## CSW182 (2021)· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

**视频·B 站 [扫码或点击链接]**
*https://www.bilibili.com/video/BV1Ff4y1n7ar*

**课件 & 代码·博客 [扫码或点击链接]**
*http://blog.showmeai.tech/berkeley-csw182*

Berkeley
循环神经网络          可视化          梯度策略
Q-Learning    风格迁移      模仿学习    元学习
计算机视觉      机器学习基础    生成模型    卷积网络

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS23In |

### # Awesome AI Courses Notes Cheatsheets· 持续更新中

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏
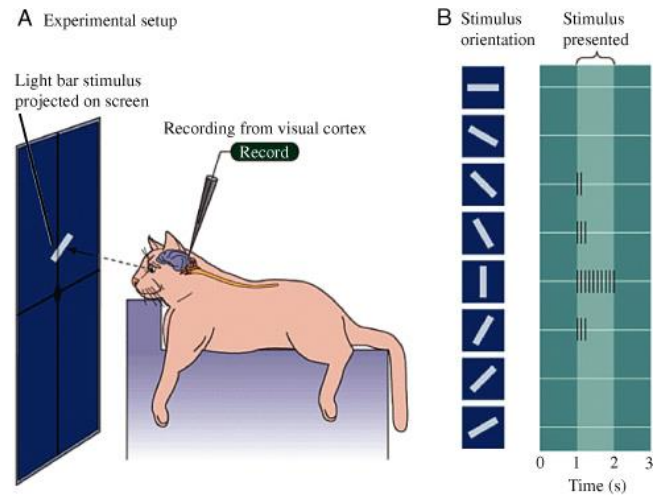
称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]
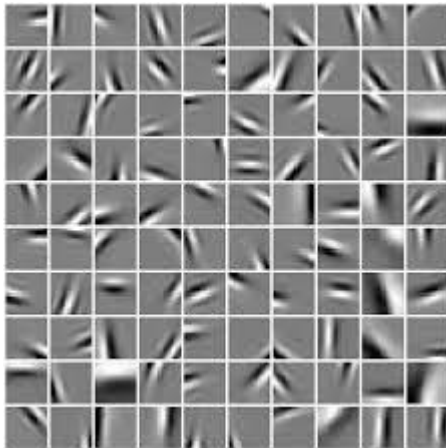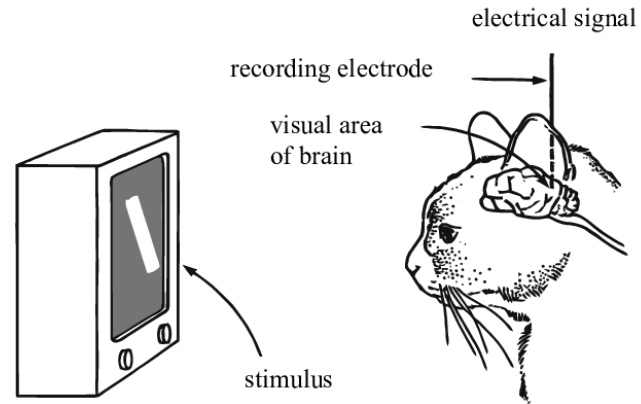
# Generating Images from CNNs

Designing, Visualizing and Understanding Deep Neural Networks
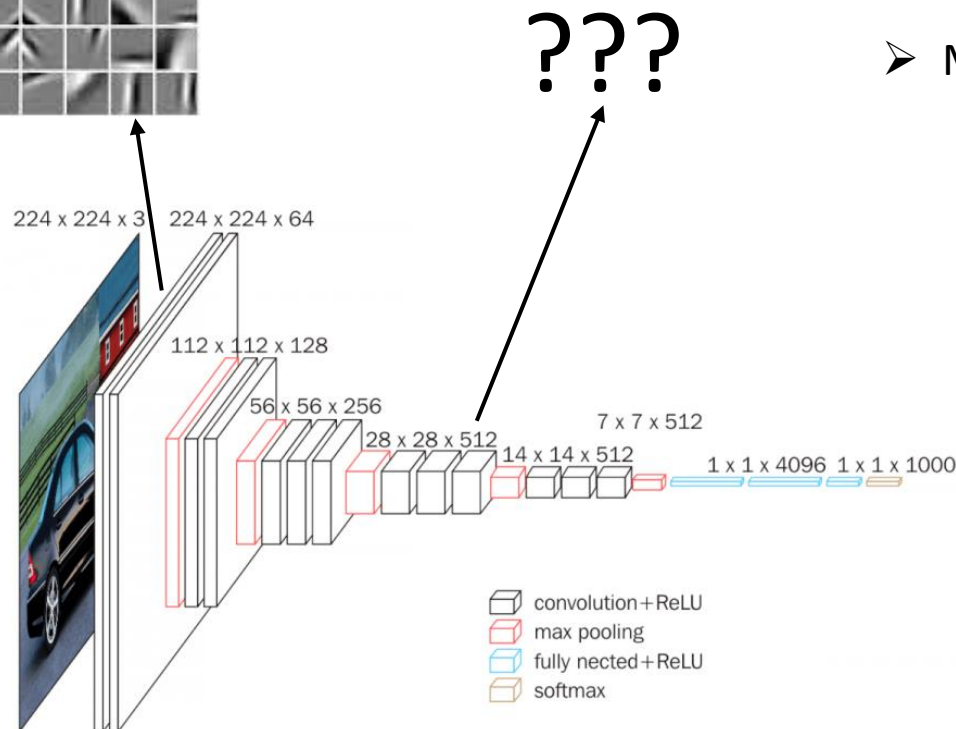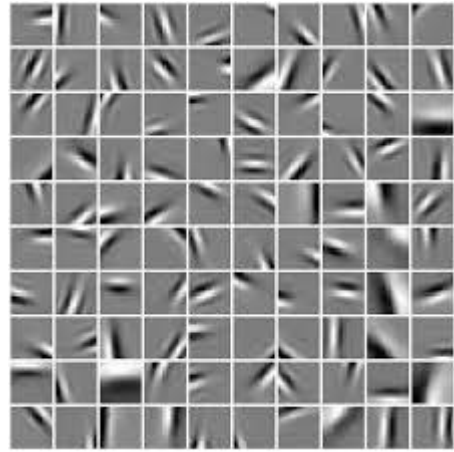
# CS W182/282A

