



# DEEP LEARNING WORKSHOP

Dublin City University  
21-22 May 2018



## Day 1 Lecture 5

# Visualization

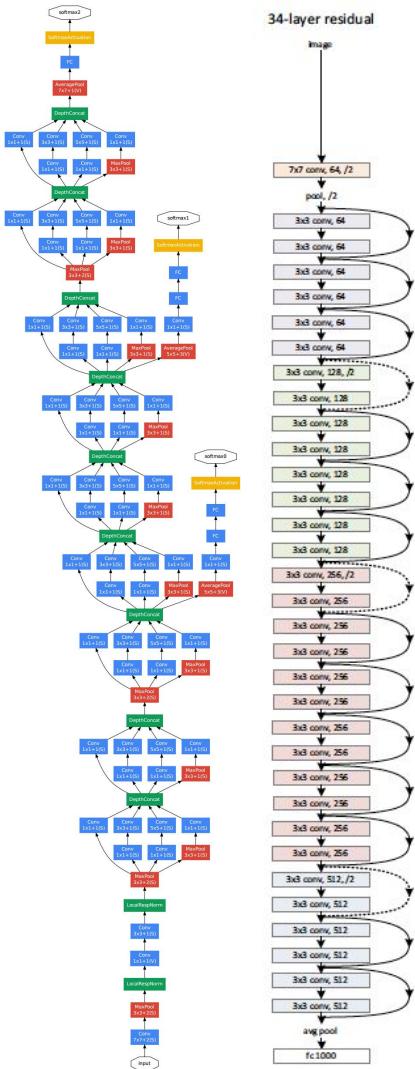


Amaia Salvador  
[amaia.salvador@upc.edu](mailto:amaia.salvador@upc.edu)  
PhD Candidate  
Universitat Politècnica de Catalunya



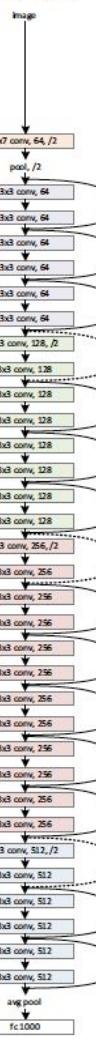
## AlexNet

<b>image</b>
<b>conv-64</b>
<b>conv-64</b>
<b>maxpool</b>
<b>conv-128</b>
<b>conv-128</b>
<b>maxpool</b>
<b>conv-256</b>
<b>conv-256</b>
<b>conv-256</b>
<b>maxpool</b>
<b>conv-512</b>
<b>conv-512</b>
<b>conv-512</b>
<b>maxpool</b>
<b>FC-4096</b>
<b>FC-4096</b>
<b>FC-1000</b>
<b>conv-512</b>
<b>conv-512</b>
<b>conv-512</b>
<b>maxpool</b>
<b>FC-4096</b>
<b>FC-4096</b>
<b>FC-1000</b>
<b>softmax</b>



34-layer residual

Image



## Revolution of Depth

152 layers



ImageNet Classification top-5 error (%)

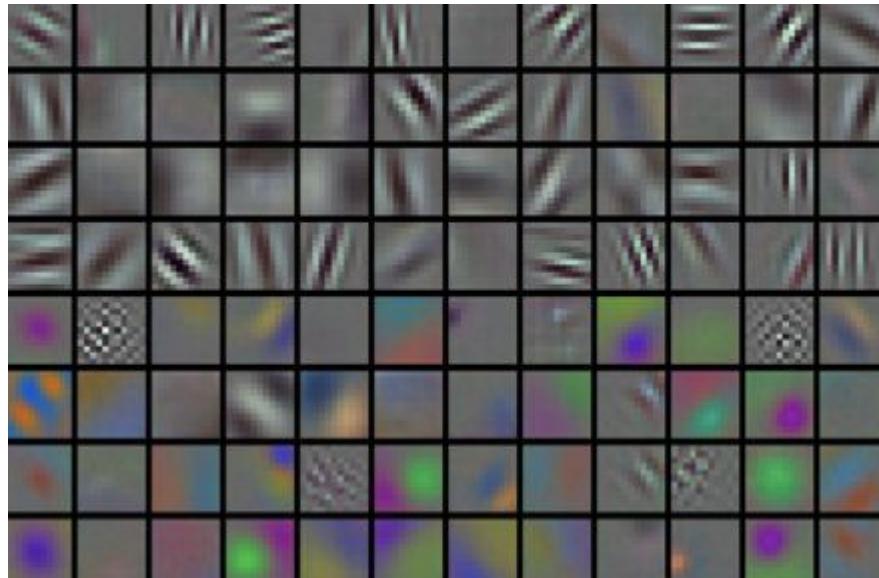
# Visualization

- Learned weights
- Activations from data
- Gradient-based
- Optimization-based

# Visualization

- **Learned weights**
- Activations from data
- Gradient-based
- Optimization-based

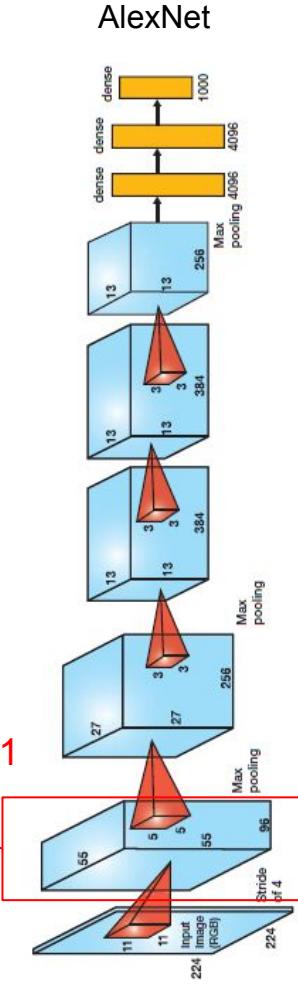
# Visualize Learned Weights



Filters are only “interpretable” on the first layer



conv1



# Visualize Learned Weights

## Weights:

## Weights:

## layer 2 weights

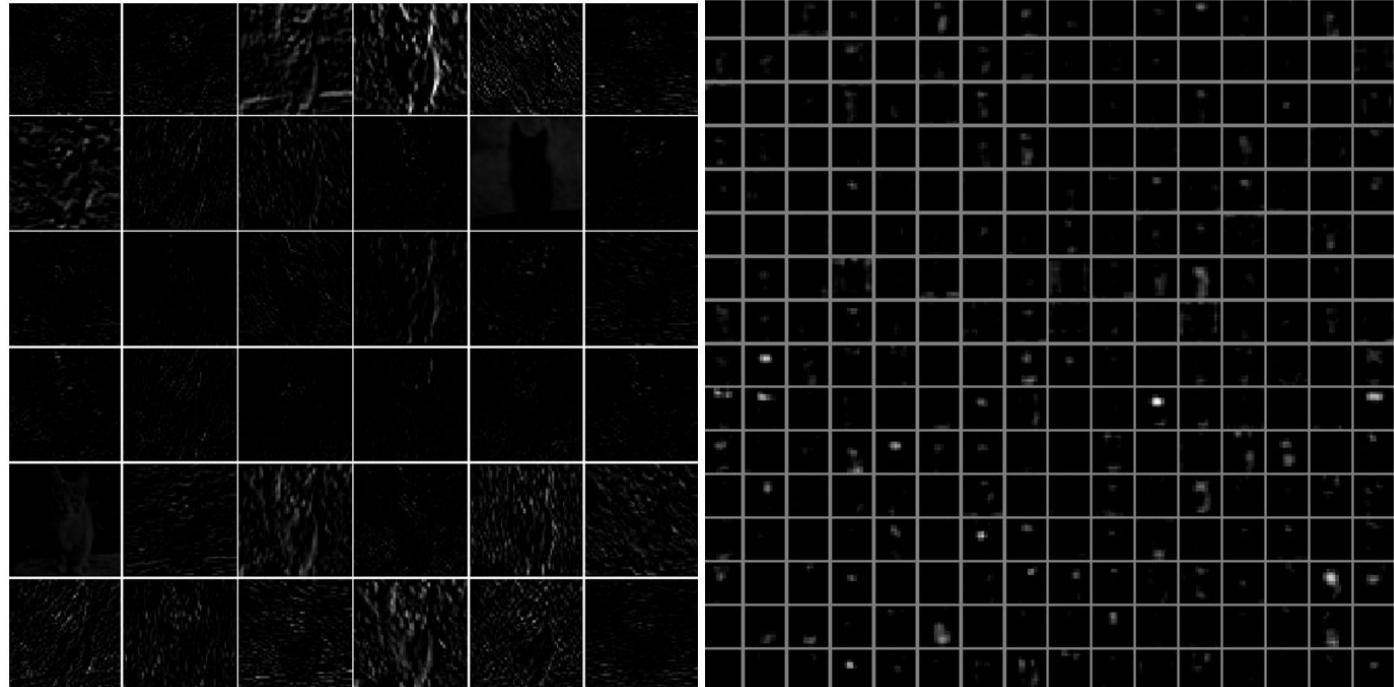
## layer 3 weights

Source: ConvnetJS

# Visualization

- Learned weights
- **Activations from data**
- Gradient-based
- Optimization-based

