

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

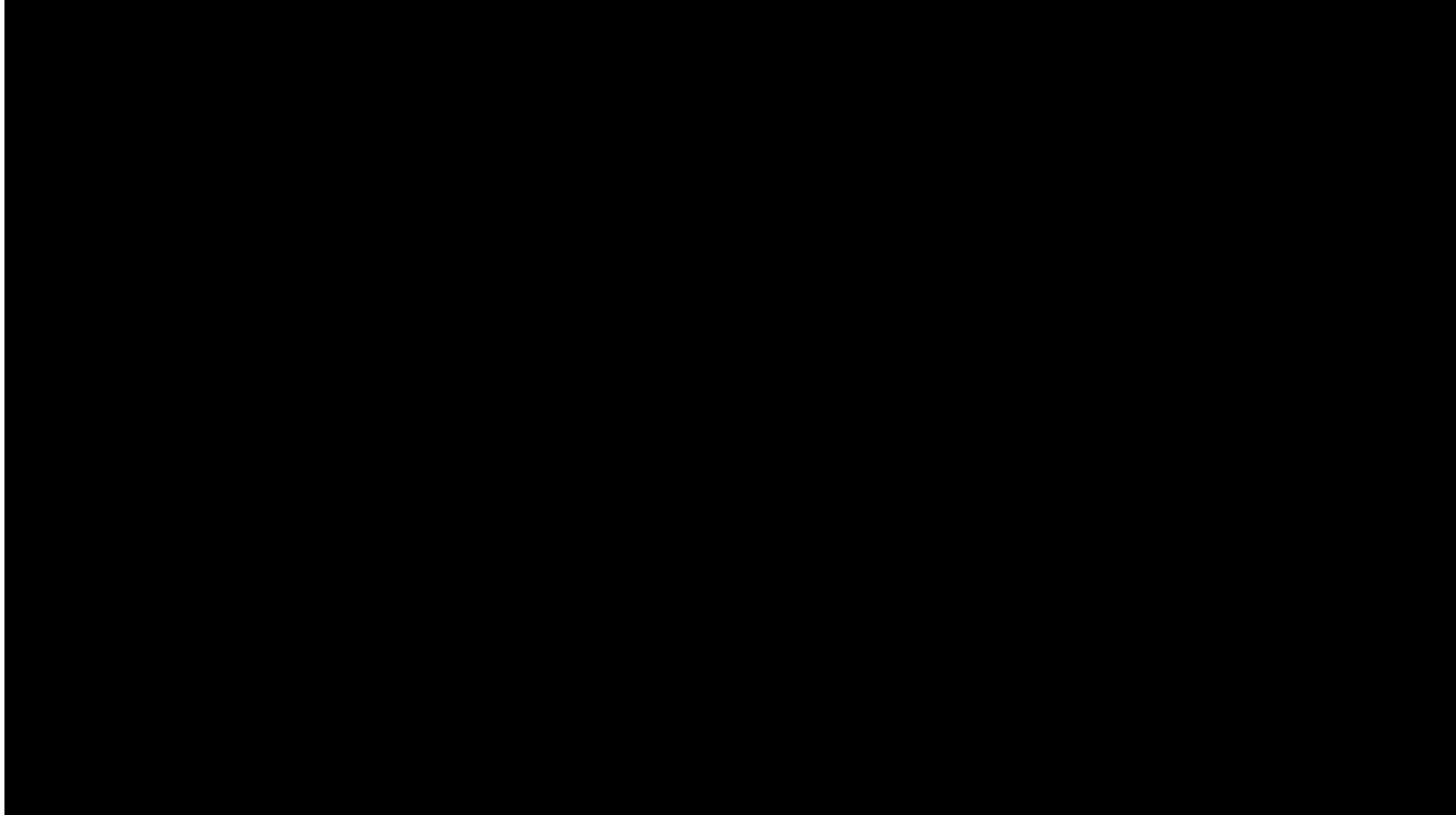


deeplearning.ai

Face recognition

What is face
recognition?

Face recognition



Face verification vs. face recognition

Face verification is a building block for face recognition and is "easier" than face recognition.