Instructor: Sergey Levine
UC Berkeley

# What does the brain see?

# What do convolutional networks "see"?

**???**

Why are we interested in this?
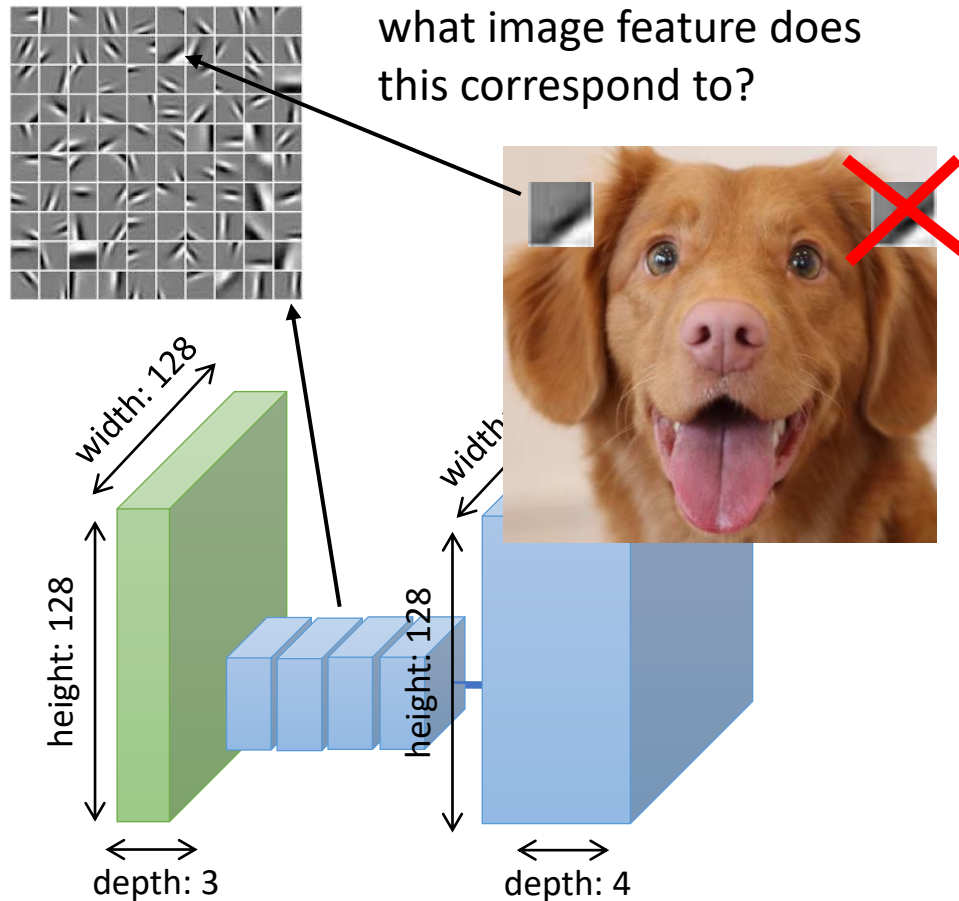
➢ Interpret what neural nets pay attention to

➢ Understand when they'll work and when they won't

➢ Compare different models and architectures
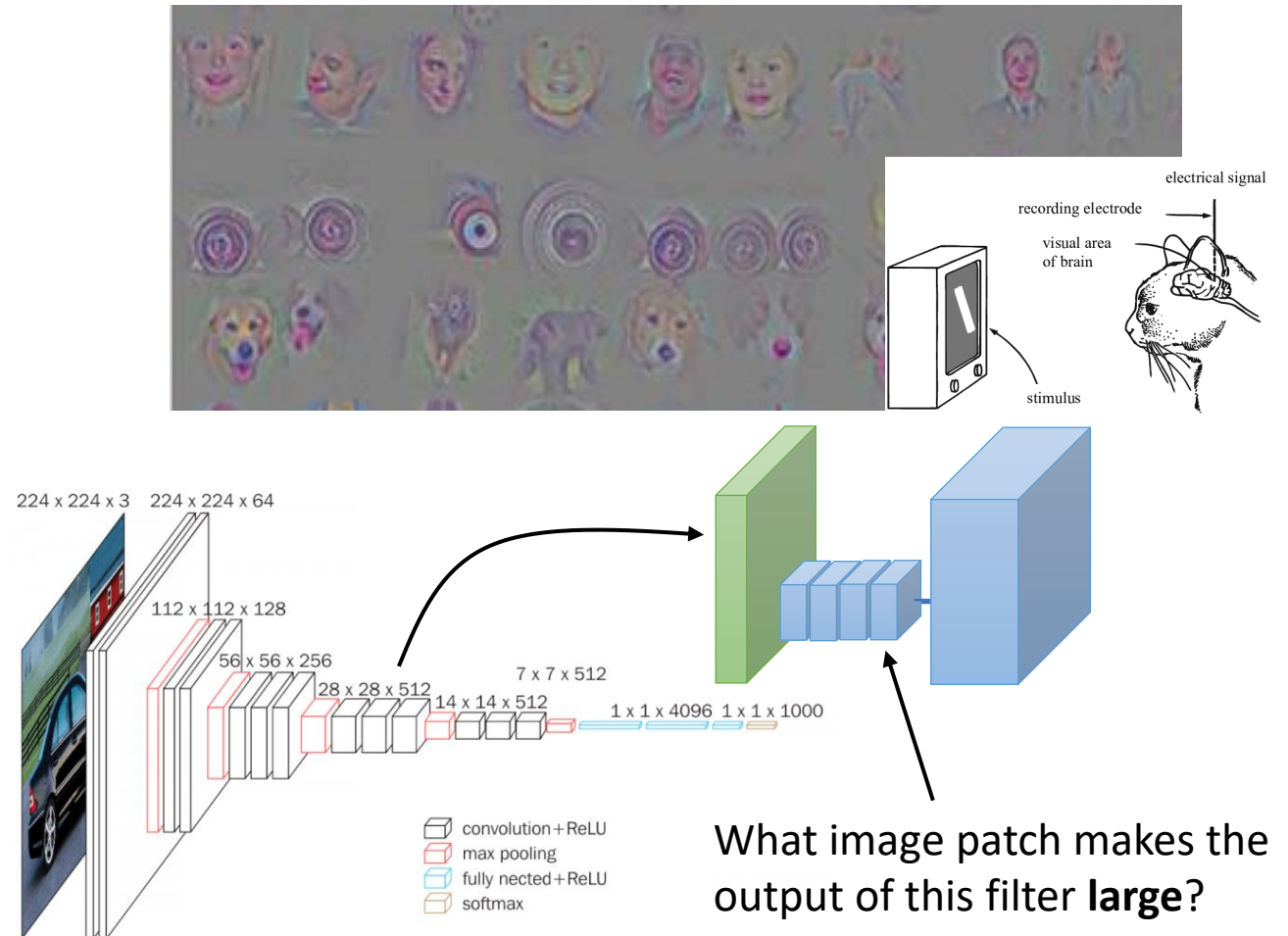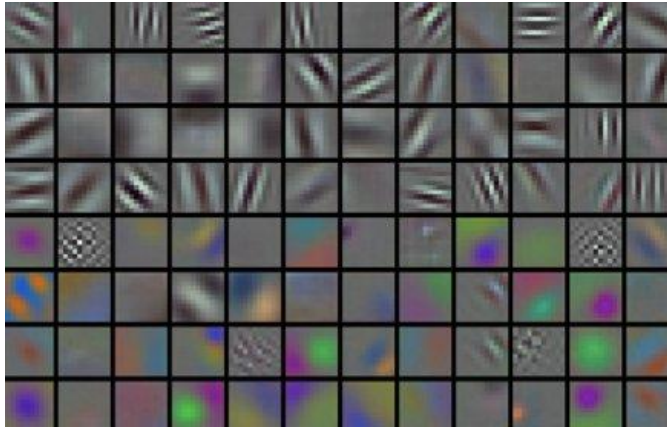
➢ Manipulate images based on conv net responses



224 x 224 x 3    224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

7 x 7 x 512

14 x 14 x 512

1 x 1 x 4096  1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

# What do we visualize?

**Option 1:** visualize the **filter** itself

what image feature does this correspond to?



depth: 3

depth: 4

width: 128

height: 128

height: 128

**Option 2:** visualize stimuli that activate a "neuron"



224 x 224 x 3

224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096  1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

electrical signal
recording electrode
visual area of brain
stimulus

What image patch makes the output of this filter **large**?

