

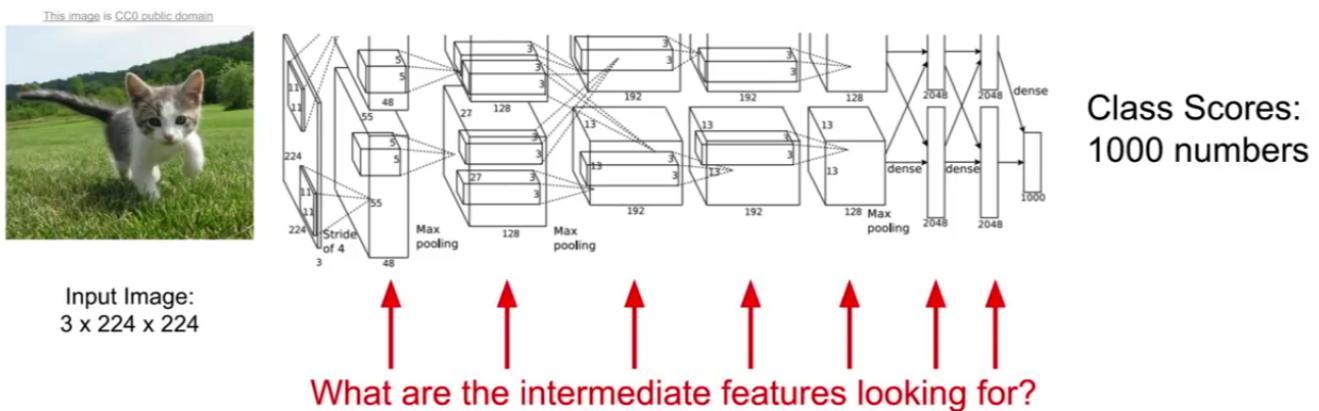
lec12

Creation Date: 15/01/2020 22:24

Last Modified Date: 16/01/2020 00:38

Lec 12: Visualizations and Understanding CNNs

What's going on inside ConvNets?



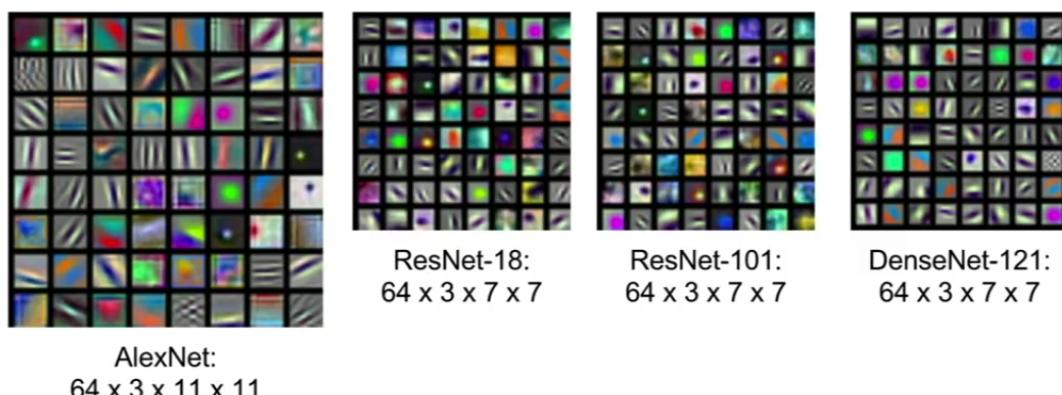
Krizhevsky et al., "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.
Figure reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 -

University, May 10, 2017

First Layer: Visualize Filters

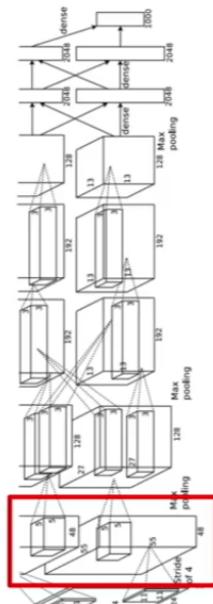


Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014
He et al, "Deep Residual Learning for Image Recognition", CVPR 2016
Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 -

University, May 10, 2017



- These are the convolutional filter weights represented as images.
- We scale the values for 0-255 range for visualizing.

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

$16 \times 3 \times 7 \times 7$

Weights:
()

layer 2 weights

$20 \times 16 \times 7 \times 7$

Weights:
()

layer 3 weights

$20 \times 20 \times 7 \times 7$

Fei-Fei Li & Justin Johnson & Serena Yeung

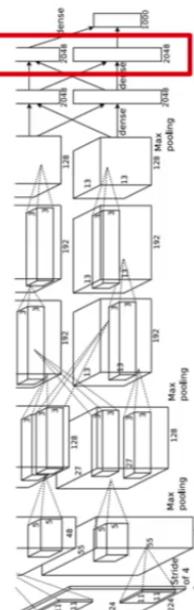
Lecture 11 -

University, May 10, 2017

- The filters after the first Conv layer will have channel depth more than 3 which makes it difficult to visualize as RGB image. So we can visualize them into list of grayscale images for each filter. They have some kind of spatial structure but there is no much interpretability.
- This is caused because the second layer activation are dependent on the activation of first layer and the second layers' activation are highest depending on the activation of first layers' output.
- We can notice that there is no much interpretability from this.

Last Layer

FC7 layer



4096-dimensional feature vector for an image
(layer immediately before the classifier)

Run the network on many images, collect the feature vectors

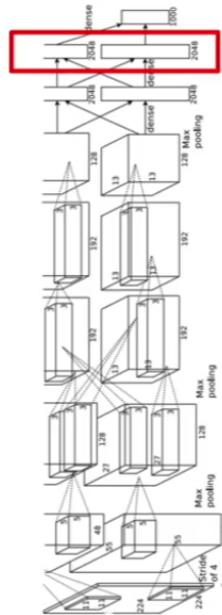
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 -

University, May 10, 2017

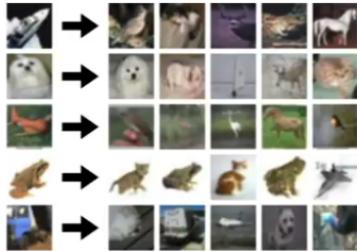
Last Layer: Nearest Neighbors

4096-dim vector



Test image L2 Nearest neighbors in feature space

Recall: Nearest neighbors in pixel space



Krizhevsky et al., "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.
Figures reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 -