→ Verification

- Input image, name/ID
- Output whether the input image is that of the claimed person

1:1

99%

99.9
~~~

## → Recognition

- Has a database of K persons
- Get an input image
- Output ID if the image is any of the K persons (or "not recognized")

1:K

K=100 ←

To perform face recognition in a scalable way the best option is to learn a similarity function between pairs of faces. The idea is:

- 1- Learn a similarity function on a certain dataset (not test dataset)
- 2- at runtime compare the test image with all the faces of the employee and see if one matches.



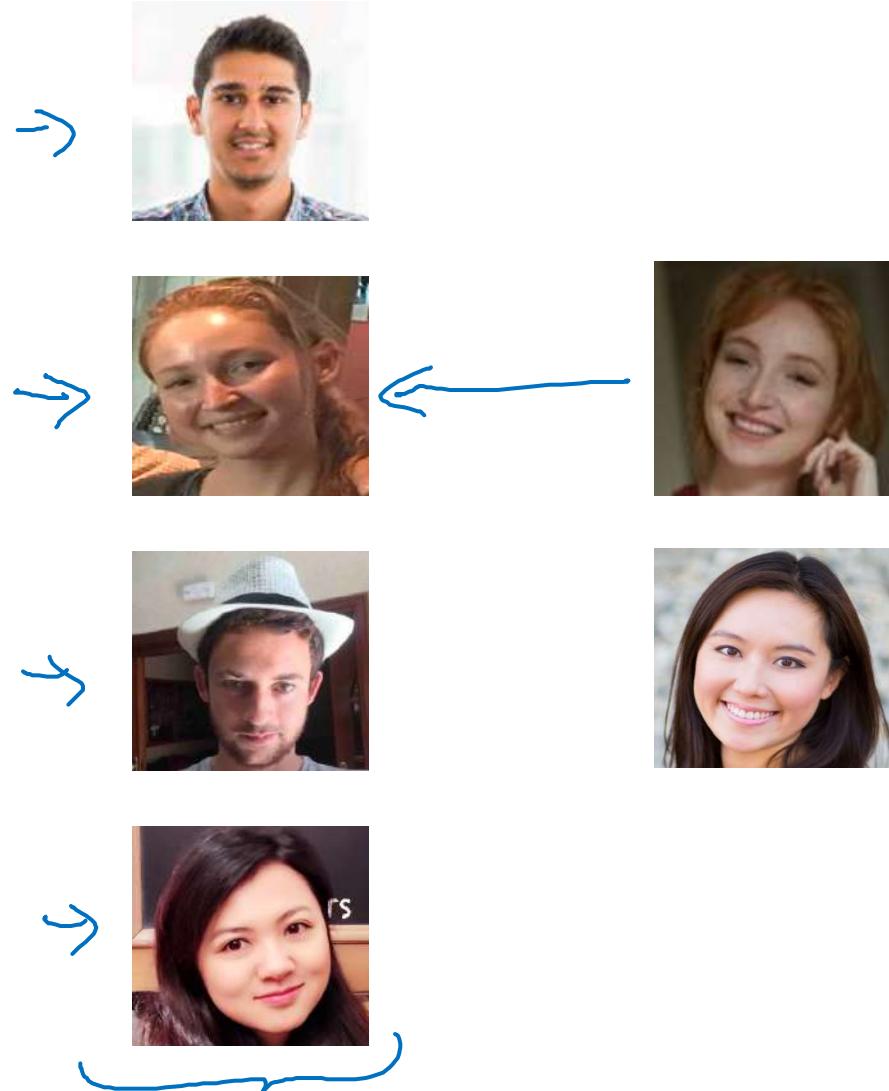
## Face recognition

This enables 1-shot learning because i can use my pretrained "similarity network" to compare pairs of my target images.

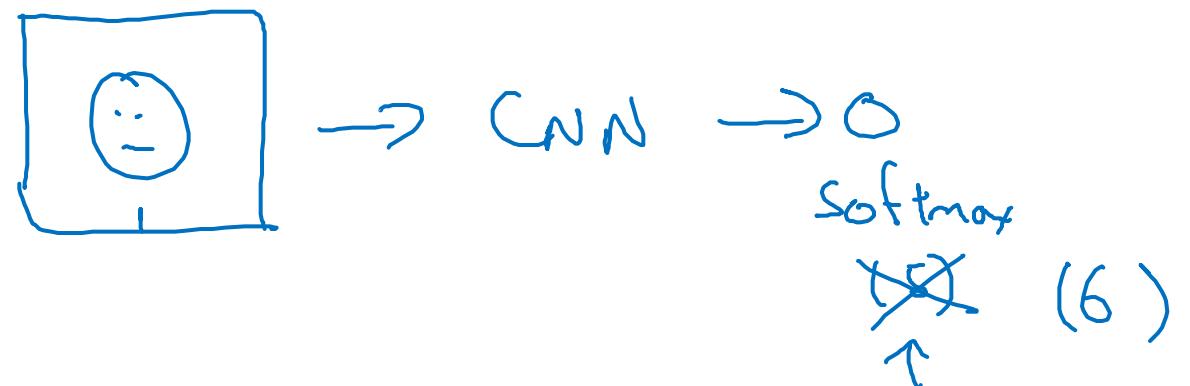
## One-shot learning

deeplearning.ai

# One-shot learning



Learning from one example to recognize the person again



# Learning a “similarity” function

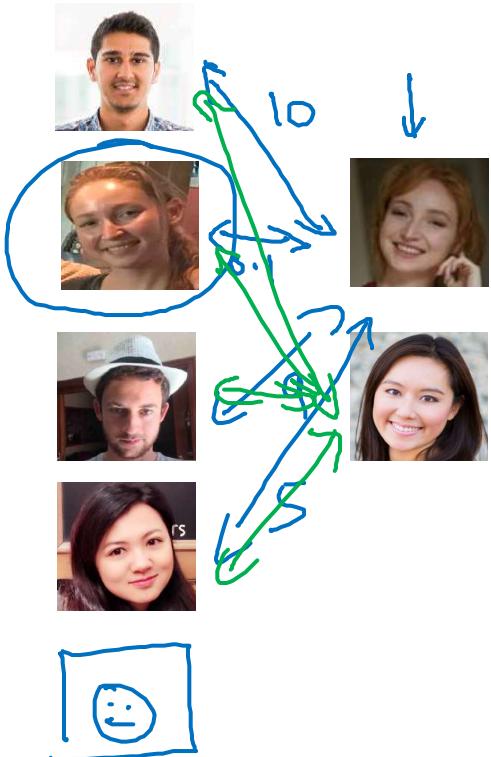
→  $d(\underline{\text{img1}}, \underline{\text{img2}})$  = degree of difference between images

If  $d(\text{img1}, \text{img2}) \leq \tau$