# Visualizing filters



224 x 224 x 3    224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096   1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

Can't really visualize higher layers in a way that makes much sense

# Visualizing neuron responses

**Idea 1:** look for images that maximally "excite" specific units

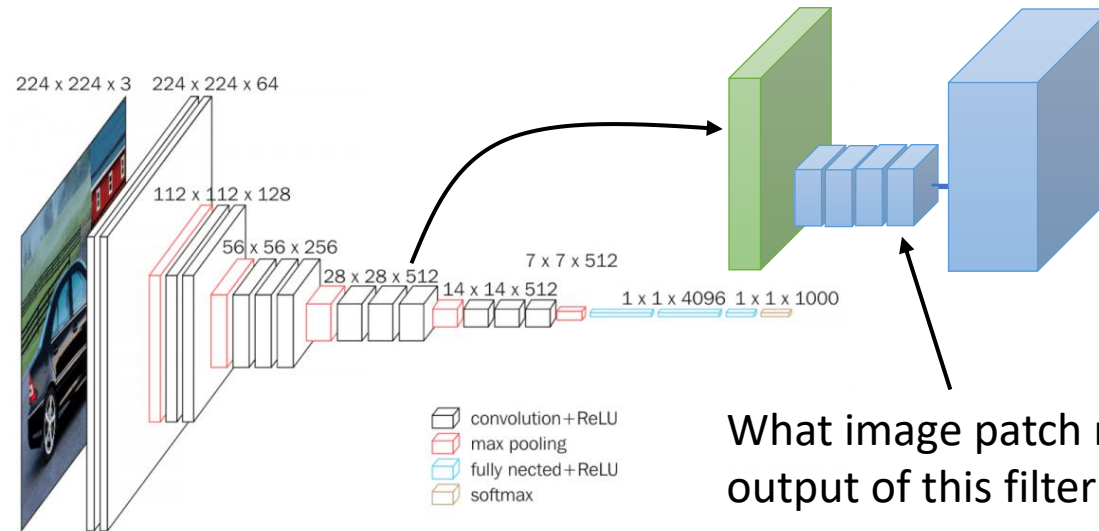12.3 37.4 17.1 21.4 42.1



224 x 224 x 3  224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

7 x 7 x 512

14 x 14 x 512

1 x 1 x 4096  1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

What image patch makes the output of this filter **large**?

# Visualizing neuron responses



**Figure 4: Top regions for six pool$_5$ 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).
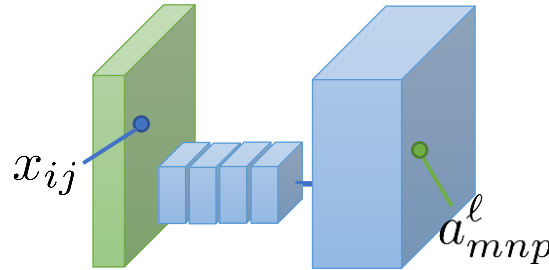
Girshick, Donahue, Darrell, Malik. **Rich feature hierarchies for accurate object detection and semantic segmentation**. 2014

# Using gradients for visualization

**Idea 1:** See which image pixels maximally influence the value at some unit

what does "influence" mean?

given a pixel $x_{ij}$ and a unit $a^\ell_{mnp}$

how much does changing $x_{ij}$ change $a^\ell_{mnp}$?



$\dfrac{da^\ell_{mnp}}{dx_{ij}}$    how do we get this quantity?

        backpropagation!

1. set $\delta$ to be same size as $a^\ell$
2. set $\delta_{mnp} = 1$, all other entries to 0
3. backprop from layer $\ell$ to the image
4. last $\delta$ gives us $\dfrac{da^\ell_{mnp}}{dx}$

forward pass: calculate each $a^{(i)}$ and $z^{(i)}$

backward pass:    $\dfrac{da^\ell_{mnp}}{da^\ell}$ $\leftarrow$ zero in each position except $mnp$ (which is 1)

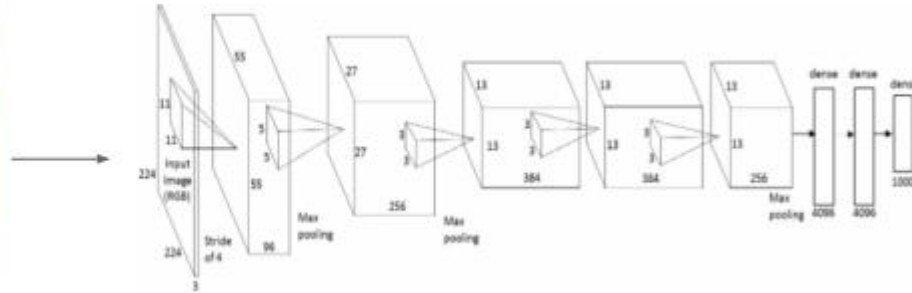initialize $\delta = \dfrac{d\mathcal{L}}{dz^{(i)}}$

for each $f$ with input $x_f$ & params $\theta_f$ from end to start:

$$\dfrac{d\mathcal{L}}{d\theta_f} \leftarrow \dfrac{df}{d\theta_f}\delta$$

$$\delta \leftarrow \dfrac{df}{dx_f}\delta$$

# Using gradients for visualization

**Idea 1:** See which image pixels maximally influence the value at some unit
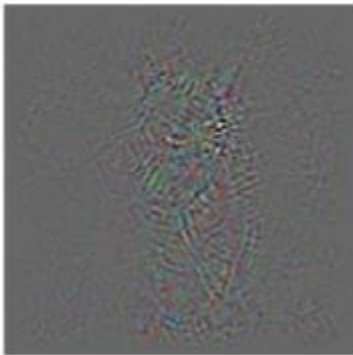
$$\frac{da^{\ell}_{mnp}}{dx_{ij}}$$



basic idea behind "guided backpropagation":

➢ Just using backprop is not very interpretable, because many units in the network contribute **positive** or **negative** gradients

➢ If we keep just the **positive** gradients, we'll avoid some of the complicated negative contributions, and get a cleaner signal

heuristically change backward step in ReLU:

for each entry $\delta_{ijk}$, set it to $\max(0, \delta_{ijk})$

"zero out negative gradients at each layer"

$$\frac{da^{\ell}_{mnp}}{dx_{ij}}$$

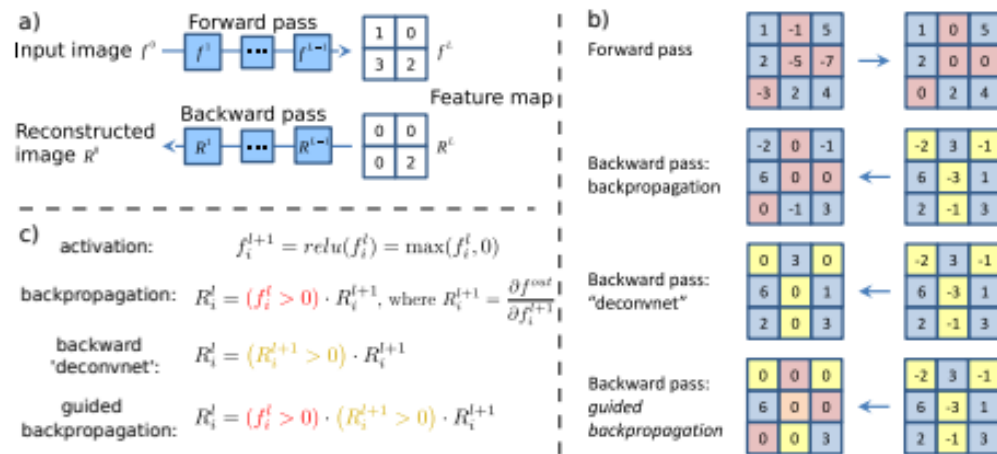slightly modified backprop gives us this:

"guided backpropagation"

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller. **Striving for Simplicity: The All Convolutional Net**. 2014.