# Activations from data

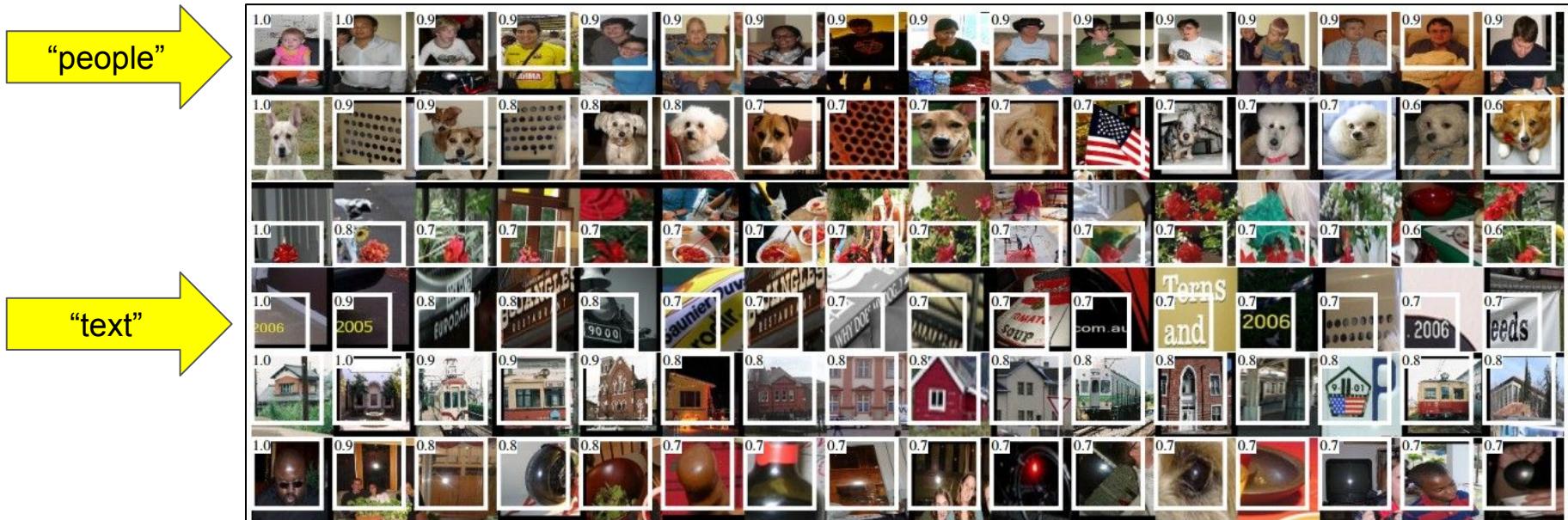


conv1

conv5

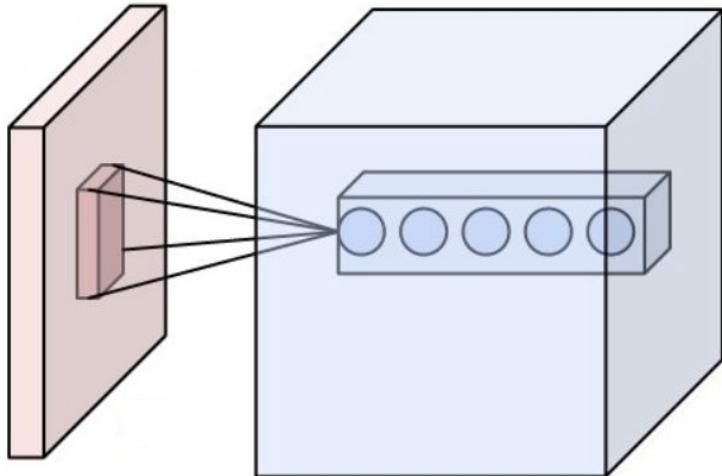
# Receptive Field

Visualize the receptive field of a neuron on those images that activate it the most



**Figure 4: Top regions for six pool<sub>5</sub> units.** Receptive fields and activation values are drawn in white. Some units are aligned to concepts, such as people (row 1) or text (4). Other units capture texture and material properties, such as dot arrays (2) and specular reflections (6).

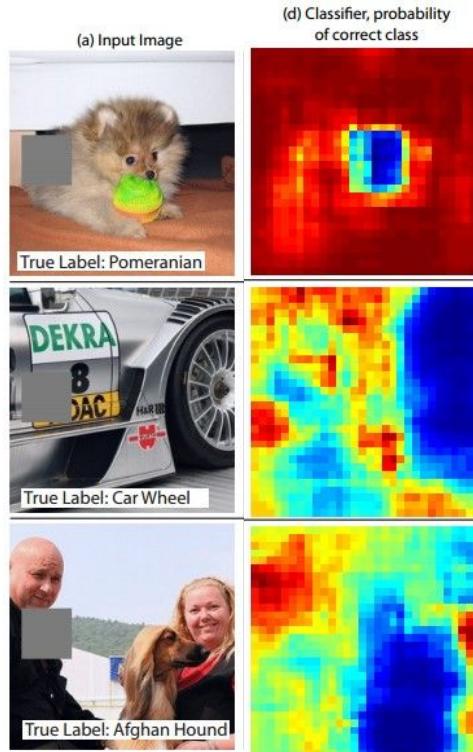
# Reminder: Receptive Field



**Receptive field:** Part of the input that is visible to a neuron. It increases as we stack more convolutional layers (i.e. neurons in deeper layers have larger receptive fields).

# Occlusion experiments

1. Iteratively forward the same image through the network, occluding a different region at a time.
2. Keep track of the probability of the correct class w.r.t. the position of the occluder

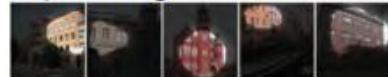


# Occlusion experiments

Can be applied to any arbitrary layer. “Manual” label assignment (AMT)