University of California, Berkeley
May 10, 2017

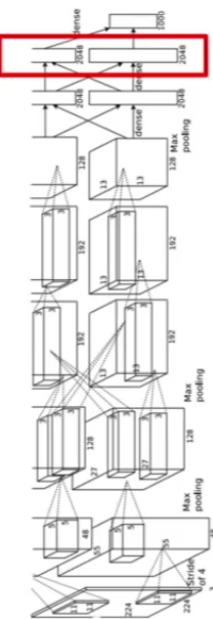
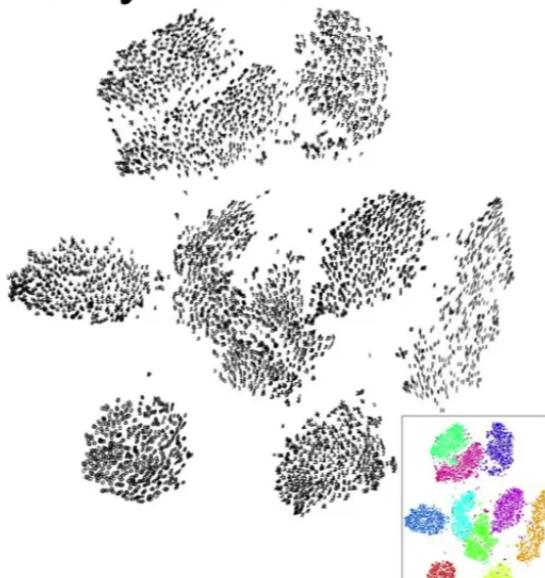
- If we just find the nearest neighbour in pixel space, we may find some images which are similar but there can be many more matches which are FP in nature.
- Finding nearest neighbour in feature space yields better neighbours. Because it has both spatial and semantic information present in the feature space.
- **Even though the CNN training did not model the images to be nearer to each other they become nearer**
- Various other approaches are used to explicitly model the loss function to have clusters. Eg. Contrastive Loss, Triplet Loss
- **For each image, we pass it through the network and record the feature vector of 4096 length. To find the nearest neighbour we use this vector to find the L2 distance.**

Last Layer: Dimensionality Reduction

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

Simple algorithm: Principle Component Analysis (PCA)

More complex: t-SNE



Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008

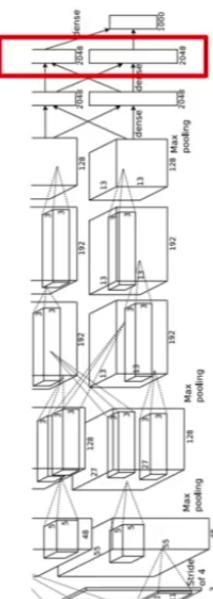
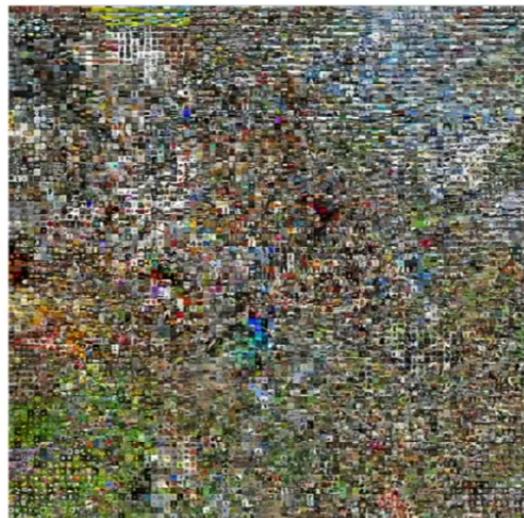
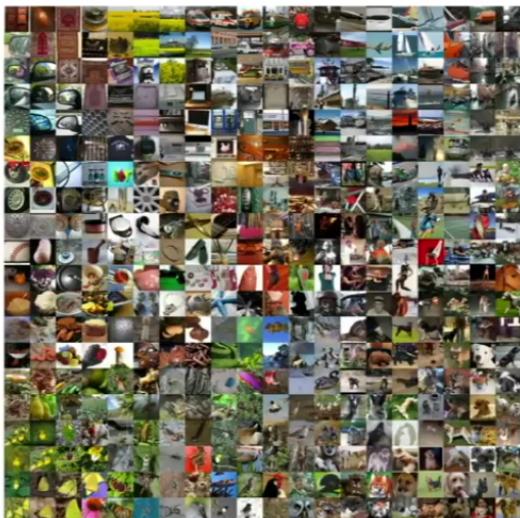
Figure copyright Laurens van der Maaten and Geoff Hinton. 2008. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 9 University, May 10, 2017

- Other method is to collect all the vectors which will be Nx4096, and then apply PCA, to reduce it to 2 or 3 dimensions.

Last Layer: Dimensionality Reduction



Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008

Krizhevsky et al., "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.

Figure reproduced with permission.

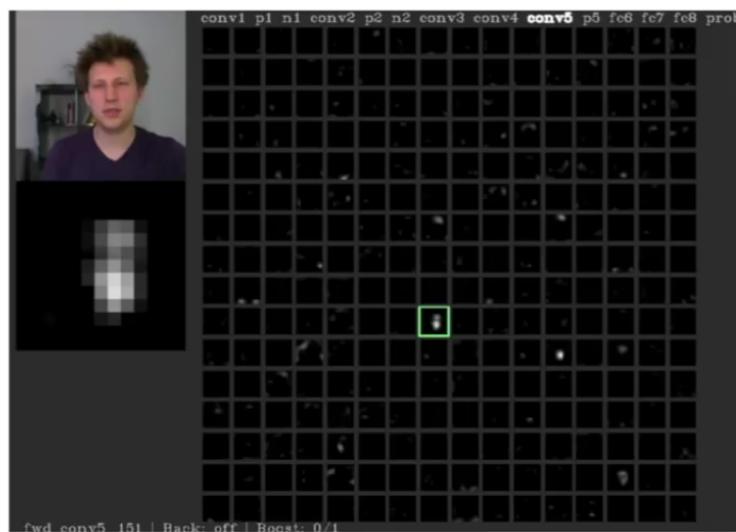
Fei-Fei Li & Justin Johnson & Serena Yeung