Images from CS231n (Fei-Fei Li, Andrej Karpathy)

# Using gradients for visualization

**Idea 1:** See which image pixels maximally influence the value at some unit

$$\frac{da^{\ell}_{mnp}}{dx_{ij}}$$



Figure 1: Schematic of visualizing the activations of high layer neurons. a) Given an input image, we perform the forward pass to the layer we are interested in, then set to zero all activations except one and propagate back to the image to get a reconstruction. b) Different methods of propagating back through a ReLU nonlinearity. c) Formal definition of different methods for propagating a output activation *out* back through a ReLU unit in layer $l$; note that the 'deconvnet' approach and guided backpropagation do not compute a true gradient but rather an imputed version.

basic idea behind "guided backpropagation":

➢ Just using backprop is not very interpretable, because many units in the network contribute **positive** or **negative** gradients

➢ If we keep just the **positive** gradients, we'll avoid some of the complicated negative contributions, and get a cleaner signal

heuristically change backward step in ReLU:

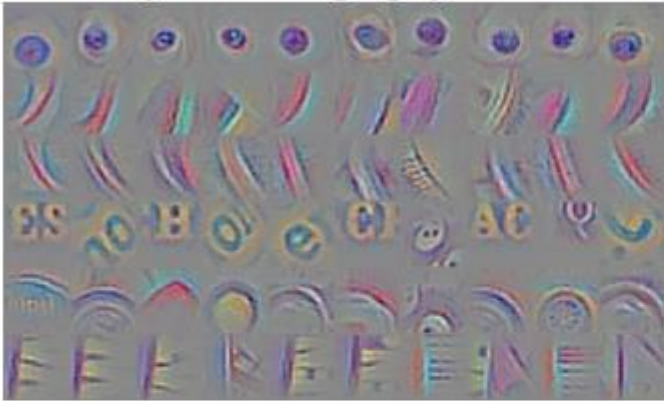for each entry $\delta_{ijk}$, set it to $\max(0, \delta_{ijk})$

"zero out negative gradients at each layer"

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller. **Striving for Simplicity: The All Convolutional Net**. 2014.

# Using gradients for visualization

**Idea 1:** See which image pixels maximally influence the value at some unit



CONV6

CONV9

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller. **Striving for Simplicity: The All Convolutional Net**. 2014.

# Visualizing Features with Backpropagation

# Using gradients for visualization

**Idea 2:** "Optimize" the image to maximally activate a particular unit



Before:   just compute   $\dfrac{da_{mnp}^{\ell}}{dx_{ij}}$

what optimization problem is this solving?

Now:   compute   $g = \dfrac{da_{mnp}^{\ell}}{dx}$

$x \leftarrow \arg\max_{x} a_{mnp}^{\ell}(x)$

$x \leftarrow x + \alpha g$

activation $a_{mnp}^{\ell}$ for image $x$

more generally:   $x \leftarrow \arg\max_{x} S(x)$

some activation or class label

# Using gradients for visualization

**Idea 2:** "Optimize" the image to maximally activate a particular unit

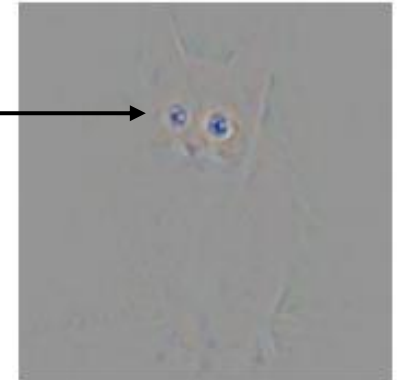$$x \leftarrow \arg\max_{x} S(x) + R(x)$$

problem: it's too easy to produce "crazy" images

making the eyes more blue increases activation

image regularizer to prevent crazy images

what if we the set blue channel to 1,000,000,000?

simple choice: $R(x) = \lambda \|x\|_2^2$

Simonyan, Vidaldi, Zisserman. **Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps.** 2013.

# Visualizing "classes"

$$x \leftarrow \arg\max_{x} S(x) + R(x)$$

simple choice: $R(x) = \lambda ||x||_2^2$
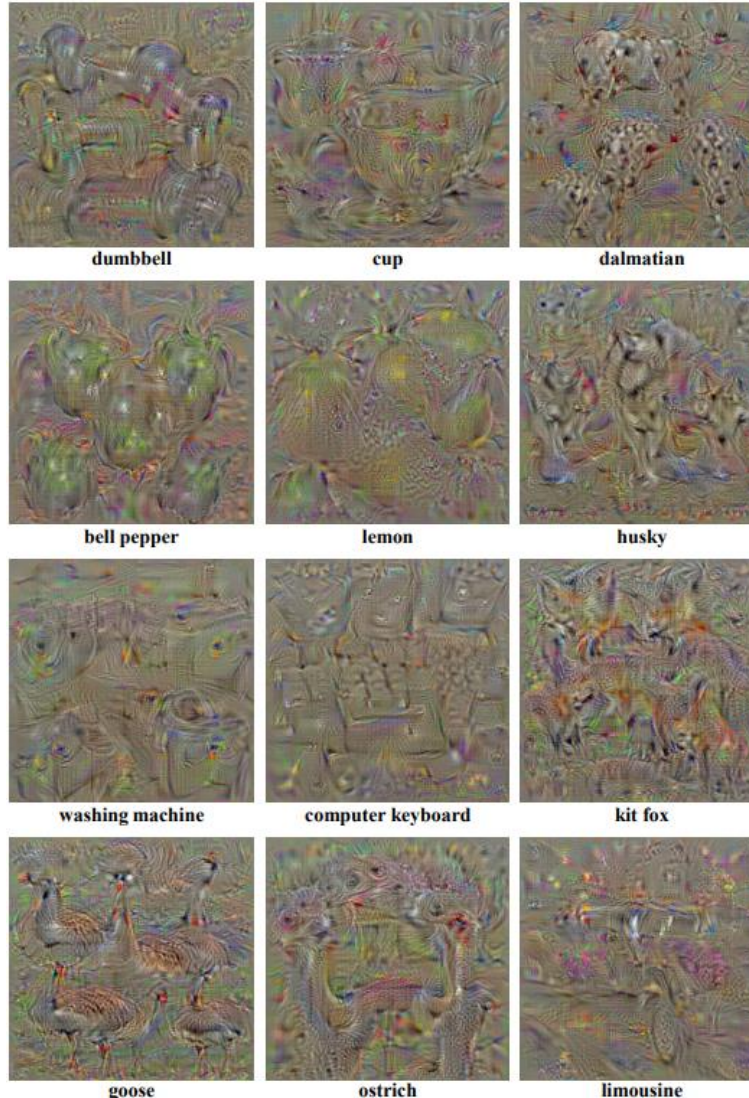
Which unit to maximize?

Let's try maximizing class labels first

**Important detail:** not maximizing **class probabilities**, but the activations **right before** the softmax

$$p(y|x) = \mathrm{softmax}(W^{\ell} a^{\ell}(x) + b^{\ell})$$

not this            this



dumbbell    cup    dalmatian

bell pepper    lemon    husky

washing machine    computer keyboard    kit fox

goose    ostrich    limousine

Simonyan, Vidaldi, Zisserman. **Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps.** 2013.
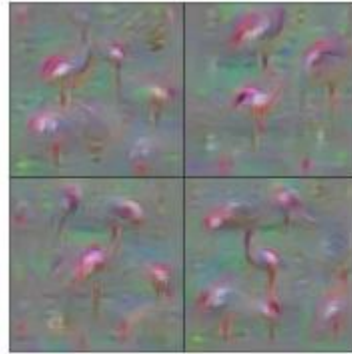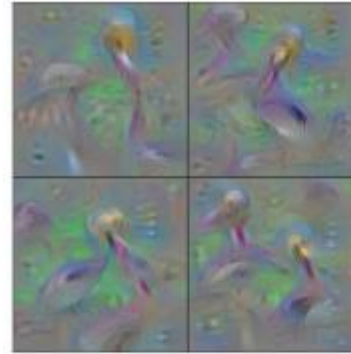
# Visualizing "classes"

$$x \leftarrow \arg\max_x S(x) + R(x)$$
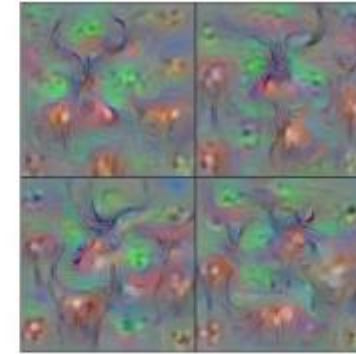
slightly more nuanced regularizer