## Buildings

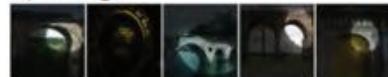
56) building



120) arcade



8) bridge



123) building



119) building



## Indoor objects

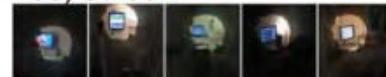
182) food



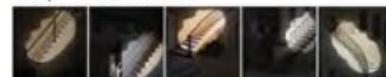
46) painting



106) screen



53) staircase



107) wardrobe

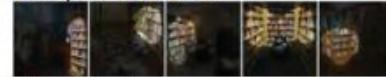


## Furniture

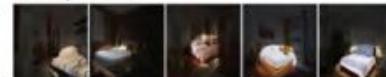
18) billiard table



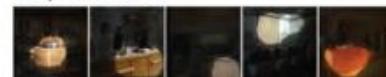
155) bookcase



116) bed



38) cabinet

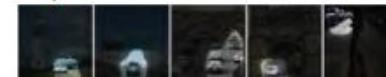


85) chair



## Outdoor objects

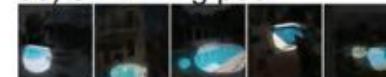
87) car



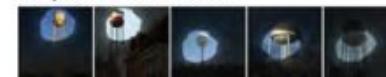
61) road



96) swimming pool



28) water tower

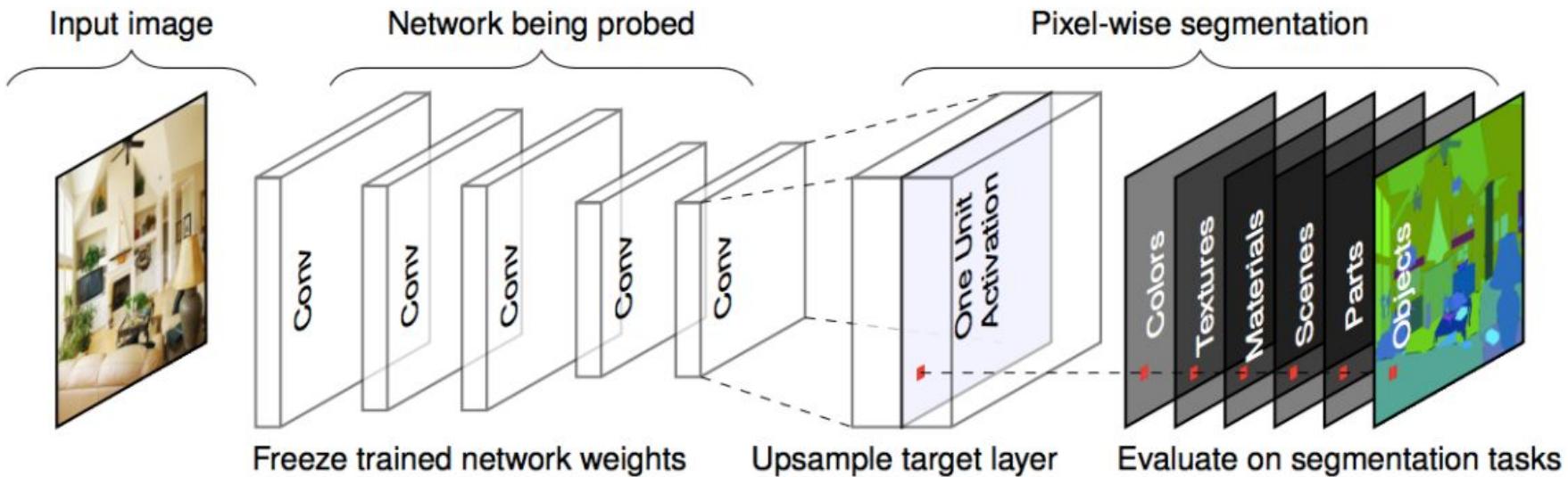


6) windmill

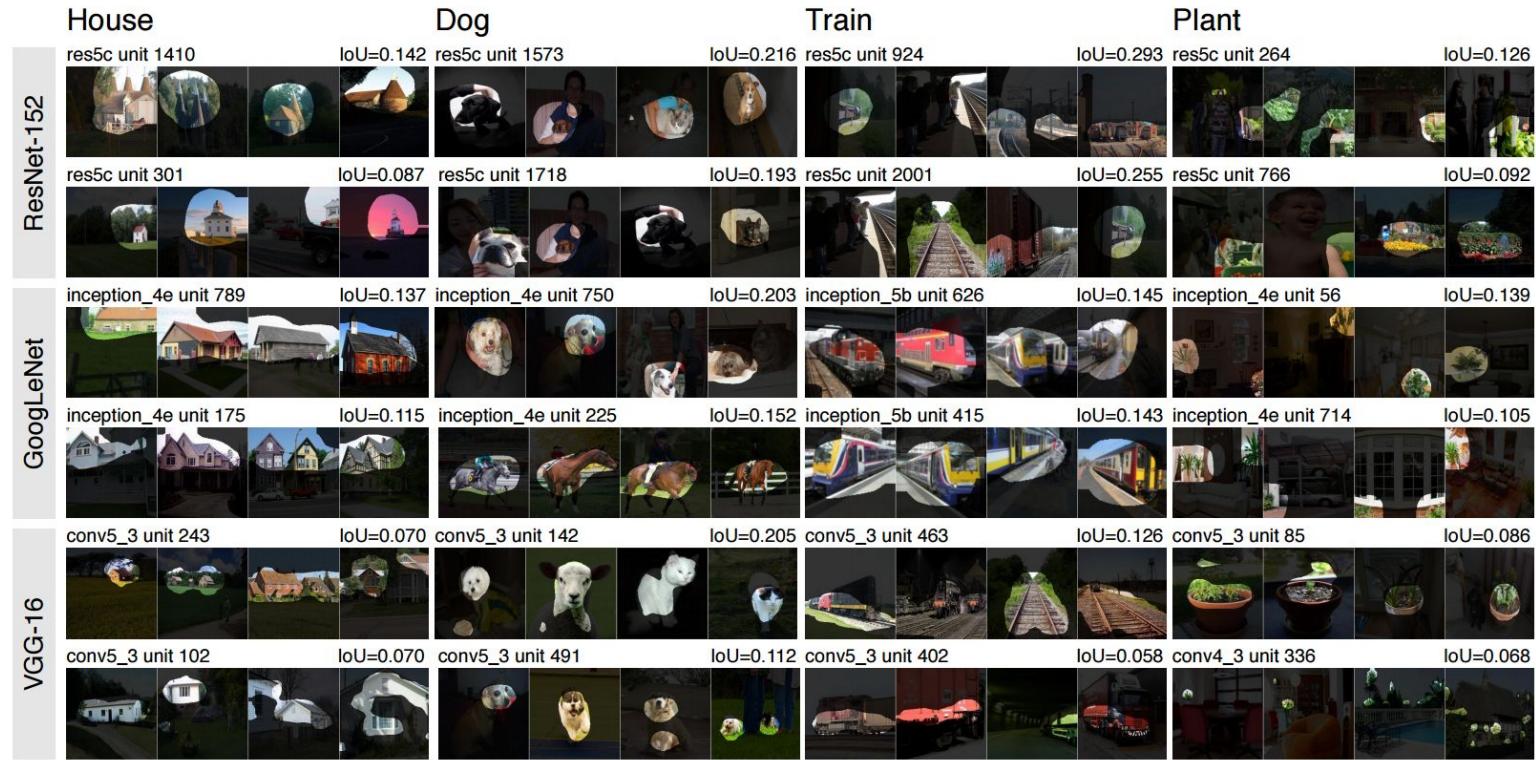


# Network Dissection

Same idea, but automatic unit labeling using densely labeled dataset (pixel-level annotations).  
The thresholded activation of each conv unit in the network is evaluated for semantic segmentation.

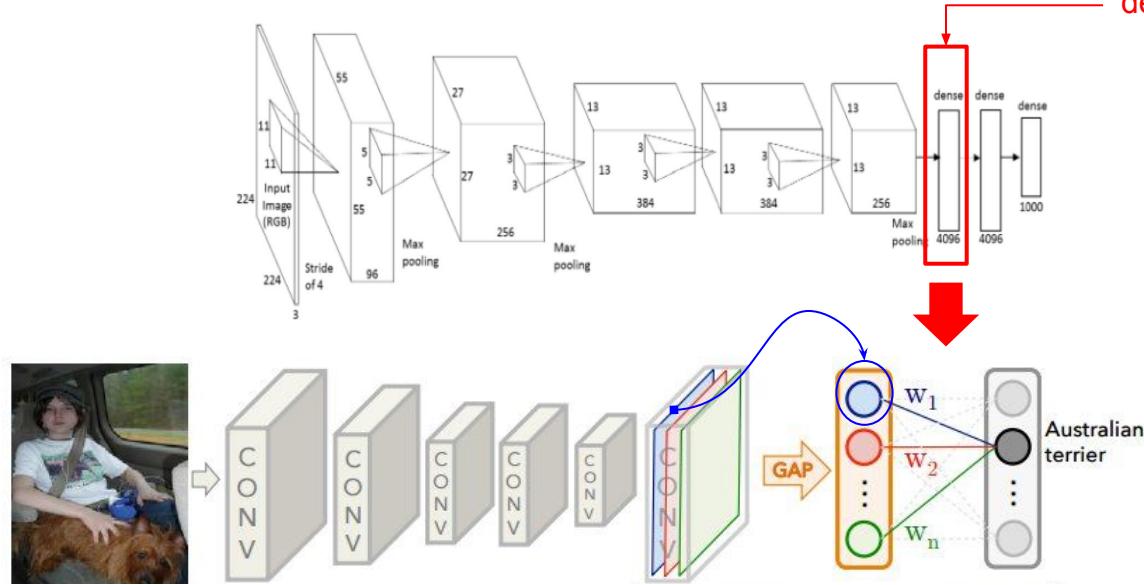


# Network Dissection



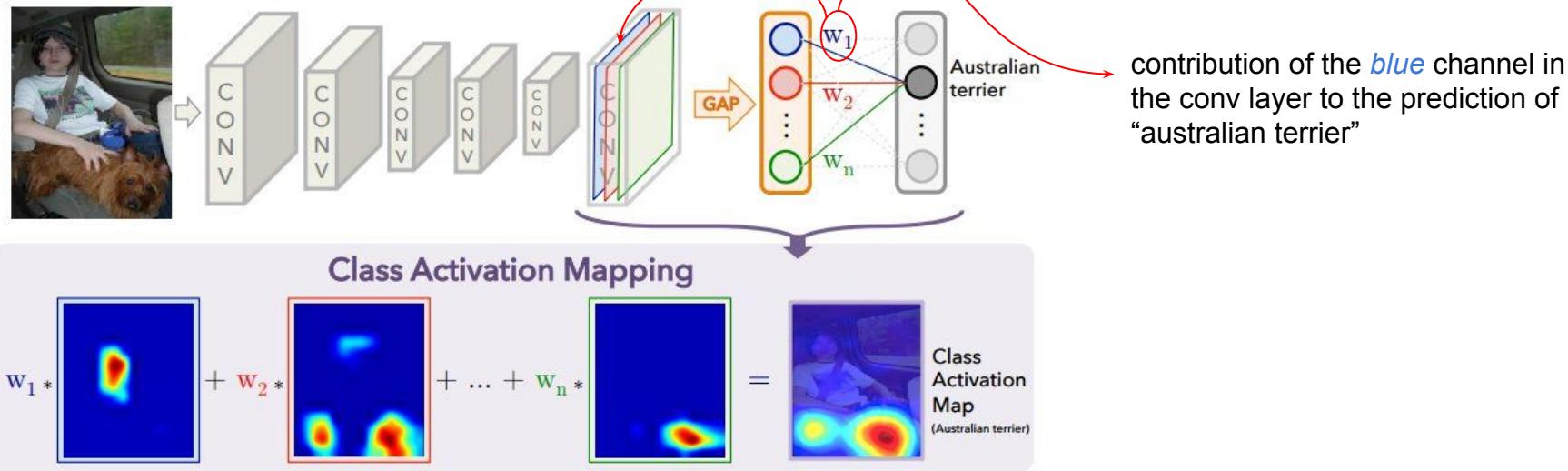
# Class Activation Maps

**Class Activation Maps (CAMs):** Replace FC layer after last conv with Global Average Pooling (GAP).  
densely connected weights



# Class Activation Maps

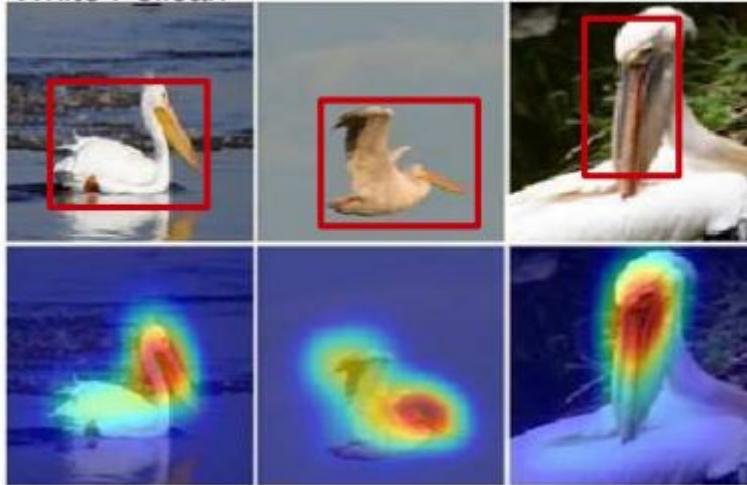
**Class Activation Maps (CAMs):** The classifier weights define the contribution of each channel in the previous layer



# Class Activation Maps

**Class Activation Maps (CAMs):** Unsupervised object localization

White Pelican



Orchard Oriole

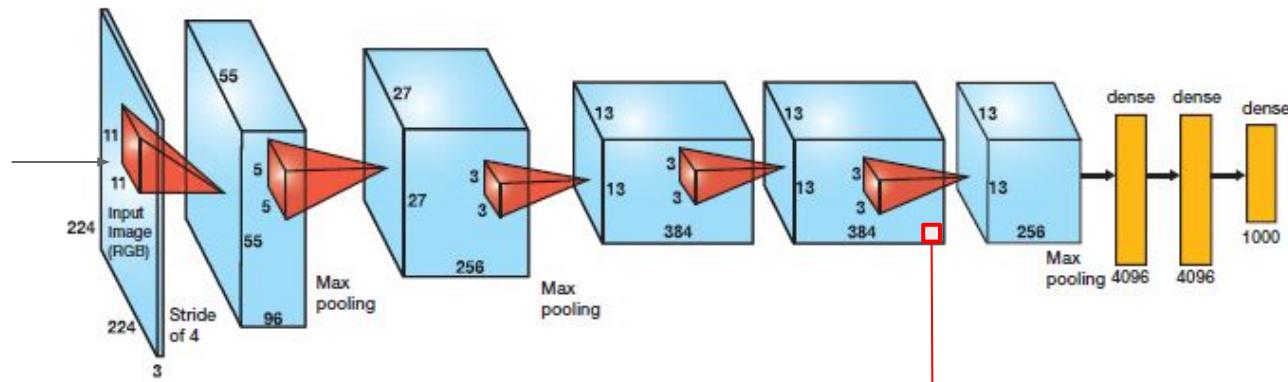


# Visualization

- Learned weights
- Activations from data
- **Gradient-based**
- Optimization-based

# Gradient-based approach

Visualize the part of an image that mostly activates a neuron



Compute the gradient of any neuron w.r.t. the image

1. Forward image up to the desired layer (e.g. conv5)
2. Set all gradients to 0
3. Set gradient for the neuron we are interested in to 1
4. Backpropagate to get reconstructed image (gradient on the image)

# Gradient-based approach

1. Forward image up to the desired layer (e.g. conv5)
2. Set all gradients to 0
3. Set gradient for the neuron we are interested in to 1
4. Backpropagate to get reconstructed image (gradient on the image)