See high-resolution versions at
<http://cs.stanford.edu/people/karpathy/cnnembed/>

Lecture 11 - 10 University, May 10, 2017

Visualizing Activations

conv5 feature map is
128x13x13; visualize
as 128 13x13
grayscale images

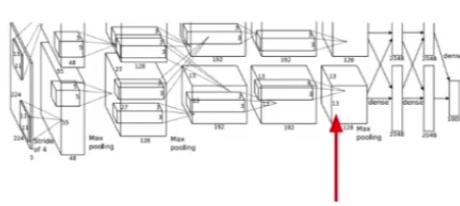


Yosinski et al., "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.
Figure copyright Jason Yosinski, 2014. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 11 University May 10, 2017

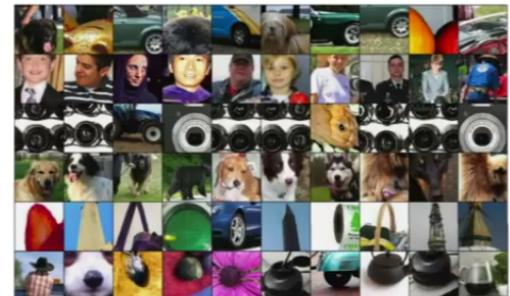
Maximally Activating Patches



Pick a layer and a channel; e.g. conv5 is
128 x 13 x 13, pick channel 17/128

Run many images through the network,
record values of chosen channel

Visualize image patches that correspond to maximal activations



Springenberg et al., "Striving for Simplicity: The All-Convolutional Net," ICLR Workshop 2015.
Figure copyright Jost Tobias Springenberg, Alexey Sutskever, Tim Salimbeni, Martin Riedmiller, 2015; reproduced with permission.

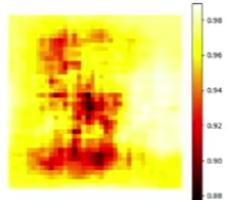
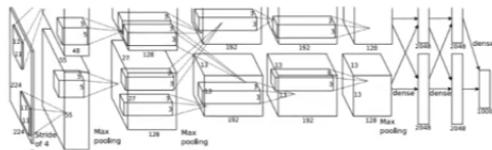
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 12 May 10, 2017

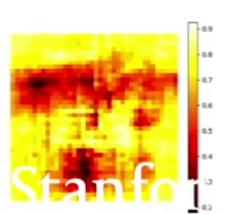
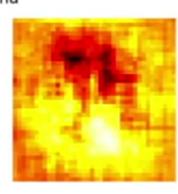
- The bottom grid of images are of deeper layer so it has large receptive fields to it has faces etc.

Occlusion Experiments

Mask part of the image before feeding to CNN, draw heatmap of probability at each mask location



African elephant, Loxodonta africana



Boat image is CC0 public domain
Elephant image is CC0 public domain
Go-Karts image is CC0 public domain

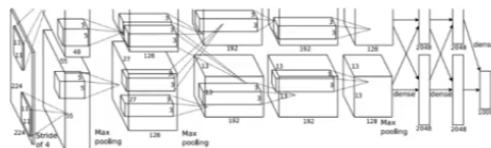
Boat image is CC0 public domain
Elephant image is CC0 public domain
Go-Karts image is CC0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

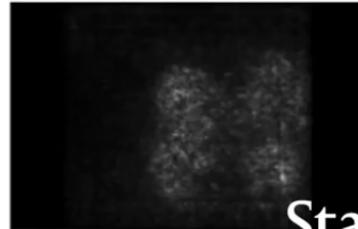
Lecture 11 - 13 University of Stanford, May 10, 2017

Saliency Maps

How to tell which pixels matter for classification?



Dog



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

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 15 University of Stanford, May 10, 2017

Saliency Maps



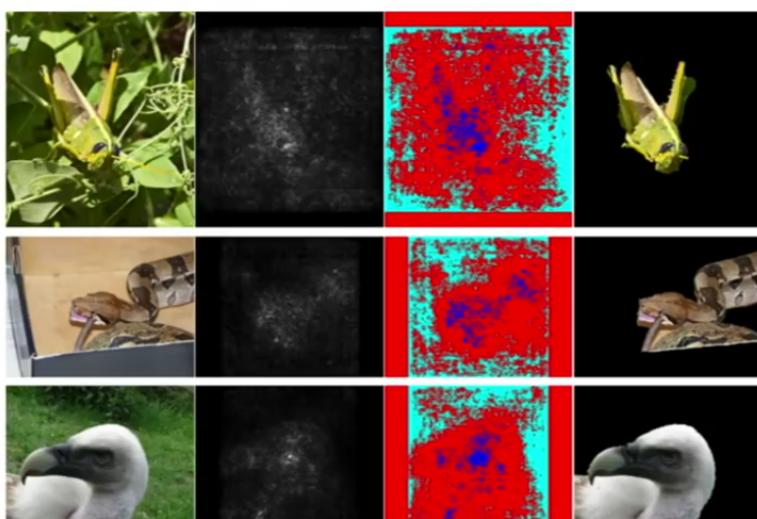
Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 16 University May 10, 2017

Saliency Maps: Segmentation without supervision

Use GrabCut on
saliency map

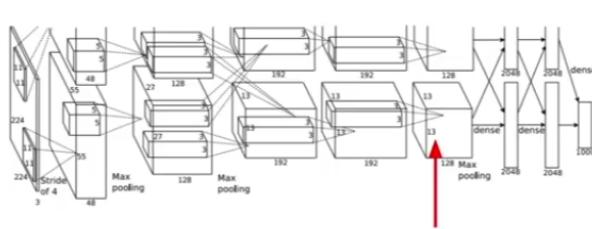


Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.
Rother et al., "Grabcut: Interactive foreground extraction using iterated graph cuts", ACM TOG 2004