1. Update image with gradient
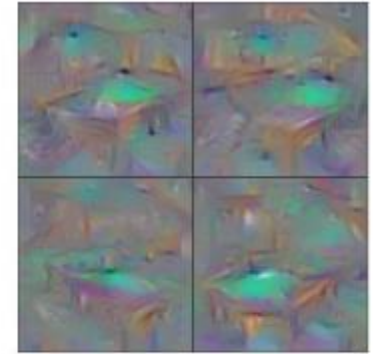2. Blur the image a little
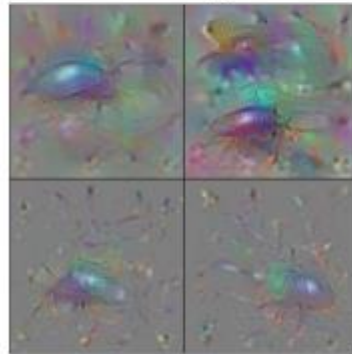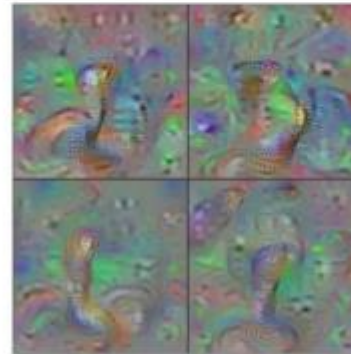3. Zero out any pixel with small value
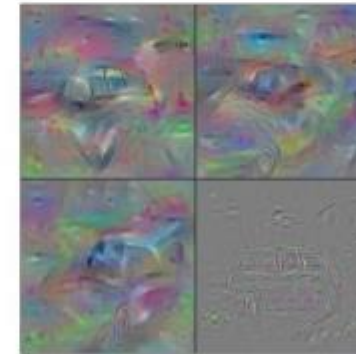4. Repeat



Flamingo

Pelican

Hartebeest

Billiard Table

Ground Beetle

Indian Cobra

Station Wagon

Black Swan

Yosinski et al. **Understanding Neural Networks Through Deep Visualization.** 2015.

# Visualizing units

$$x \leftarrow \arg\max_x S(x) + R(x)$$

slightly more nuanced regularizer

1. Update image with gradient
2. Blur the image a little
3. Zero out any pixel with small value
4. Repeat



Layer 8

Pirate Ship    Rocking Chair    Teddy Bear    Windsor Tie    Pitcher

Layer 7

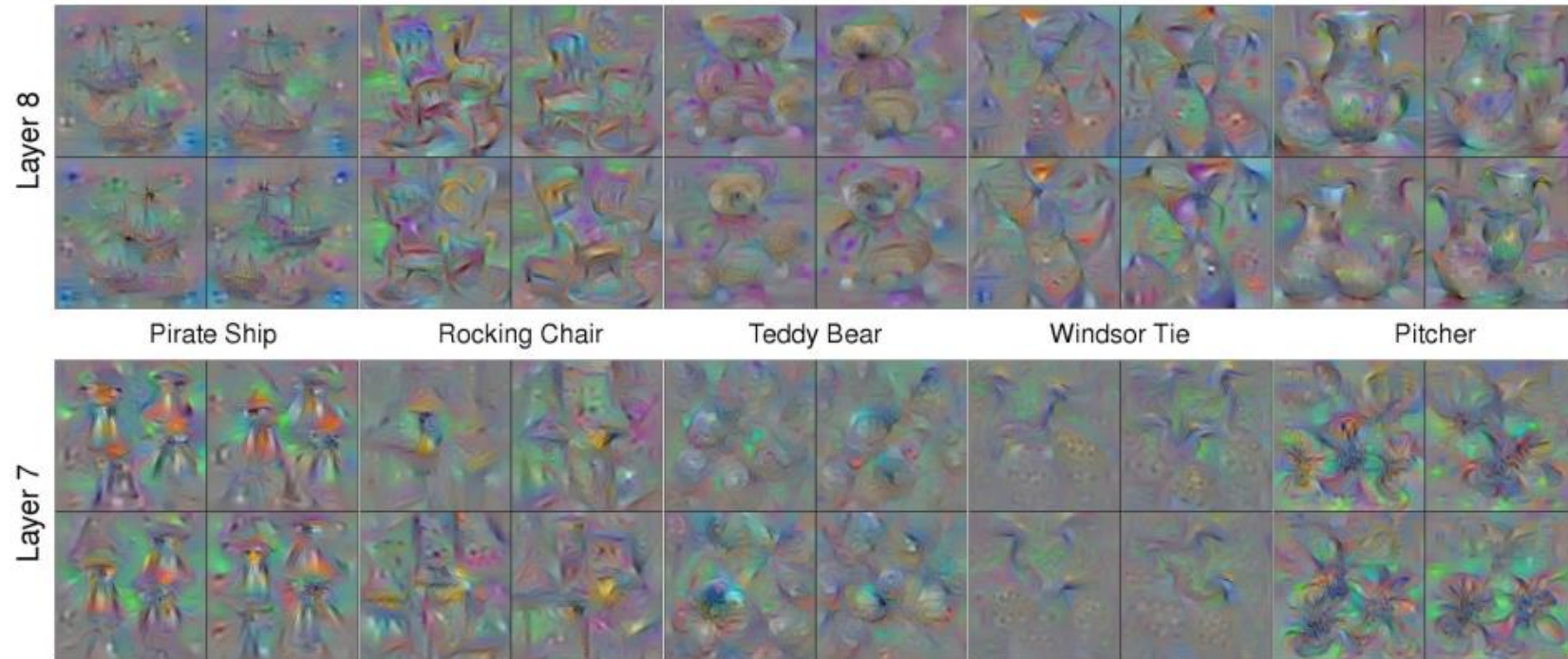Yosinski et al. **Understanding Neural Networks Through Deep Visualization.** 2015.

# Visualizing units

$$x \leftarrow \arg\max_{x} S(x) + R(x)$$

slightly more nuanced
regularizer

1. Update image with gradient
2. Blur the image a little
3. Zero out any pixel with small value
4. Repeat



Slide based on CS231n, Fei-Fei Li & Andrej Karpathy

Yosinski et al. **Understanding Neural Networks Through Deep Visualization.** 2015.