# Gradient-based approach

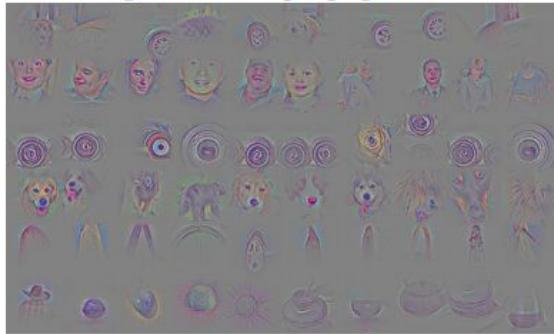
guided backpropagation



corresponding image crops



guided backpropagation



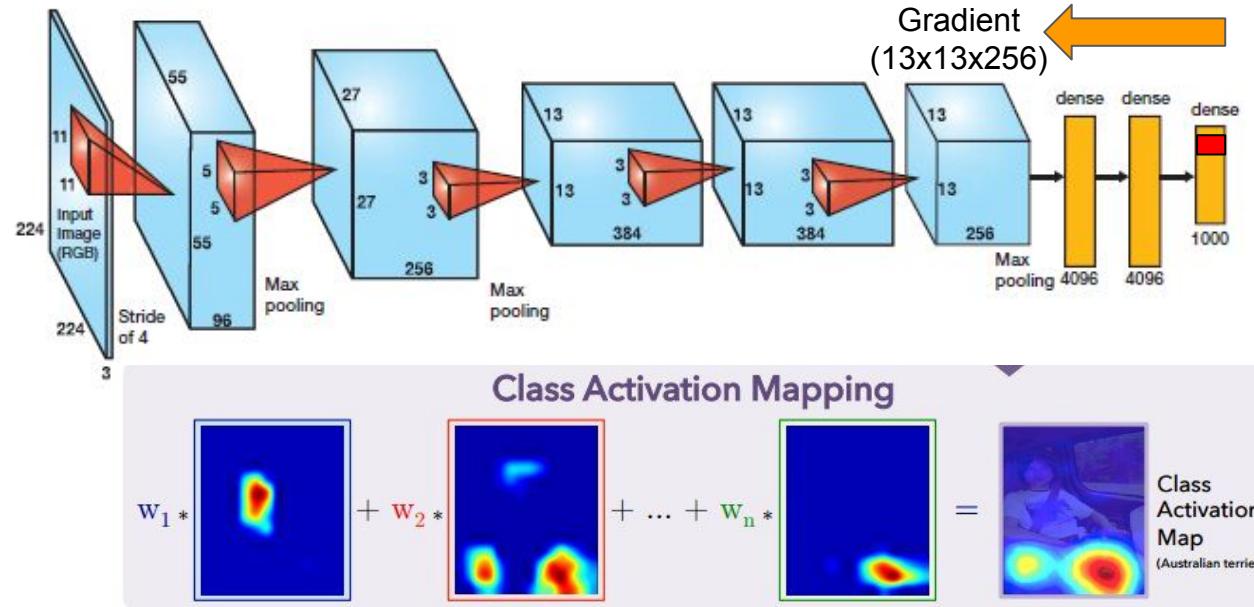
corresponding image crops



# Grad-CAM

Same idea as CAM, but:

- No modifications required to the network
- Weight feature map by the gradient of the class wrt each channel

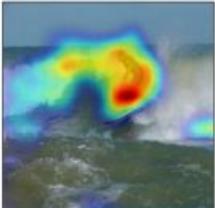


Feature maps weighted by  
the mean gradient

# Grad-CAM



What is the man doing?



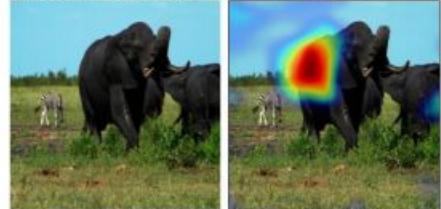
Surfing



What is she holding?



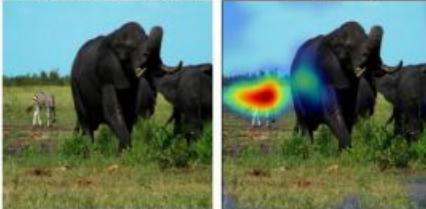
Baseball bat



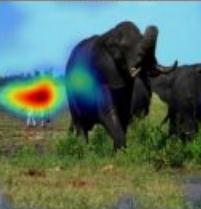
What is that?



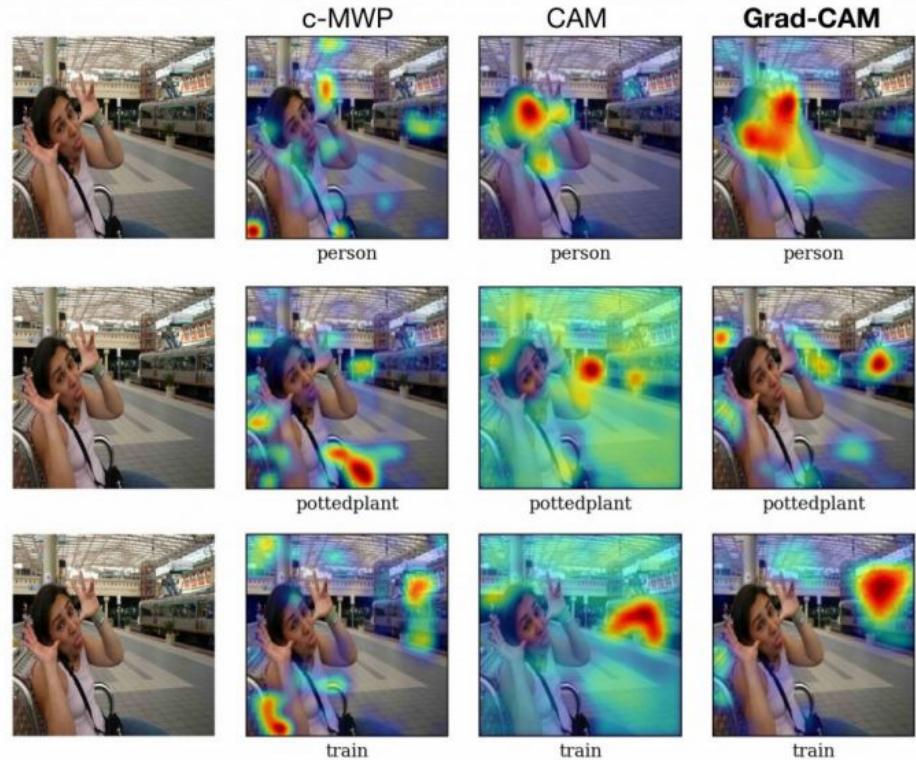
Elephant



What is that?



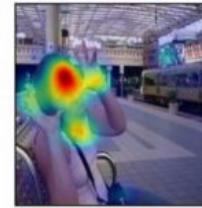
Zebra



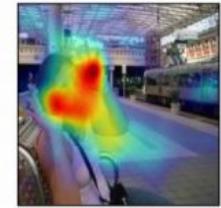
c-MWP



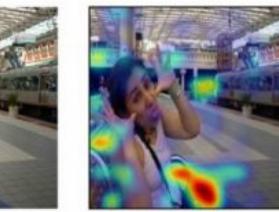
person



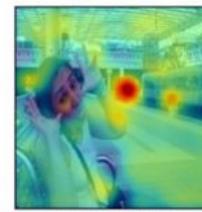
person



person



pottedplant



pottedplant



train



train



train



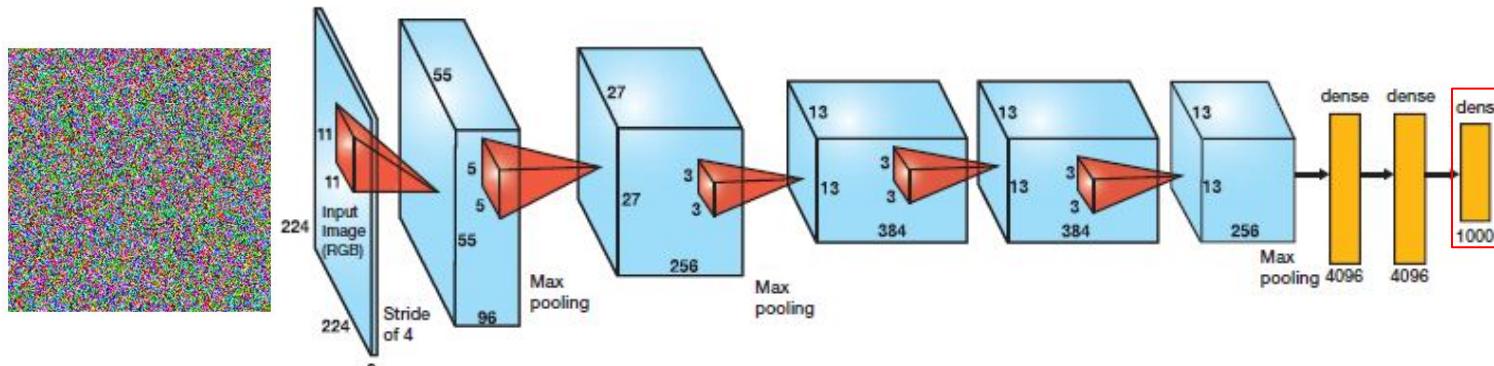
train

# Visualization

- Learned weights
- Activations from data
- Gradient-based
- **Optimization-based**

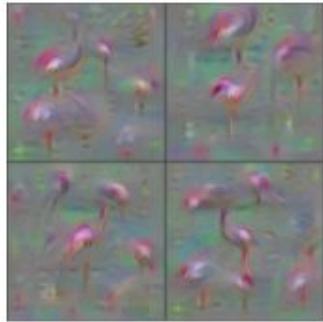
# Optimization approach

Obtain the image that maximizes a class score (or a neuron activation)

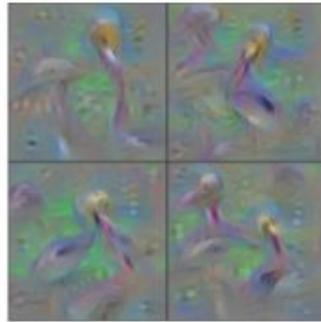


1. Forward random image
2. Set the gradient of the scores vector to be  $[0,0,0\dots,1,\dots,0,0]$
3. Backprop to get gradient on the image
4. Update image (small step in the gradient direction)
5. Repeat

