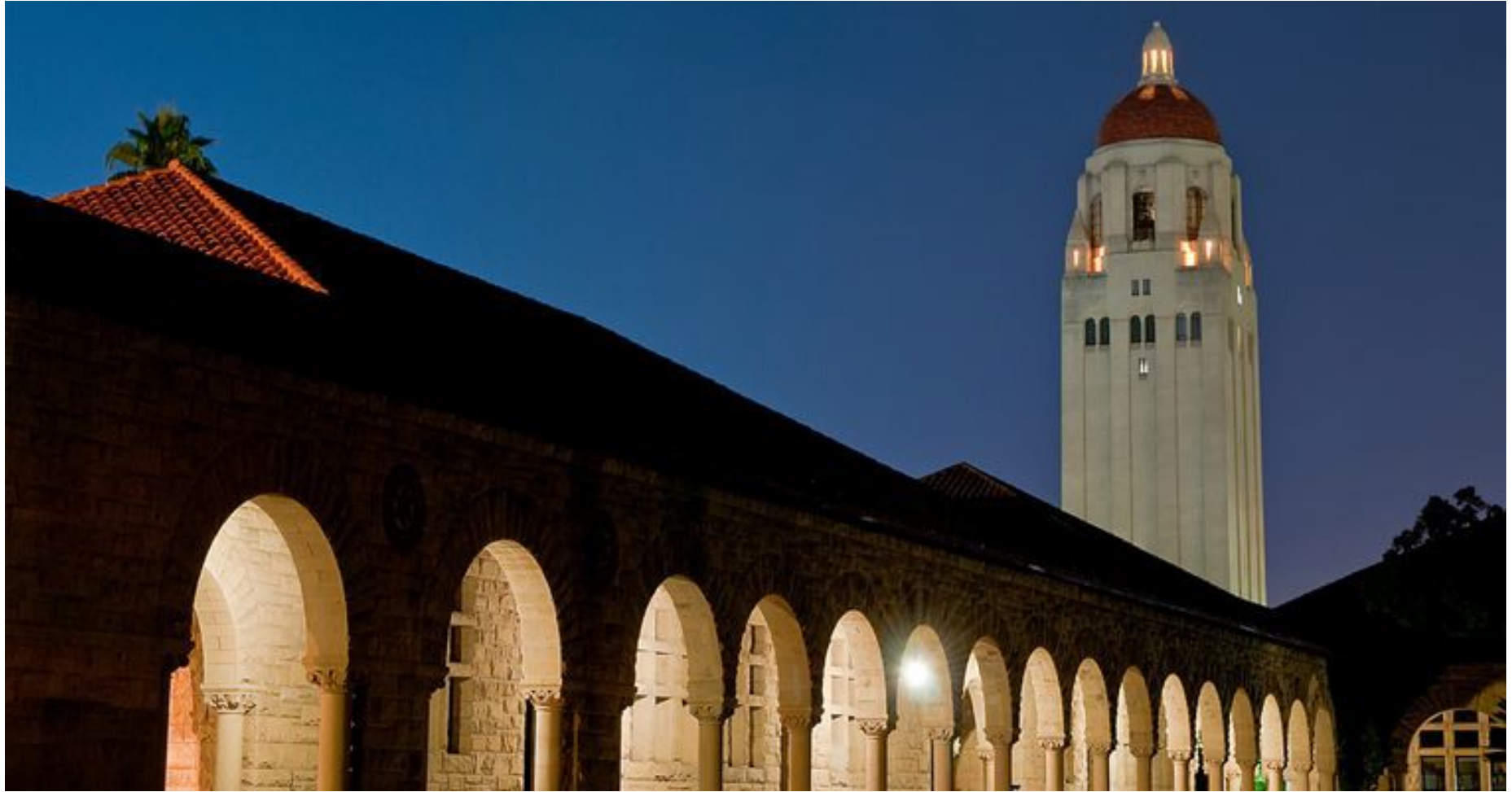


Smaller style
image

Neural Style Transfer: Multiple Style Images

Mix style from multiple images by taking a weighted average of Gram matrices





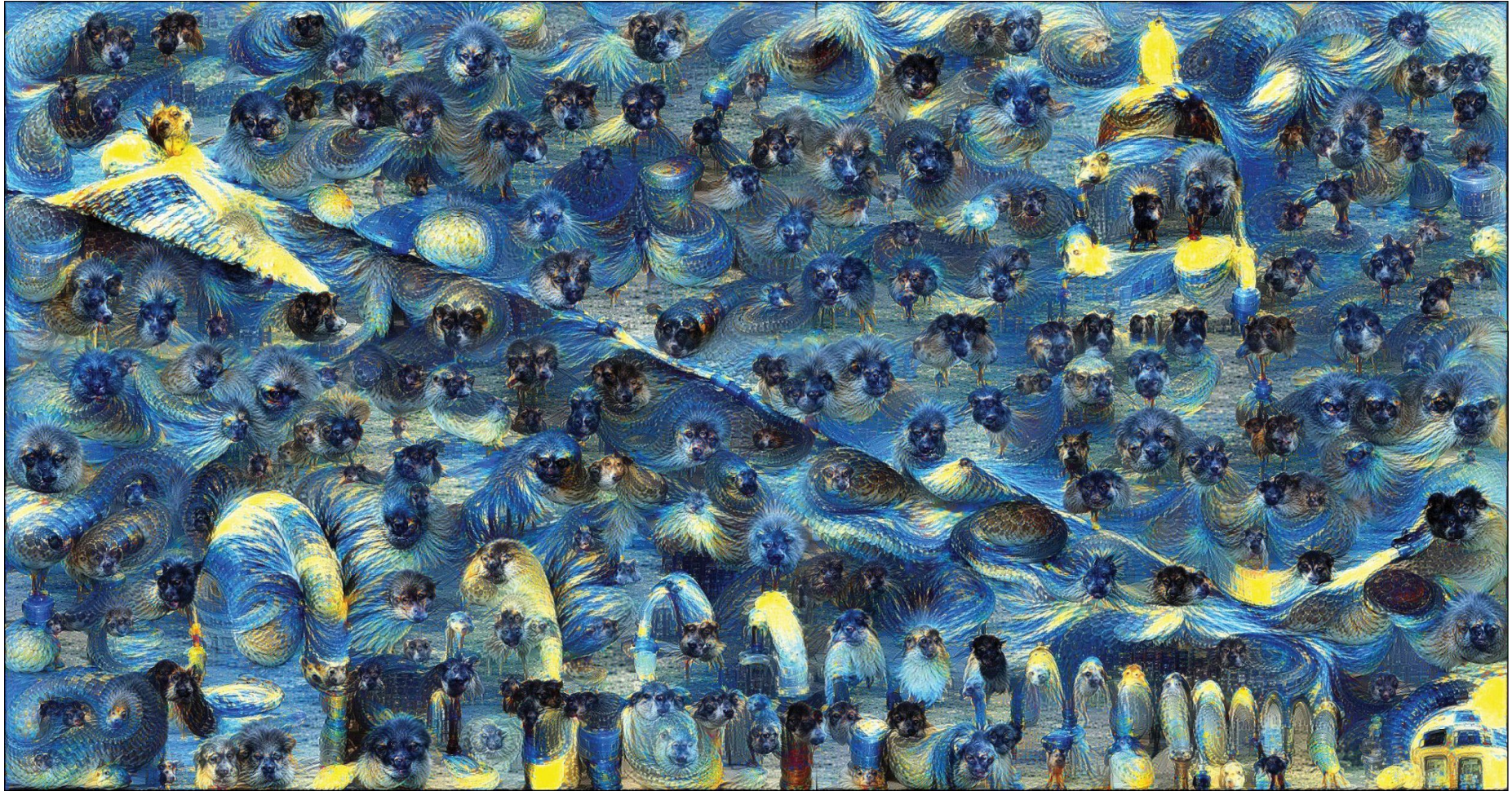
<https://github.com/jcjohnson/fast-neural-style>



<https://github.com/jcjohnson/fast-neural-style>



<https://github.com/jcjohnson/fast-neural-style>



<https://github.com/jcjohnson/fast-neural-style>