# Visualizing units

$$x \leftarrow \arg\max_{x} S(x) + R(x)$$

slightly more nuanced
regularizer

1. Update image with gradient
2. Blur the image a little
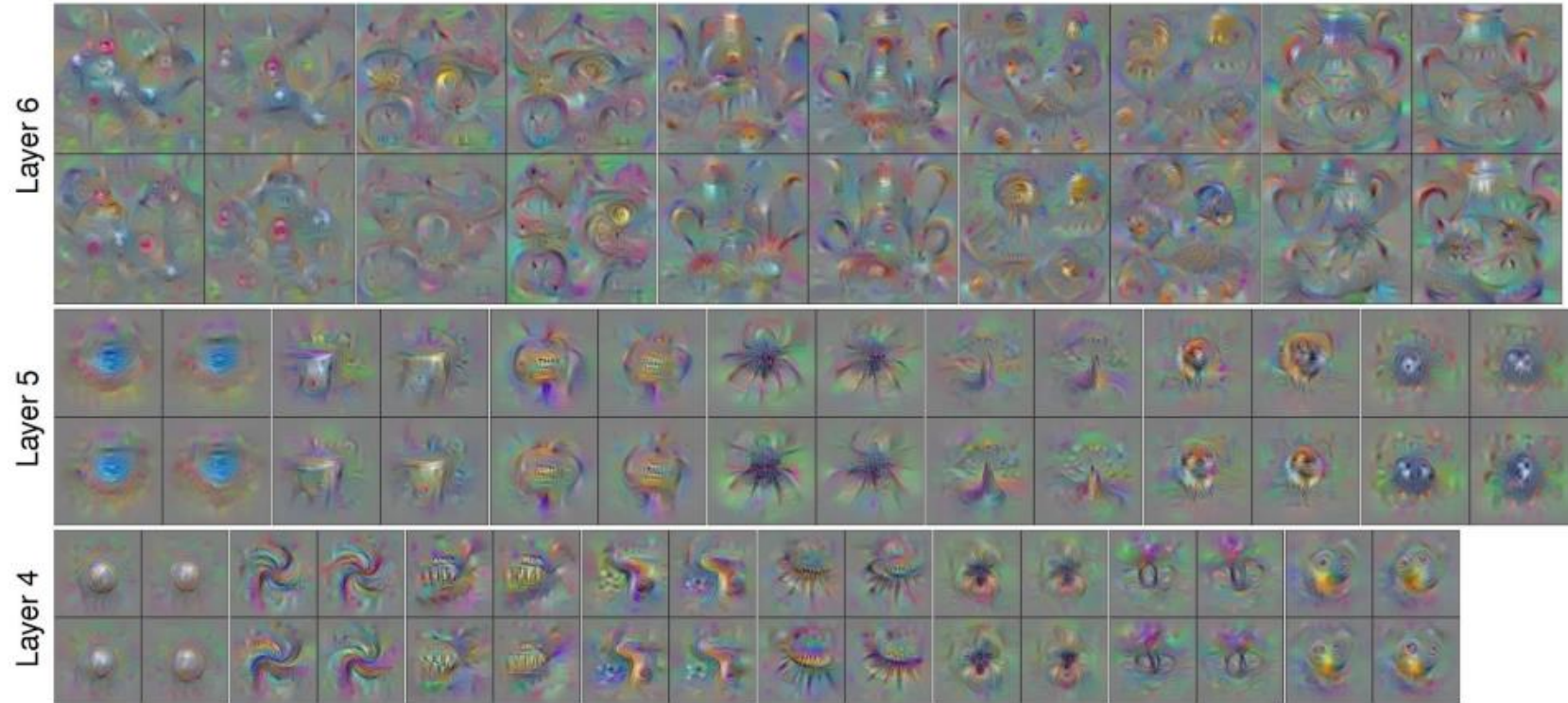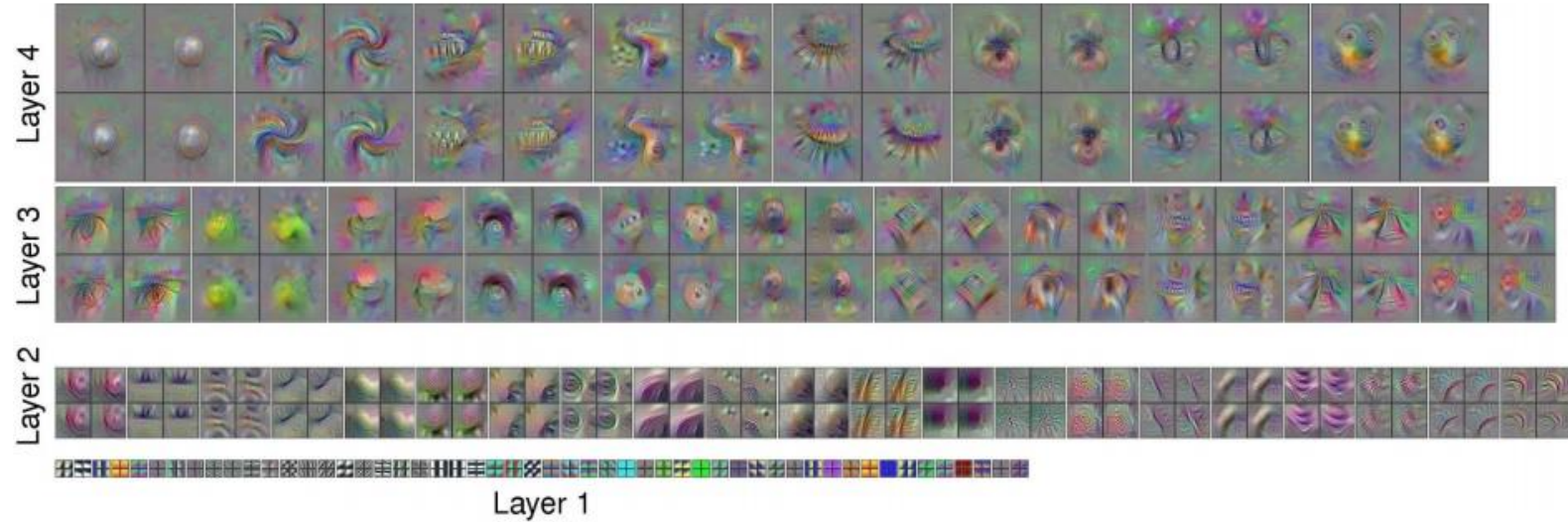3. Zero out any pixel with small value
4. Repeat

# Deep Dream & Style Transfer

# Using backprop to *modify* pictures?

Before:

Now:



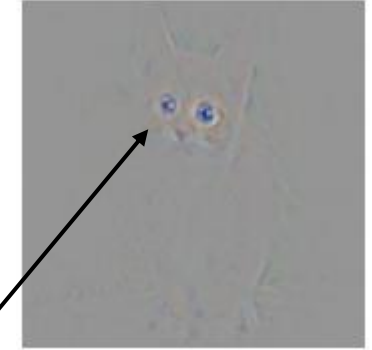What is channel 17 in layer conv5 looking at?



What would it look like if it were a Monet painting?



How are these questions related?

**Intuition:** Monet paintings have particular feature distributions, we can "transport" these distributions to other images!

??

# Looking for patterns in clouds



making the eyes more blue will make this more cat-like

224 x 224 x 3  224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096  1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

"looks kind of like a dog"

"where is the dog?"

"let me show you!"

# DeepDream



inception_4c/output

1. Pick a layer

2. Run forward pass to compute activations at that layer

3. Set delta to be **equal** to the activations

4. Backprop and apply the gradient

5. Repeat

Images from CS231n (Fei-Fei Li, Andrej Karpathy)

# DeepDream

inception_4c/output



Images from CS231n (Fei-Fei Li, Andrej Karpathy)

# DeepDream

```python
def objective_L2(dst):
    dst.diff[:] = dst.data

def make_step(net, step_size=1.5, end='inception_4c/output',
              jitter=32, clip=True, objective=objective_L2):
    '''Basic gradient ascent step.'''

    src = net.blobs['data'] # input image is stored in Net's 'data' blob
    dst = net.blobs[end]

    ox, oy = np.random.randint(-jitter, jitter+1, 2)
    src.data[0] = np.roll(np.roll(src.data[0], ox, -1), oy, -2) # apply jitter shift

    net.forward(end=end)
    objective(dst)  # specify the optimization objective
    net.backward(start=end)
    g = src.diff[0]
    # apply normalized ascent step to the input image
    src.data[:] += step_size/np.abs(g).mean() * g

    src.data[0] = np.roll(np.roll(src.data[0], -ox, -1), -oy, -2) # unshift image

    if clip:
        bias = net.transformer.mean['data']
        src.data[:] = np.clip(src.data, -bias, 255-bias)
```

DeepDream:  set dx = x :)

jitter regularizer

"image update"