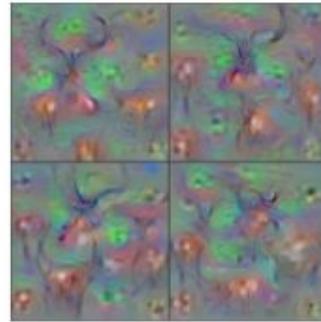
# Optimization approach



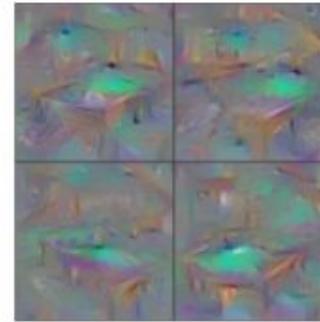
Flamingo



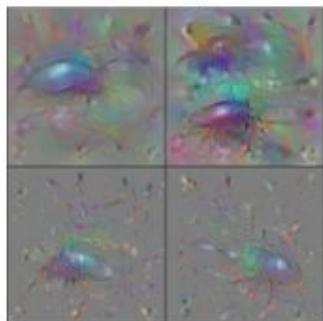
Pelican



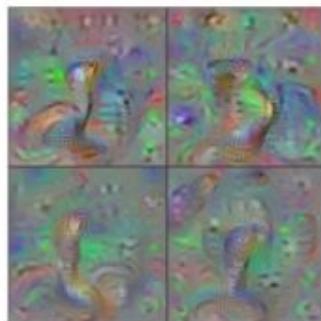
Hartebeest



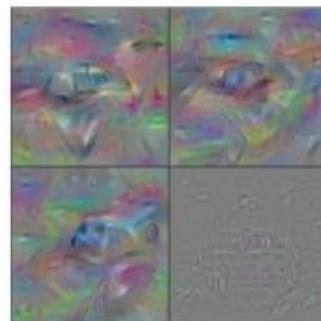
Billiard Table



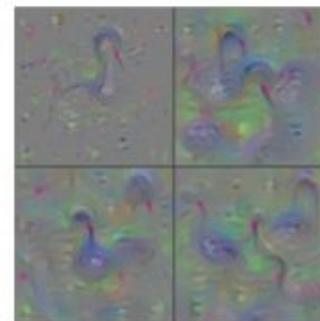
Ground Beetle



Indian Cobra



Station Wagon

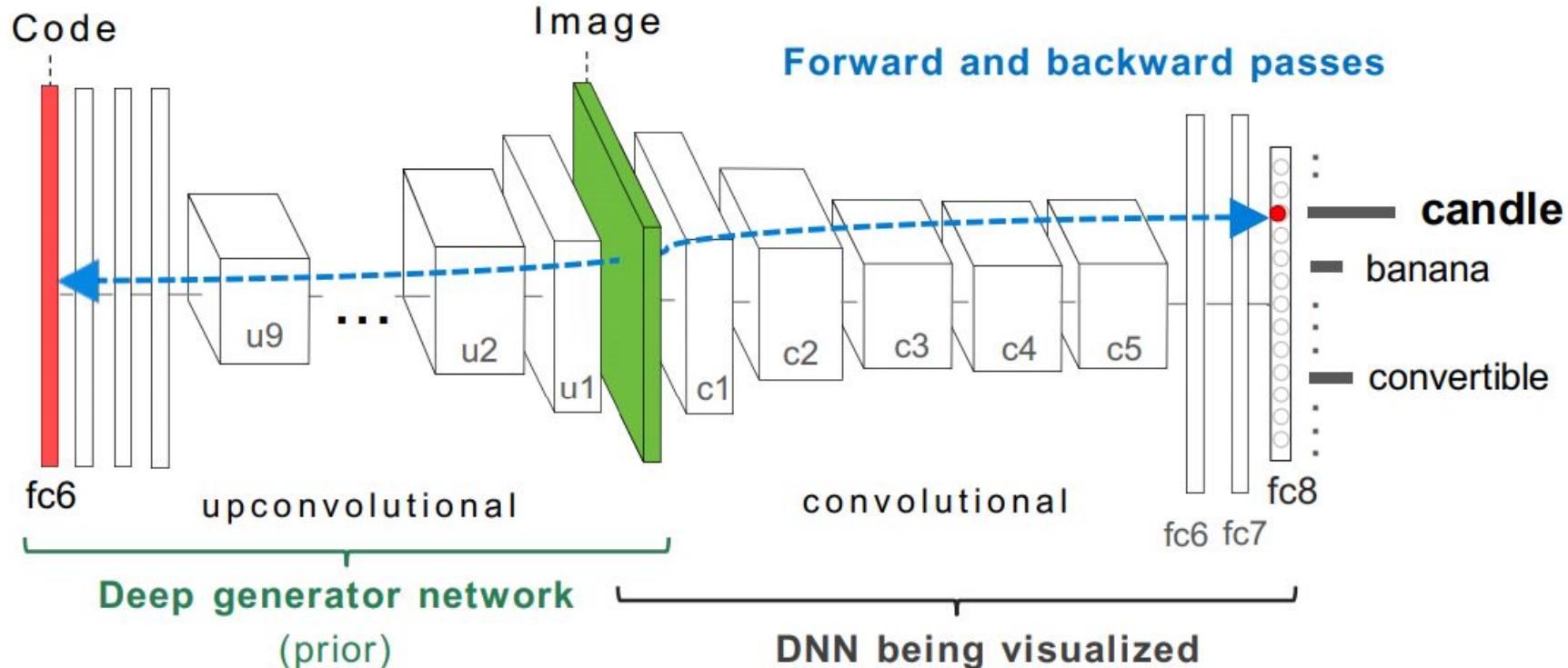


Black Swan

# Regularization tricks

	Unregularized	Frequency Penalization	Transformation Robustness	Learned Prior	Dataset Examples
 <b>Erhan, et al., 2009 [3]</b> Introduced core idea. Minimal regularization.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 <b>Szegedy, et al., 2013 [11]</b> Adversarial examples. Visualizes with dataset examples.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
 <b>Mahendran &amp; Vedaldi, 2015 [7]</b> Introduces total variation regularizer. Reconstructs input from representation.	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 <b>Nguyen, et al., 2015 [14]</b> Explores counterexamples. Introduces image blurring.	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 <b>Mordvintsev, et al., 2015 [4]</b> Introduced jitter & multi-scale. Explored GMM priors for classes.	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
 <b>Oygard, et al., 2015 [15]</b> Introduces gradient blurring. (Also uses jitter.)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 <b>Tyka, et al., 2016 [16]</b> Regularizes with bilateral filters. (Also uses jitter.)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 <b>Mordvintsev, et al., 2016 [17]</b> Normalizes gradient frequencies. (Also uses jitter.)	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
 <b>Nguyen, et al., 2016 [18]</b> Paramaterizes images with GAN generator.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
 <b>Nguyen, et al., 2016 [19]</b> Uses denoising autoencoder prior to make a generative model.	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

# Regularization tricks: Generative prior



Nguyen et al. [Synthesizing the preferred inputs for neurons in neural networks via deep generator networks](#), NIPS 2016

# Regularization tricks



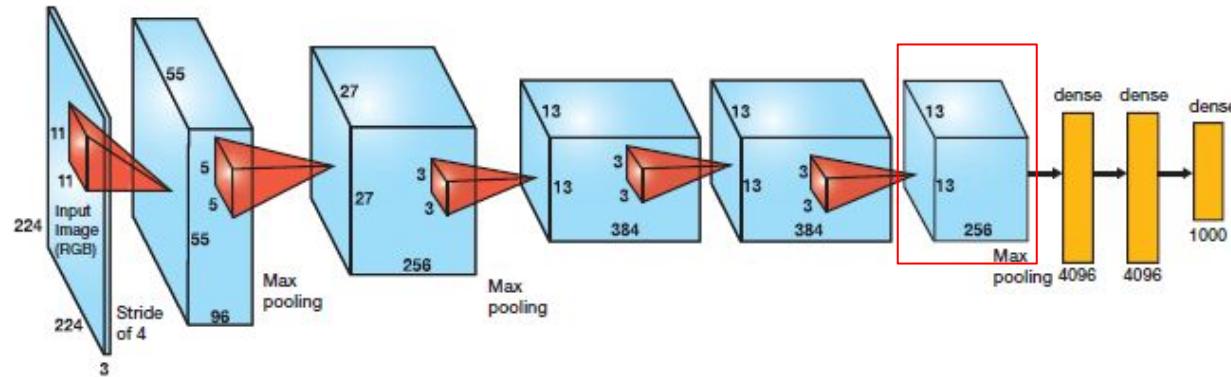
Nguyen et al. [Synthesizing the preferred inputs for neurons in neural networks via deep generator networks](#), NIPS 2016

# Extra: DeepDream



<https://github.com/google/deepdream>

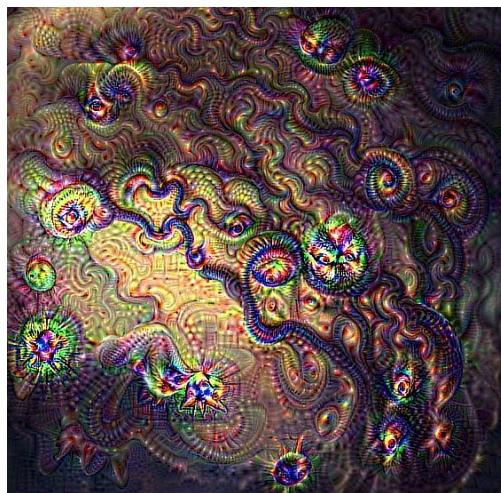
# Extra: DeepDream



1. Forward image up to some layer (e.g. conv5)
2. Set the gradients to equal the layer activations
3. Backprop to get gradient on the image
4. Update image (small step in the gradient direction)
5. Repeat

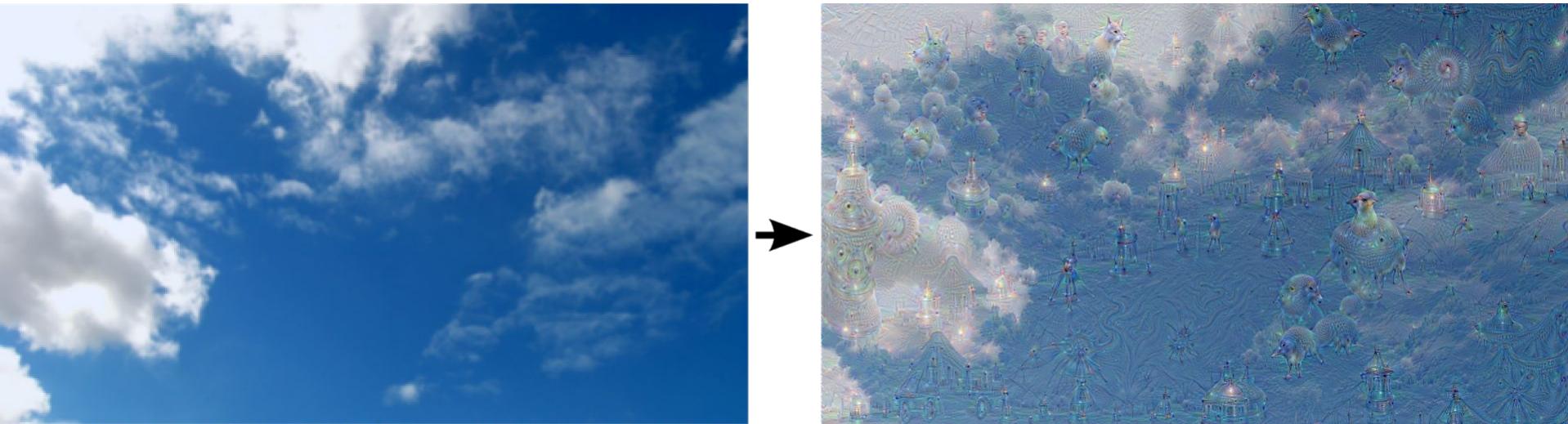
# Extra: DeepDream

1. Forward image up to some layer (e.g. conv5)
2. **Set the gradients to equal the layer activations**
3. Backprop to get gradient on the image
4. Update image (small step in the gradient direction)
5. Repeat