Fei-Fei Li & Justin Johnson & Serena Yeung

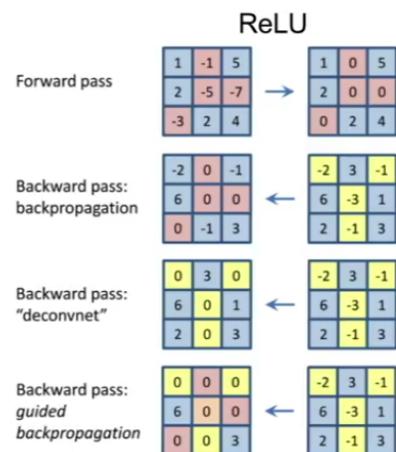
Lecture 11 - 17 University May 10, 2017

Intermediate features via (guided) backprop



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

Compute gradient of neuron value with respect to image pixels



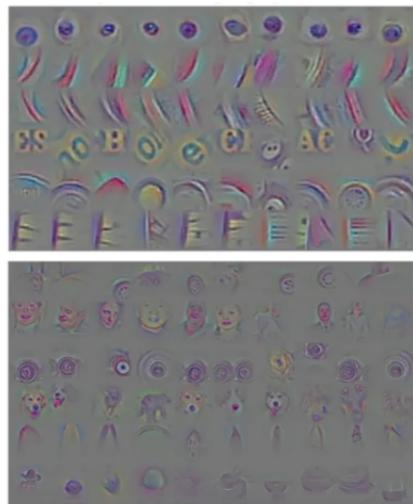
Images come out nicer if you only backprop positive gradients through each ReLU (guided backprop)

Figure copyright Jost Tobias Springenberg, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 19 University May 10, 2017

Intermediate features via (guided) backprop



Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 20 University May 10, 2017

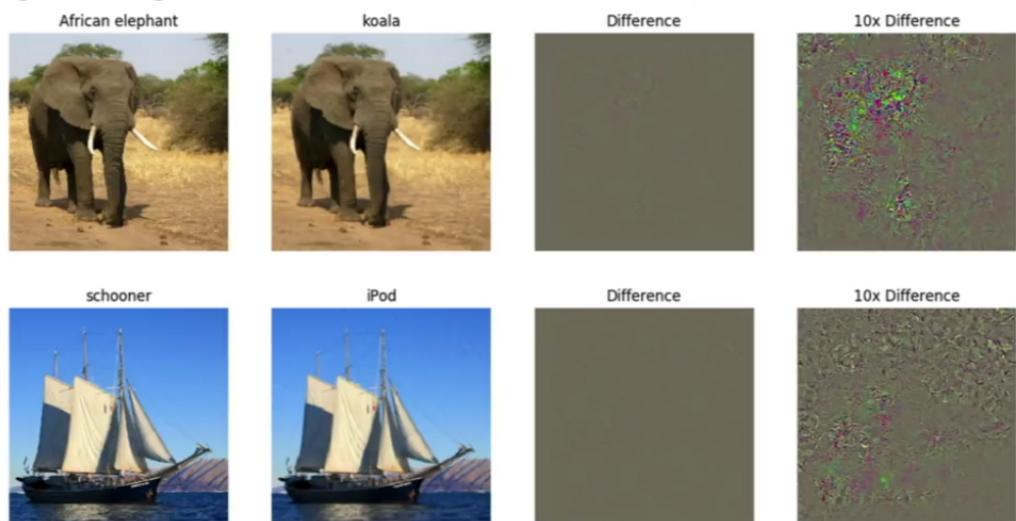
Fooling Images / Adversarial Examples

- (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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 34 University, May 10, 2017

Fooling Images / Adversarial Examples



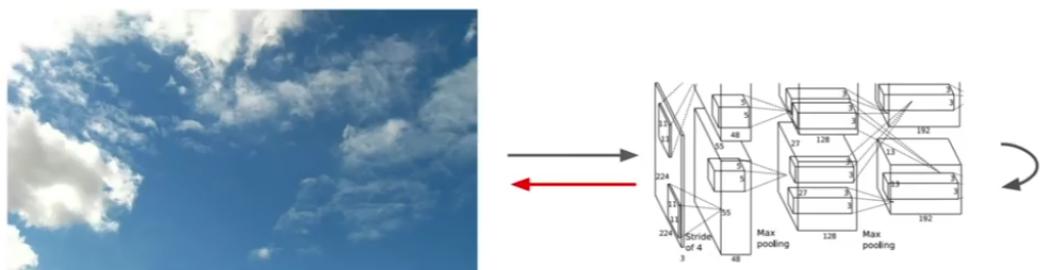
Boat image is CC0 public domain
Elephant image is CC0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 35 University, May 10, 2017

DeepDream: Amplify existing features

Rather than synthesizing an image to maximize a specific neuron, instead try to **amplify** the neuron activations at some layer in the network



Choose an image and a layer in a CNN; repeat:

1. Forward: compute activations at chosen layer
2. Set gradient of chosen layer *equal to its activation*
3. Backward: Compute gradient on image
4. Update image

Equivalent to:

$$I^* = \arg \max_I \sum_i f_i(I)^2$$

Mordvintsev, Olah, and Tygina, "Inceptionism: Seeing Depth in Neural Networks", Google Research Blog, images are licensed under CC-BY

Feature Inversion

Given a CNN feature vector for an image, find a new image that:

- Matches the given feature vector
- “looks natural” (image prior regularization)

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Given feature vector

Features of new image

$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

$$\mathcal{R}_{V^\beta}(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{i,j})^2 + (x_{i+1,j} - x_{i,j})^2 \right)^{\frac{\beta}{2}}$$

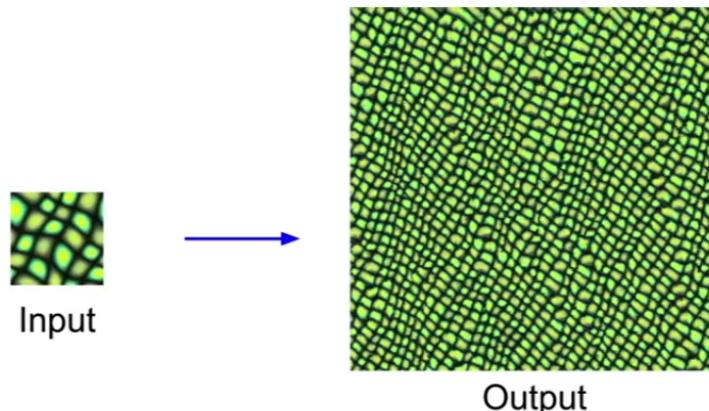
Total Variation regularizer
(encourages spatial smoothness)

Mahendran and Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR 2015

- Texture Synthesis

Texture Synthesis

Given a sample patch of some texture, can we generate a bigger image of the same texture?