Remember this:

$$x \leftarrow \arg\max_{x} S(x) + R(x)$$

# DeepDream



DeepDream https://github.com/google/deepdream
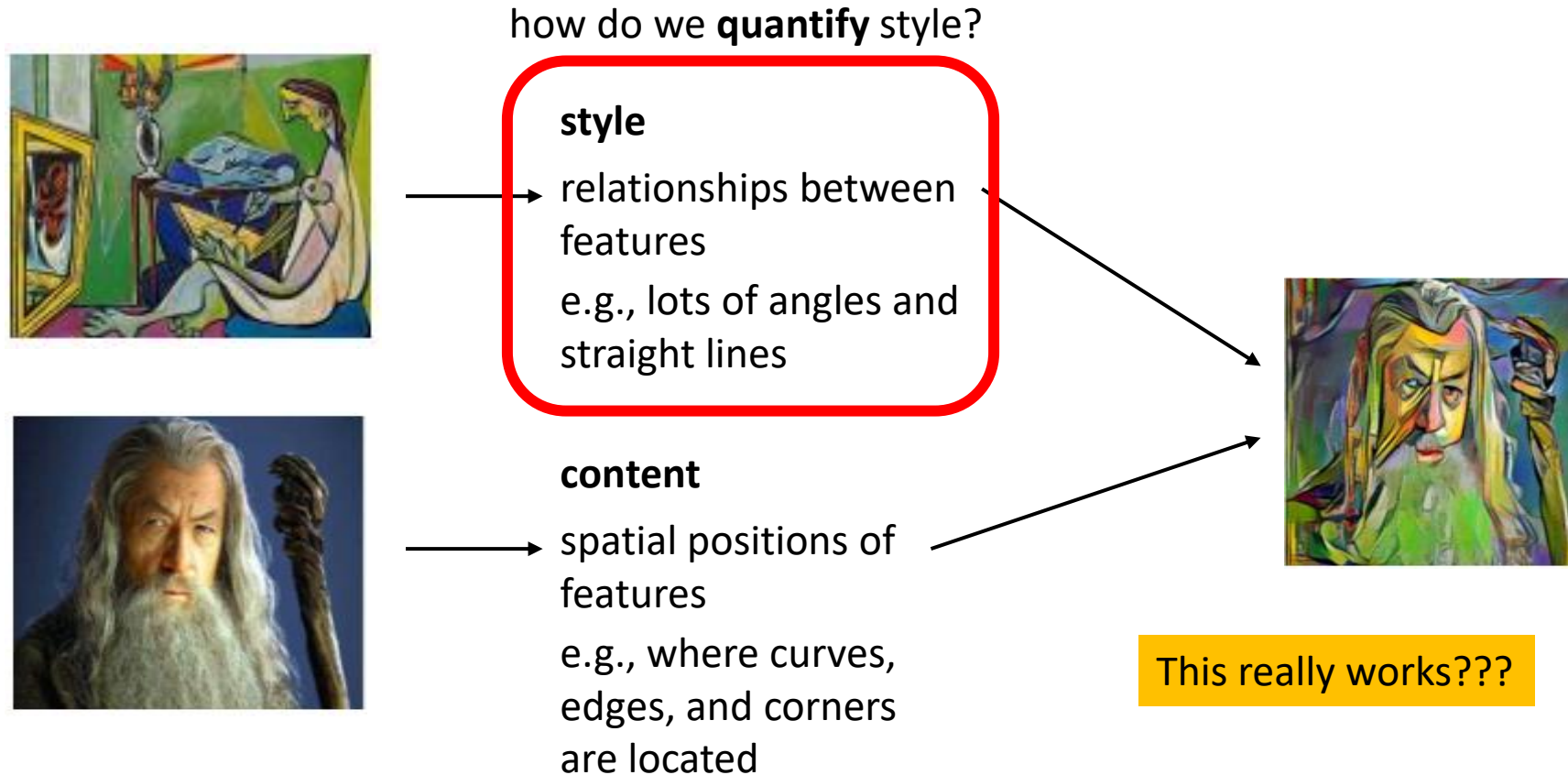
# Another idea

**Another idea:** instead of exaggerating the features in a single image, what if we make feature of one image look more like the features in **another**?
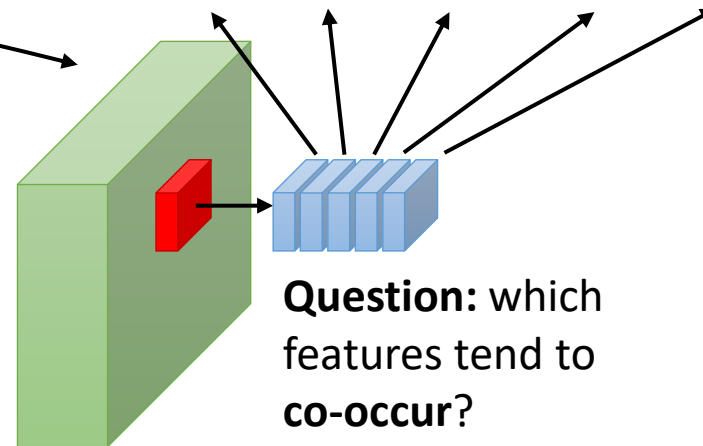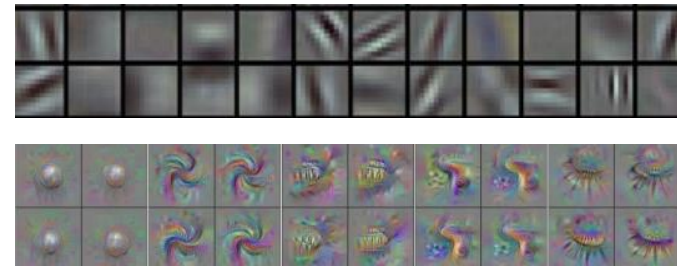
how do we **quantify** style?



**style**

relationships between features

e.g., lots of angles and straight lines

**content**

spatial positions of features

e.g., where curves, edges, and corners are located

This really works???

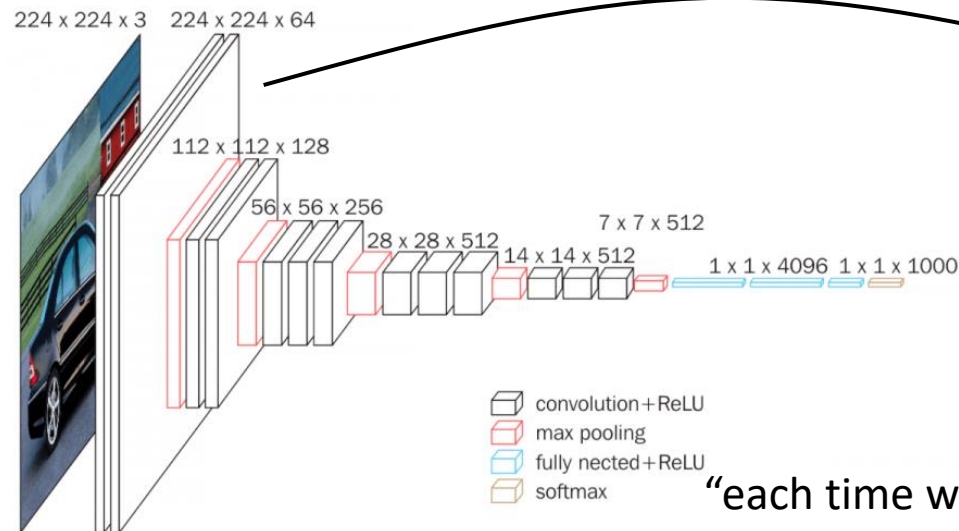# How do we quantify **style**?

**style**

relationships between features

e.g., lots of angles and straight lines



**Question:** which features tend to **co-occur**?

"each time we see a straight line, we also see this texture"

"each time we see a curve, we also see flat shading"

How do we quantify this?
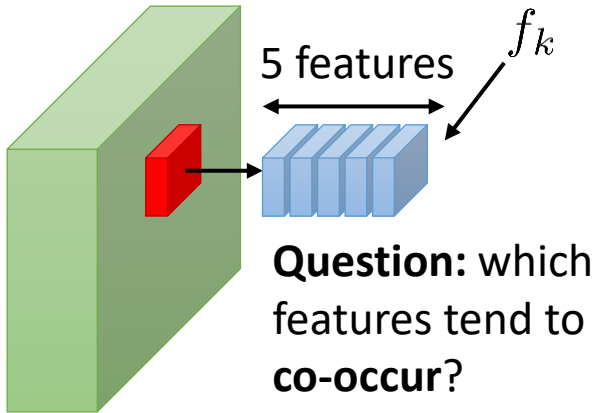
# How do we quantify **style**?