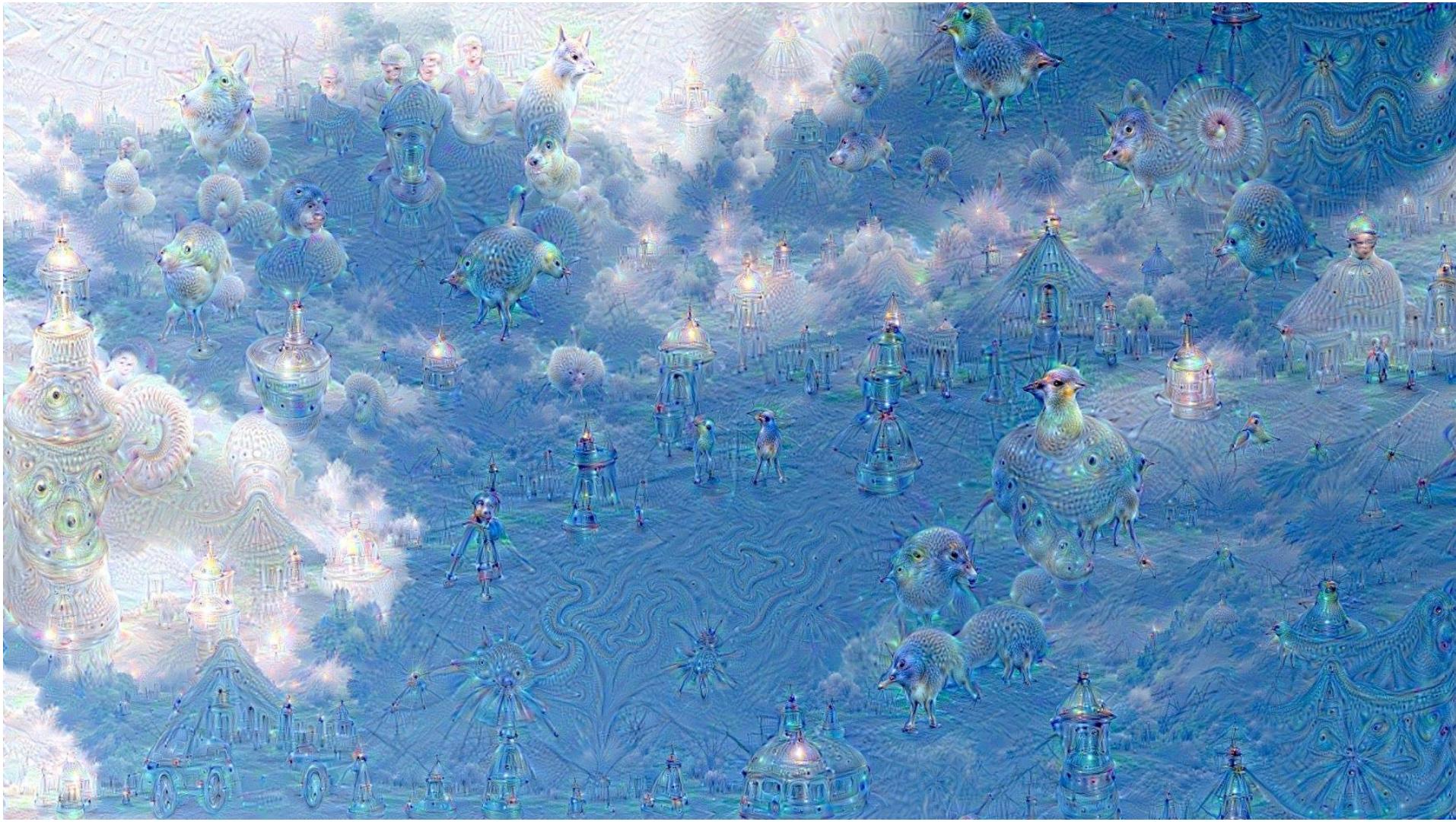


At each iteration, the image is updated to **boost all features** that activated in that layer in the forward pass.

# Extra: DeepDream



More examples [here](#)



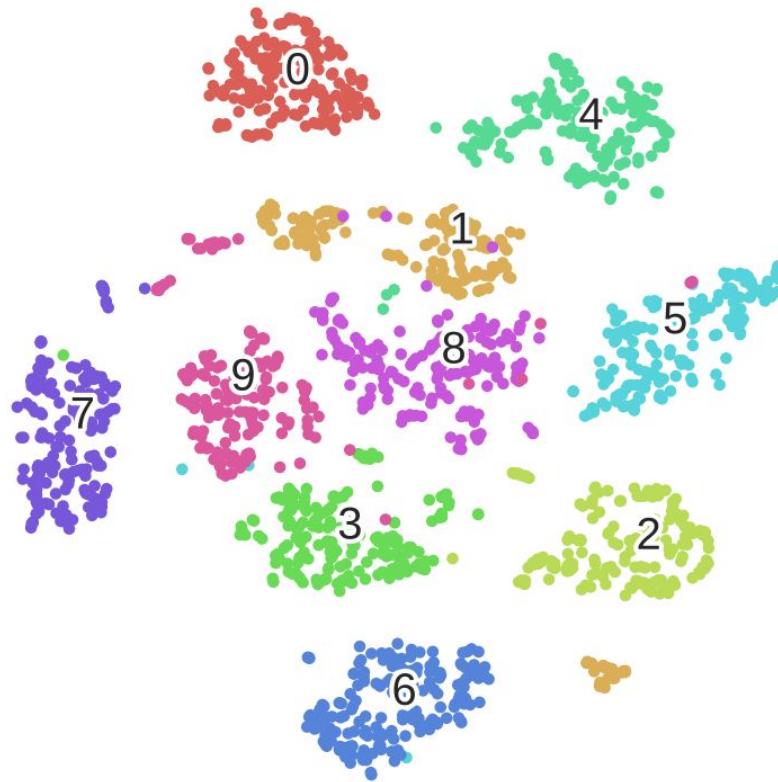
# Visualization

- Learned weights
- Activations from data
- Gradient-based
- Optimization-based

# Questions?

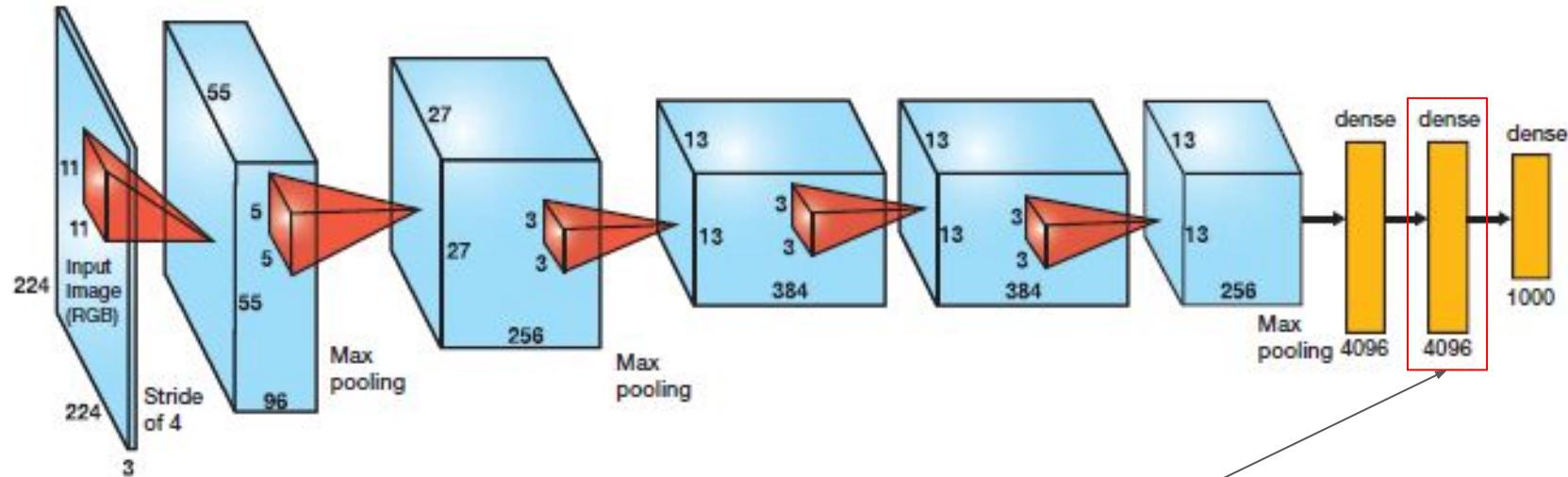
# t-SNE

Embed high dimensional data points (i.e. feature codes) so that pairwise distances are preserved in local neighborhoods.



Maaten & Hinton. [Visualizing High-Dimensional Data using t-SNE](#).  
Journal of Machine Learning Research (2008).