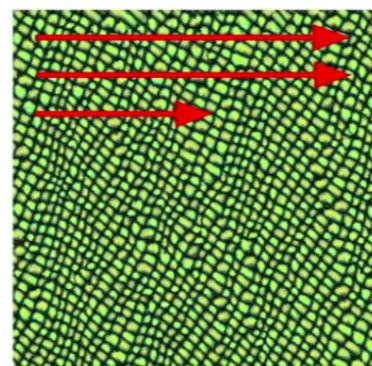
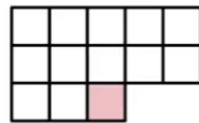
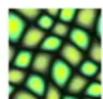
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 51 University of California, Berkeley, May 10, 2017

MIT license

Texture Synthesis: Nearest Neighbor

Generate pixels one at a time in scanline order; form neighborhood of already generated pixels and copy nearest neighbor from input

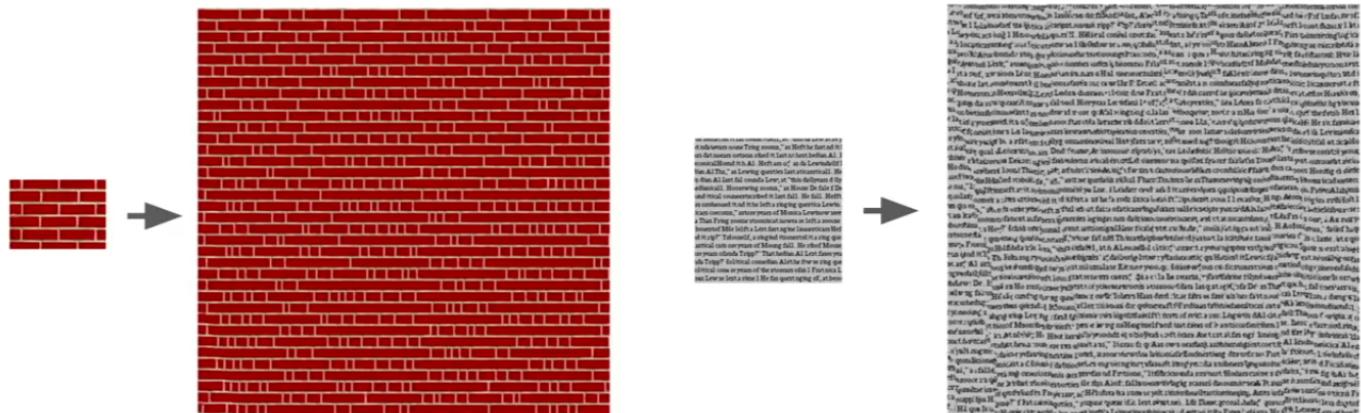


Wei and Levoy, "Fast Texture Synthesis using Tree-structured Vector Quantization", SIGGRAPH 2000
Efros and Leung, "Texture Synthesis by Non-parametric Sampling", ICCV 1999

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 52 University of California, Berkeley, May 10, 2017

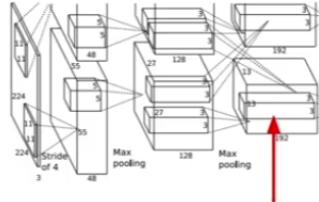
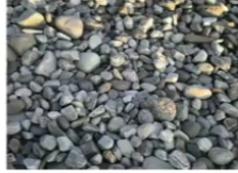
Texture Synthesis: Nearest Neighbor



Fei-Fei Li & Justin Johnson & Serena Yeung

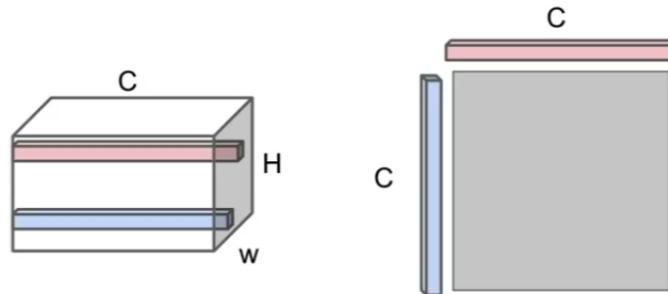
Lecture 11 - 53 University of Texas at Austin, May 10, 2017

Neural Texture Synthesis: Gram Matrix



Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

Outer product of two C -dimensional vectors gives $C \times C$ matrix measuring co-occurrence



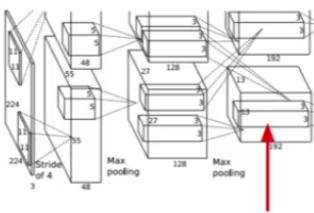
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 55 University of Texas at Austin, May 10, 2017

Neural Texture Synthesis: Gram Matrix



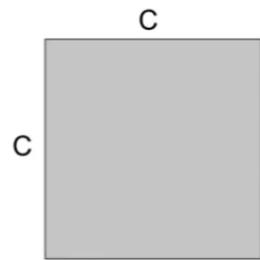
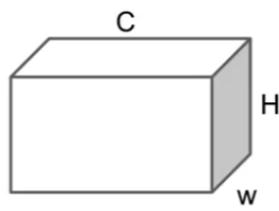
This image is in the public domain.



Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

Outer product of two C -dimensional vectors gives $C \times C$ matrix measuring co-occurrence

Average over all HW pairs of vectors, giving **Gram matrix** of shape $C \times C$



Efficient to compute; reshape features from

$C \times H \times W$ to $=C \times HW$

then compute $G = FF^T$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 57 University, May 10, 2017

- We don't use Covariance matrix to compute because it is computationally expensive.