**style**

relationships between features

e.g., lots of angles and straight lines

estimate by averaging over all positions in the image

5 features

$f_k$

feature covariance: $\mathrm{Cov}_{km} = E[f_k f_m]$

form Gram matrix: $G_{km} = \mathrm{Cov}_{km}$

If features have this covariance, then we have the right style

Surprising but true!

**Question:** which features tend to **co-occur**?

"each time we see a straight line, we also see this texture"

"each time we see a curve, we also see flat shading"

# How do we quantify **style**?

**style**

relationships between features

e.g., lots of angles and straight lines

this comes from the style source image

feature covariance: $\text{Cov}_{km} = E[f_k f_m]$

form Gram matrix: $G_{km} = \text{Cov}_{km}$

new image: $x \leftarrow \arg \min_x \mathcal{L}_{\text{style}}(x) + \mathcal{L}_{\text{content}}(x)$

$G^\ell$: source image Gram matrix at layer $\ell$

$A^\ell(x)$: new image Gram matrix at layer $\ell$

$$\mathcal{L}_{\text{style}}(x) = \sum_\ell \sum_{km} \left( G^\ell_{km} - A^\ell_{km}(x) \right)^2 w_\ell$$

Different weight on each layer, to prioritize relative contribution to (desired) style of different levels of abstraction
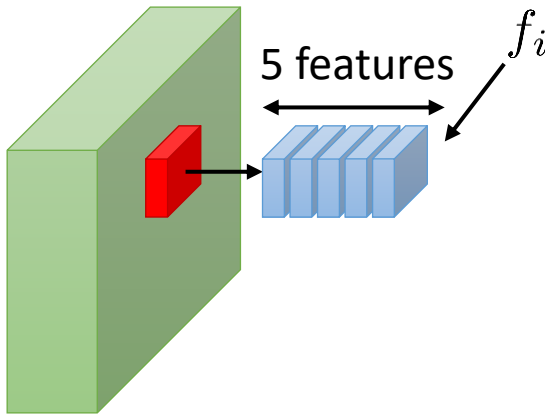
# How do we quantify **content**?



**content**

spatial positions of features

e.g., where curves, edges, and corners are located

Pick specific layer for matching content

5 features    $f_i$

Just directly match the features!

$$\mathcal{L}_{\text{content}}(x) = \sum_{ij} \sum_{k} \left( f^{\ell}_{ijk}(x_{\text{content}}) - f^{\ell}_{ijk}(x) \right)^2$$

$$\mathcal{L}_{\text{style}}(x) = \sum_{\ell} \sum_{km} \left( G^{\ell}_{km} - A^{\ell}_{km}(x) \right)^2 w_{\ell}$$
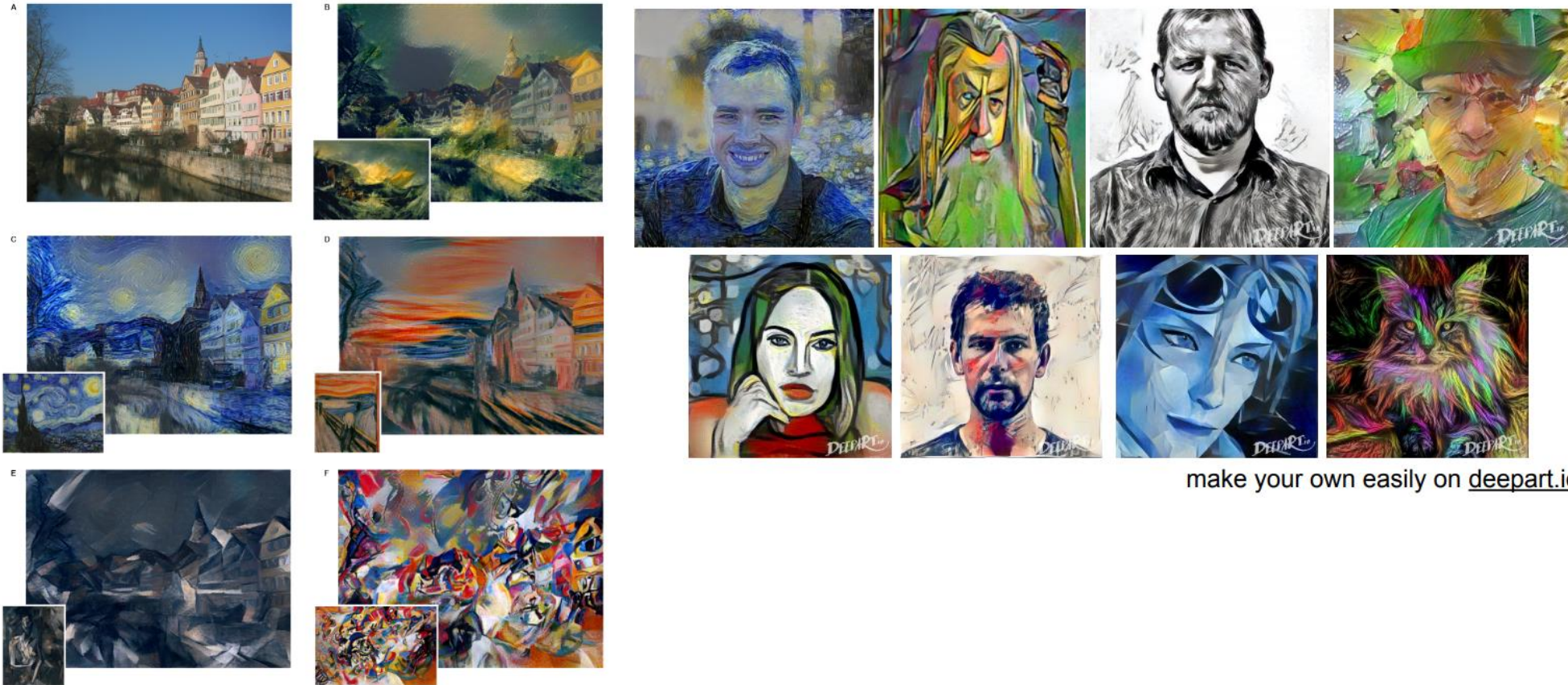
new image: $x \leftarrow \arg\min_x \mathcal{L}_{\text{style}}(x) + \mathcal{L}_{\text{content}}(x)$

# Style transfer

new image: $x \leftarrow \arg\min_x \mathcal{L}_{\text{style}}(x) + \mathcal{L}_{\text{content}}(x)$

$$\mathcal{L}_{\text{content}}(x) = \sum_{ij}\sum_{k}\left(f^{\ell}_{ijk}(x_{\text{content}}) - f^{\ell}_{ijk}(x)\right)^2$$

$$\mathcal{L}_{\text{style}}(x) = \sum_{\ell}\sum_{km}\left(G^{\ell}_{km} - A^{\ell}_{km}(x)\right)^2 w_{\ell}$$



make your own easily on deepart.io

Gatys et al. A Neural Algorithm of Artist Style, 2015.

## CSW182 (2021)· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

**视频·B 站 [ 扫码或点击链接 ]**
*https://www.bilibili.com/video/BV1Ff4y1n7ar*

**课件 & 代码·博客 [ 扫码或点击链接 ]**
*http://blog.showmeai.tech/berkeley-csw182*

Berkeley
循环神经网络
可视化
梯度策略
Q-Learning
风格迁移
模仿学习
元学习
计算机视觉
机器学习基础
生成模型
卷积网络

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS23In |

### # Awesome AI Courses Notes Cheatsheets· 持续更新中

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏

称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]