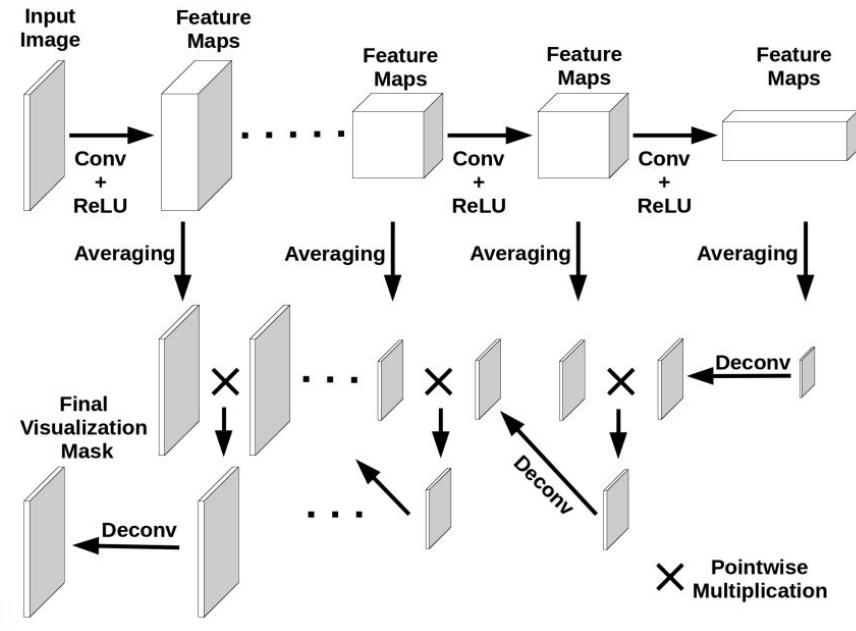
# t-SNE



Can be used with features from layer before classification

# Optimization approach

## Value-based backpropagation



# Neural Style



Style image

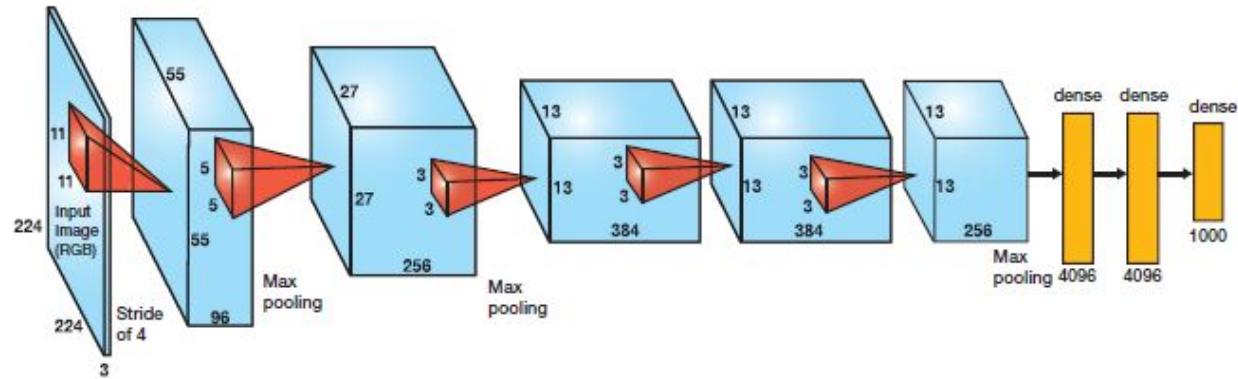


Content image



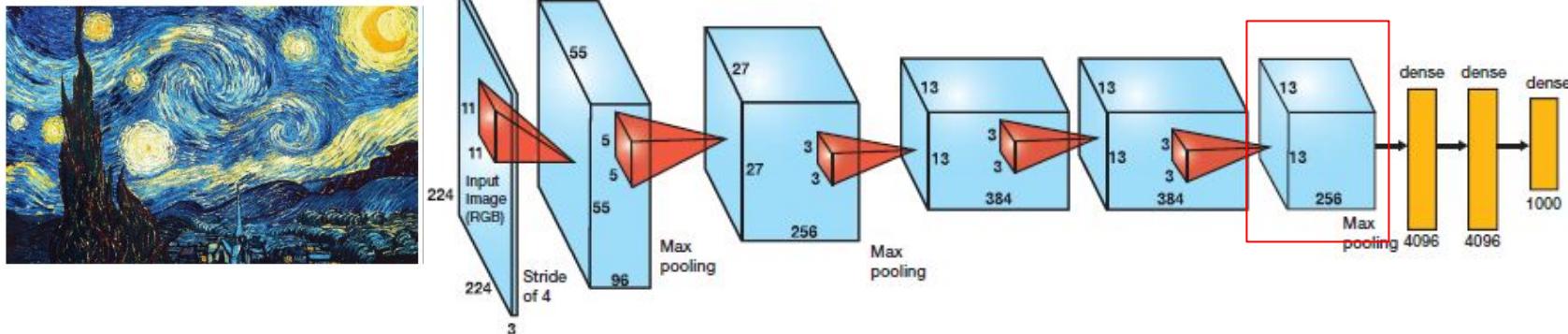
Result

# Neural Style



Extract raw activations in all layers. These activations will represent the contents of the image.

# Neural Style



- Activations are also extracted from the style image for all layers.
  - Instead of the raw activations, gram matrices ( $G$ ) are computed at each layer to represent the style.

E.g. at conv5 [13x13x256], reshape to:

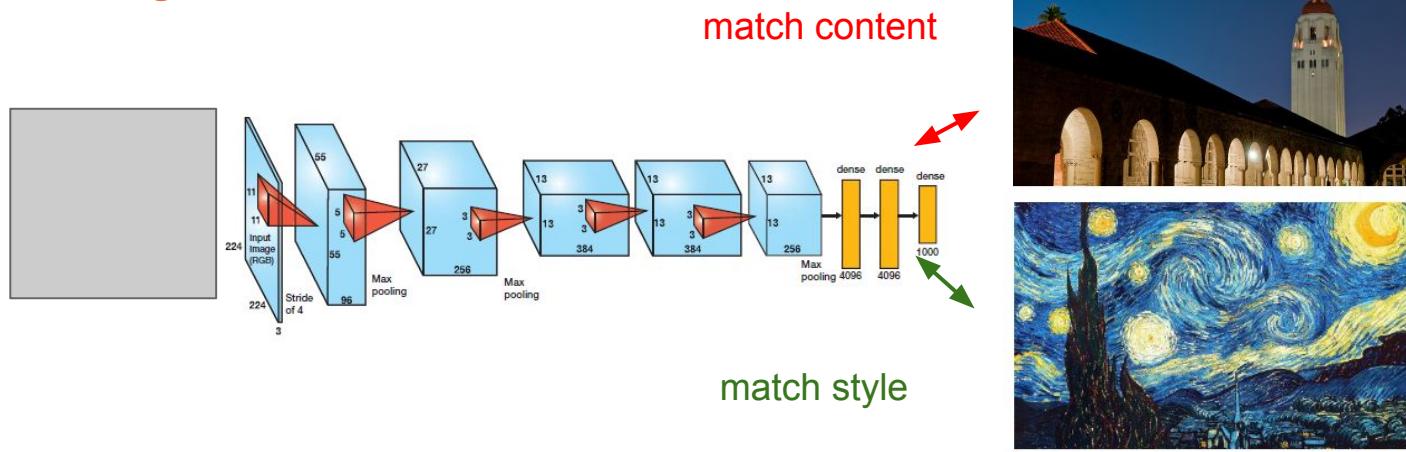
$V =$

...

$G = V^T V$

The Gram matrix  $G$  gives the correlations between filter responses.

# Neural Style

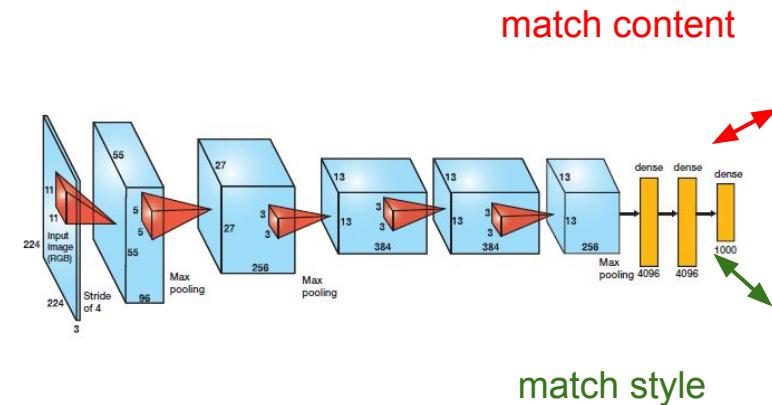


$$\mathcal{L}_{total}(\vec{p}, \vec{a}, \vec{x}) = \alpha \mathcal{L}_{content}(\vec{p}, \vec{x}) + \beta \mathcal{L}_{style}(\vec{a}, \vec{x})$$

Match activations  
from content image

Match gram matrices  
from style image

# Neural Style

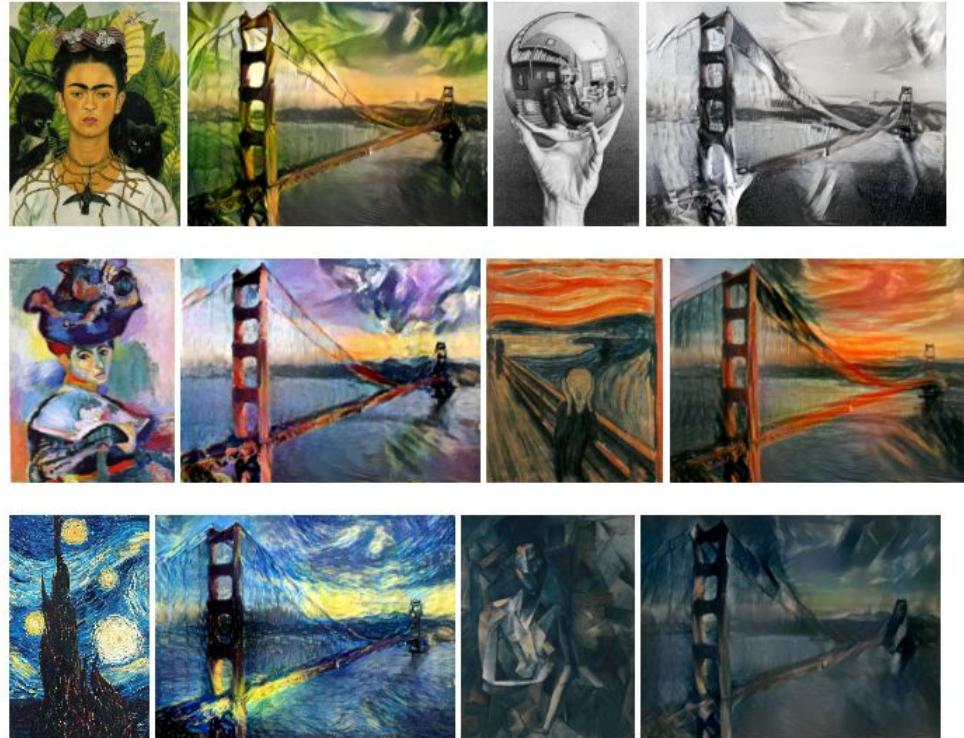


$$\mathcal{L}_{total}(\vec{p}, \vec{a}, \vec{x}) = \alpha \mathcal{L}_{content}(\vec{p}, \vec{x}) + \beta \mathcal{L}_{style}(\vec{a}, \vec{x})$$