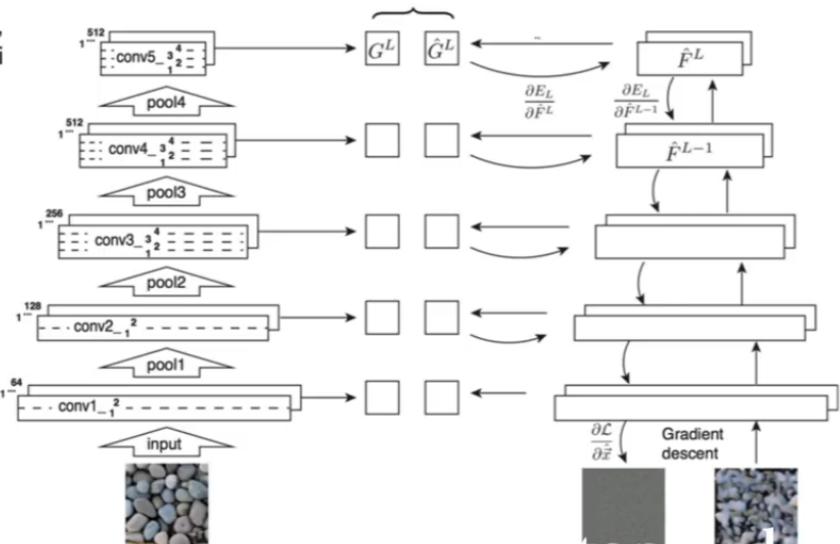
Neural Texture Synthesis

- Pretrain a CNN on ImageNet (VGG-19)
- Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
- At each layer compute the *Gram matrix* giving outer product of features:

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l \text{ (shape } C_i \times C_i\text{)}$$

- Initialize generated image from random noise
- Pass generated image through CNN, compute Gram matrix on each layer
- Compute loss: weighted sum of L2 distance between Gram matrices
- Backprop to get gradient on image
- Make gradient step on image
- GOTO 5

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



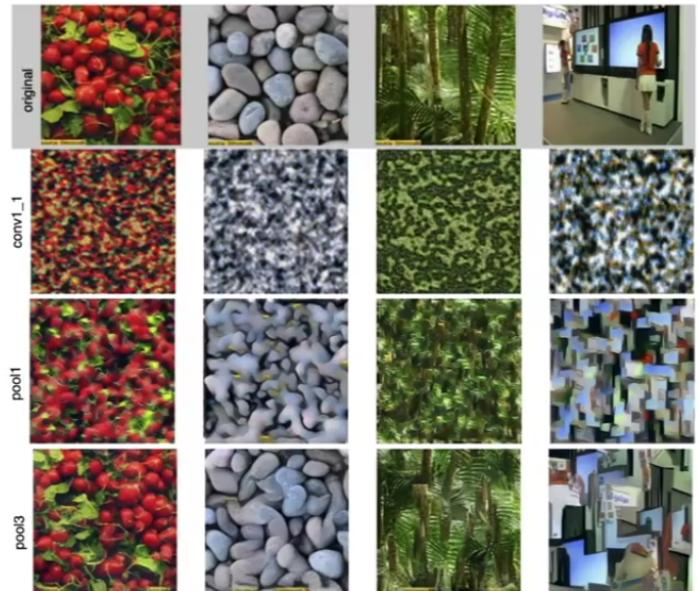
Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015
Figure copyright Leon Gatys, Alexander S. Ecker, and Matthias Bethge, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 61 University, May 10, 2017

Neural Texture Synthesis

Reconstructing texture from higher layers recovers larger features from the input texture



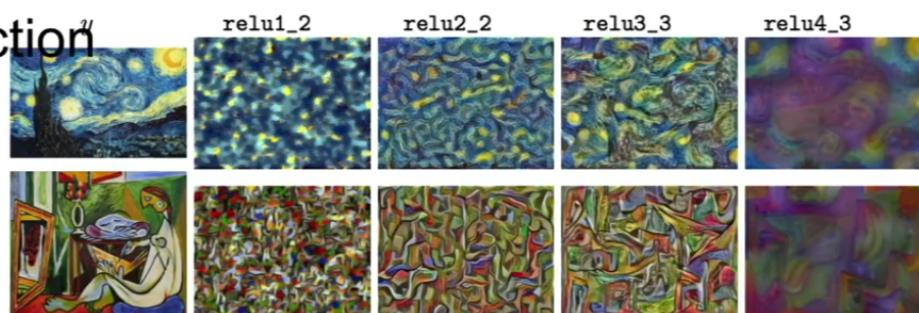
Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015
Figure copyright Leon Gatys, Alexander S. Ecker, and Matthias Bethge, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 62 University, May 10, 2017

Neural Style Transfer: Feature + Gram Reconstruction

Texture synthesis
(Gram reconstruction)



Feature reconstruction



Figure from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016.
Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 64 University, May 10, 2017