Match activations  
from content image

Match gram matrices  
from style image

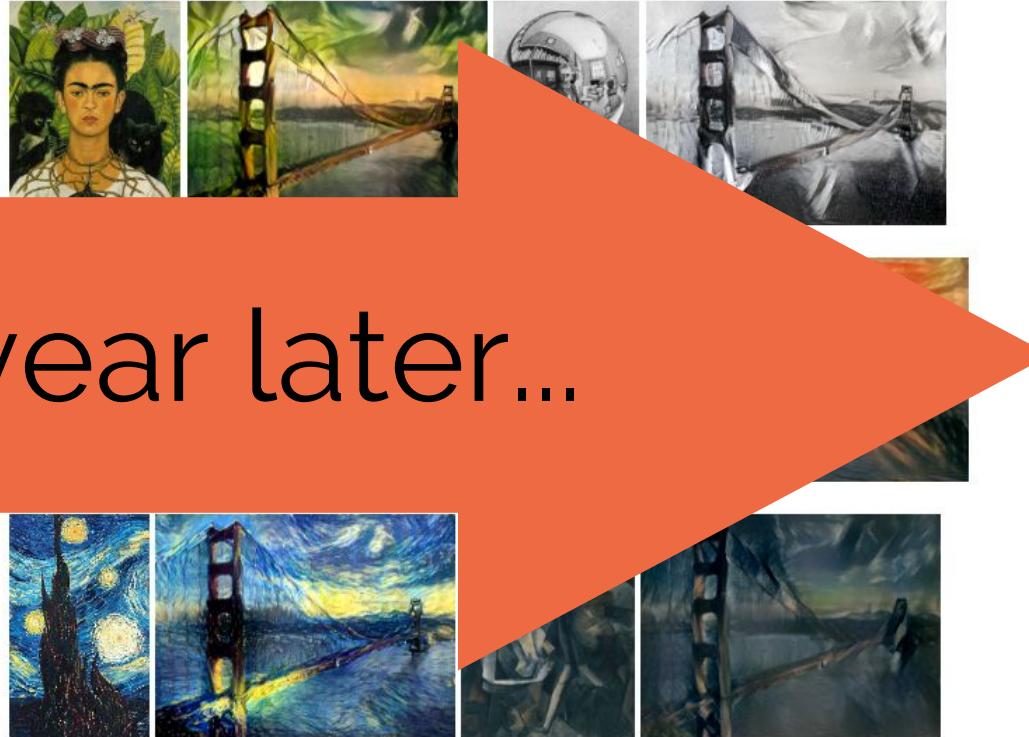
# Neural Style

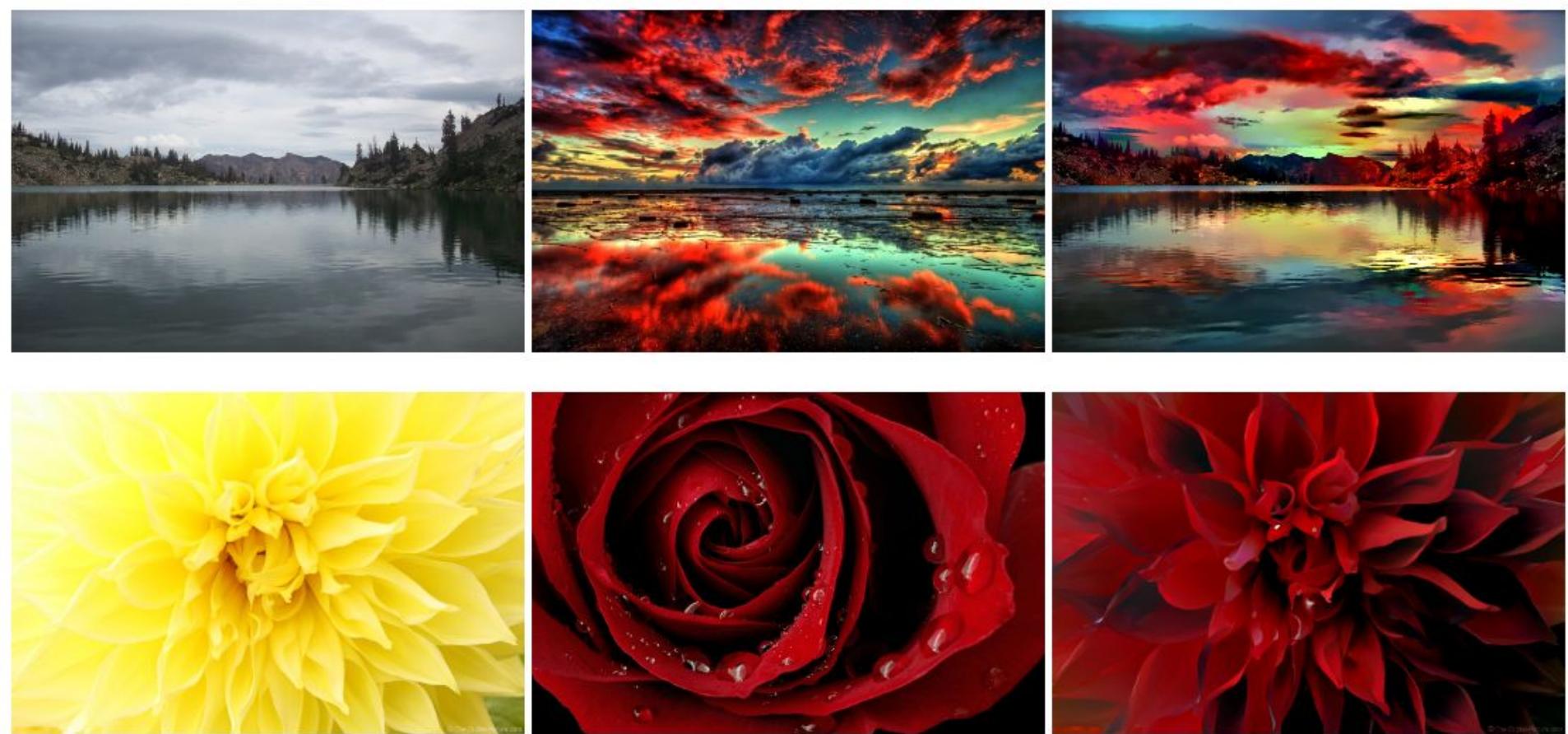


# Neural Style



A year later...

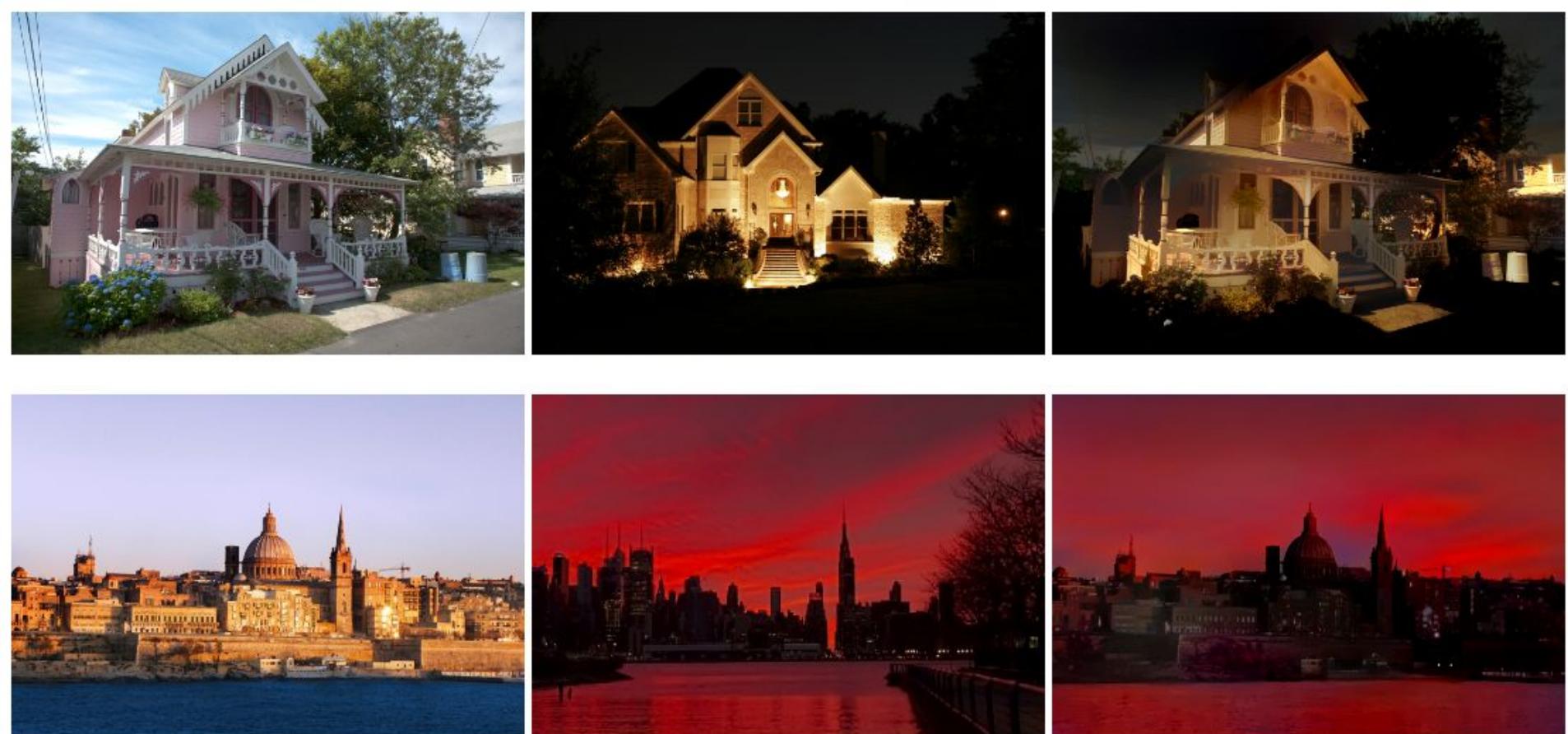




Content Image

Style Image

Result



Content Image

Style Image

Result