Stanford

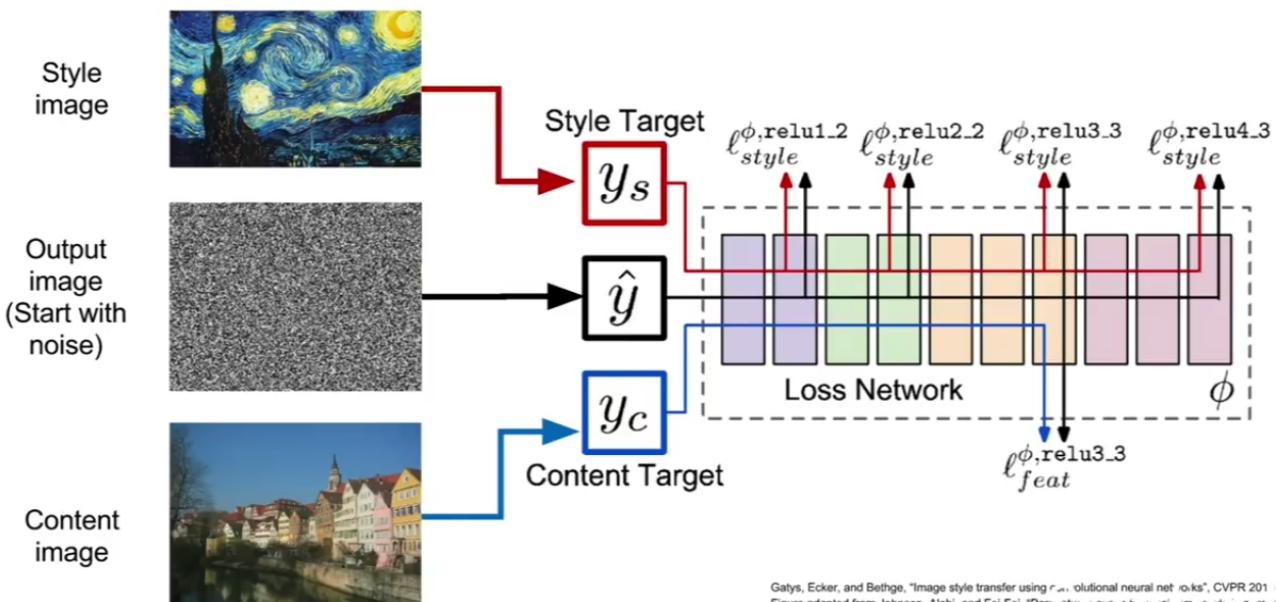
Neural Style Transfer



Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

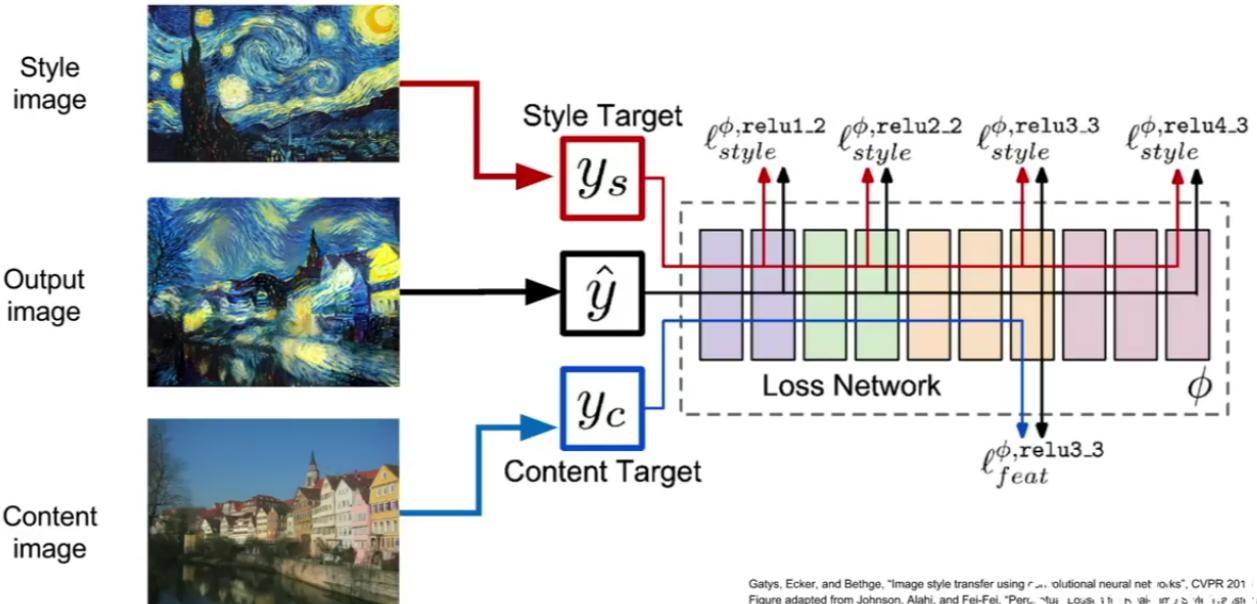
Lecture 11 - 66 University, May 10, 2017



Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016
Figure adapted from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Image Recoloring and Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 67 University, May 10, 2017



Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2015
Figure adapted from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Image Recognition and Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced with permission.

Neural Style Transfer



Neural Style Transfer

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



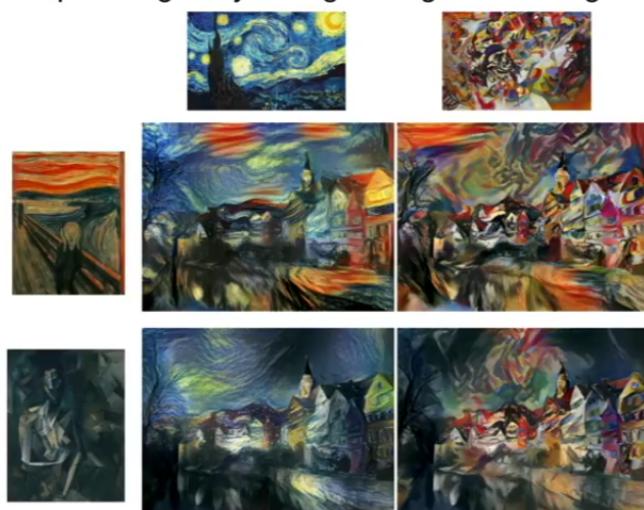
Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016
Figure copyright Justin Johnson, 2015.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 71 University, May 10, 2017

Neural Style Transfer: Multiple Style Images

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



Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016
Figure copyright Justin Johnson, 2015.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 72 University, May 10, 2017

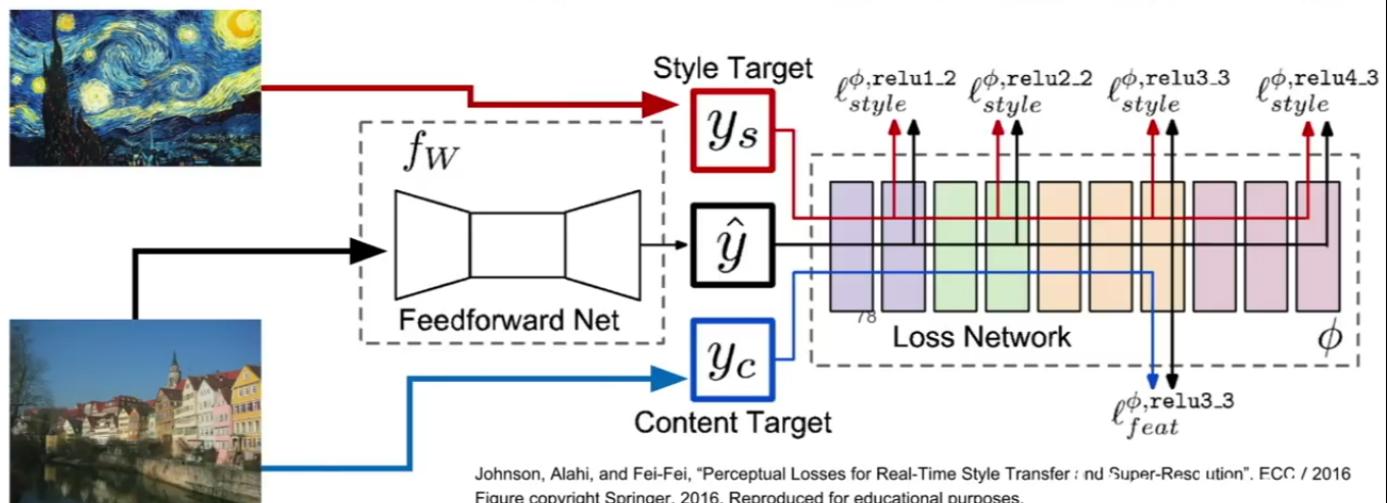
Neural Style Transfer

Problem: Style transfer requires many forward / backward passes through VGG; very slow!

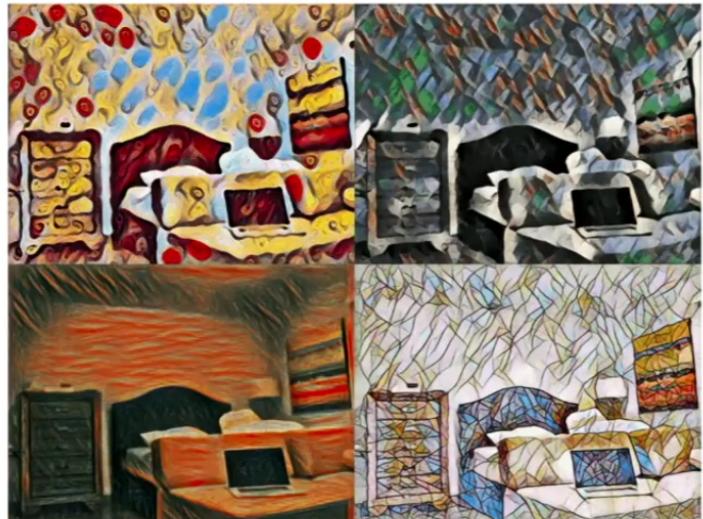
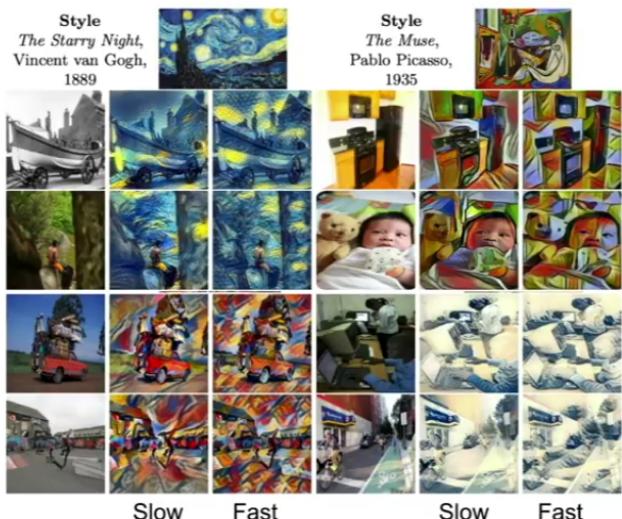
Solution: Train another neural network to perform style transfer for us!

Fast Style Transfer

- (1) Train a feedforward network for each style
- (2) Use pretrained CNN to compute same losses as before
- (3) After training, stylize images using a single forward pass



Fast Style Transfer



Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016
Figure copyright Springer, 2016. Reproduced for educational purposes.

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 79 University, May 10, 2017

Fast Style Transfer



Replacing batch normalization with Instance Normalization improves results

Ulyanov et al., "Texture Networks: Feed-forward Synthesis of Textures and Stylized Images", ICML 2016
Ulyanov et al., "Instance Normalization: The Missing Ingredient for Fast Stylization", arXiv 2016
Figures copyright Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, and Victor Lempitsky, 2016. Reproduced with

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 81 University, May 10, 2017

One Network, Many Styles



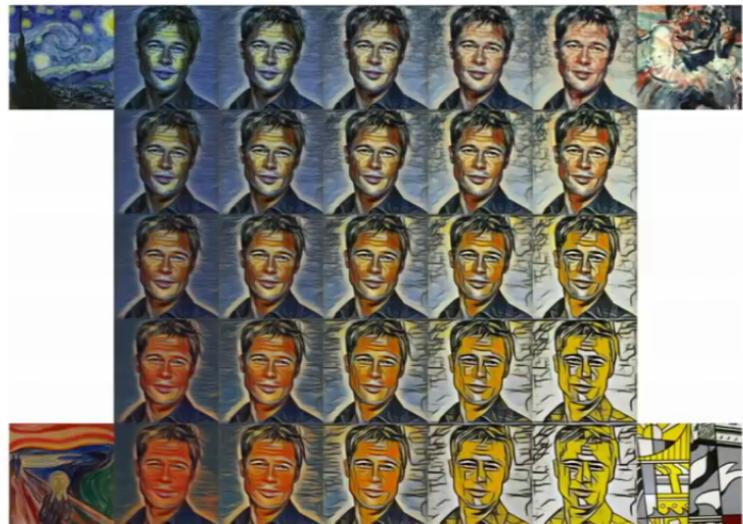
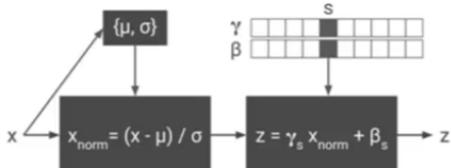
Dumoulin, Shlens, and Kudlur, "A Learned Representation for Artistic Style", ICLR 2017.
Figure copyright Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur, 2016; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 82 University, May 10, 2017

One Network, Many Styles

Use the same network for multiple styles using conditional instance normalization: learn separate scale and shift parameters per style



Single network can blend styles after training

Dumoulin, Shlens, and Kudlur, "A Learned Representation for Artistic Style", ICLR 2017.
Figure copyright Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur, 2016; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 83 University, May 10, 2017

Summary

Many methods for understanding CNN representations

Activations: Nearest neighbors, Dimensionality reduction, maximal patches, occlusion

Gradients: Saliency maps, class visualization, fooling images, feature inversion

Fun: DeepDream, Style Transfer.

Next time: **Unsupervised Learning**

Autoencoders

Variational Autoencoders

Generative Adversarial Networks