$> \tau$

“some”  
“different”

Verification.



$d(\text{img1}, \text{img2})$

To learn a face similarity function one uses Siamese networks. These networks have  $n$  replicated branches with shared weights, which are joined at the end by a similarity loss. In case of face recognition



deeplearning.ai

## Face recognition

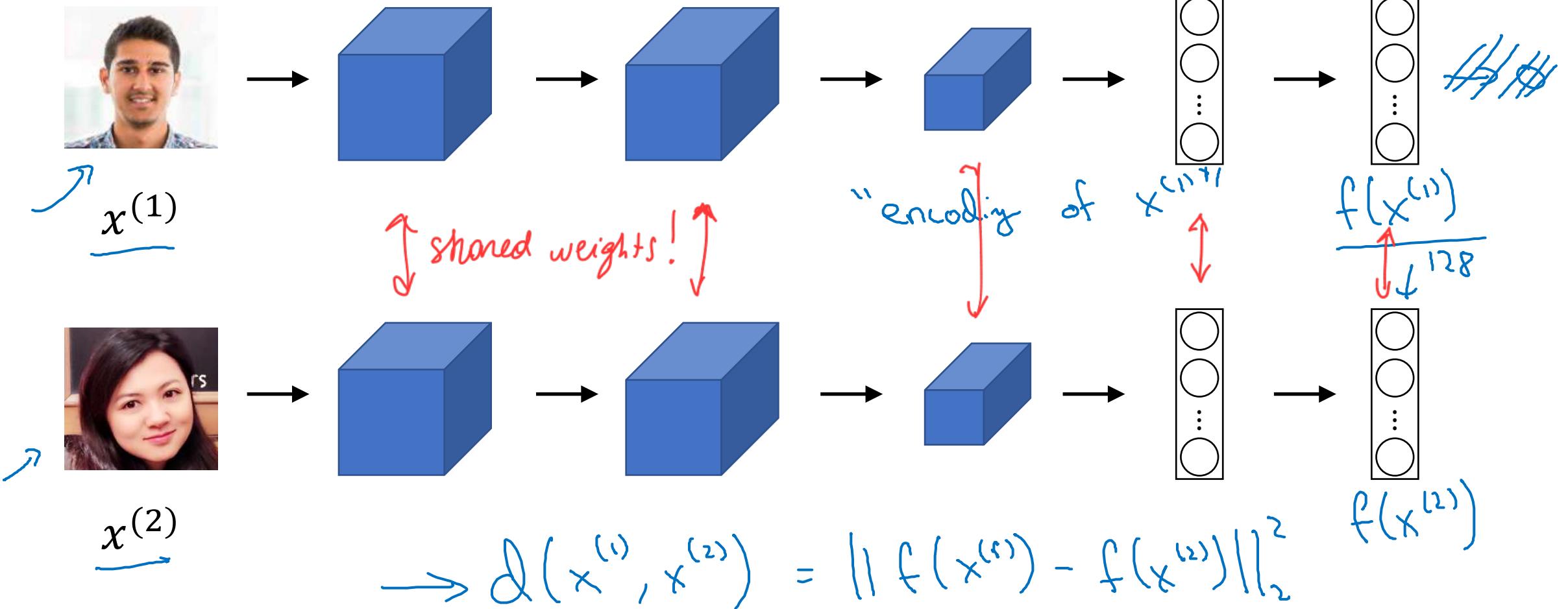
• "triplet loss" is typically used to learn high similarity (small distance) between two images of the same person and low similarity compared to a third image of another person. Of the 3 images:

## Siamese network

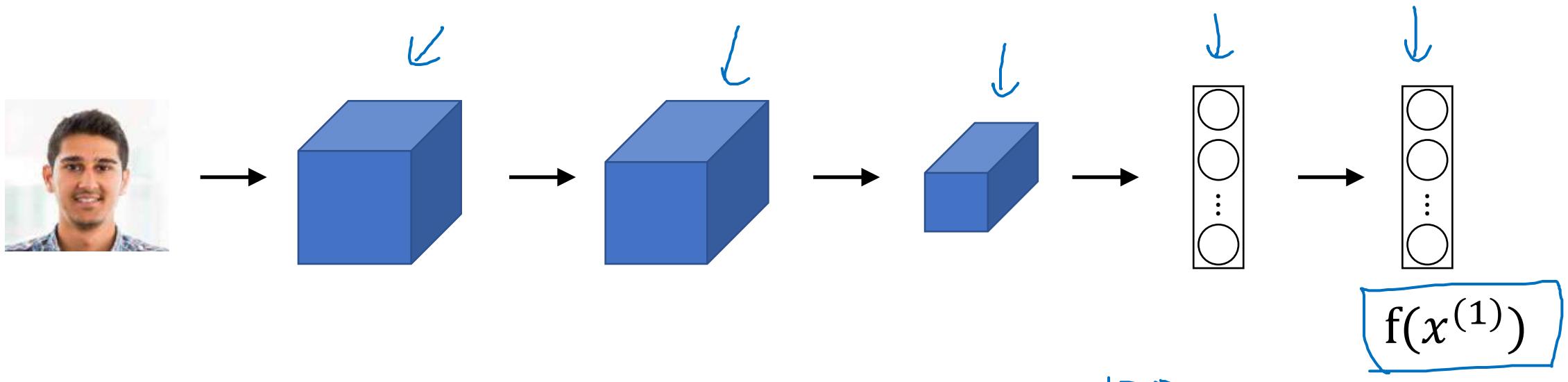
- one "anchor" image belongs to the target person
- one "positive" image belongs to the same person
- one "negative" image belongs to another person.

# Siamese network

NOTE: When using triplet loss you need 3 branches.



# Goal of learning



Parameters of NN define an encoding  $f(x^{(i)})$  128

Learn parameters so that:

If  $x^{(i)}, x^{(j)}$  are the same person,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is small.

If  $x^{(i)}, x^{(j)}$  are different persons,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is large.

The triplet loss can be used as a loss function in siamese networks to push the networks to learn a similarity function between couples of images. This loss takes 3 images at a time :

- an anchor = the person to be recognized
- a positive example = same person
- a negative example = other person



# Face recognition

---

## Triplet loss

deeplearning.ai

# Learning Objective



idea: you want that  $d(A, P) \leq d(A, N)$ .

- to avoid the trivial solution (all encodings = 0) actually impose a margin between the two values  $d(A, P) + \alpha \leq d(A, N)$ .

| <u>Anchor</u>                                                                                      | <u>Positive</u>                          | <u>Negative</u>                                       |
|----------------------------------------------------------------------------------------------------|------------------------------------------|-------------------------------------------------------|
| <u>A</u>                                                                                           | <u><math>P</math></u><br>$d(A, P) = 0.5$ | <u><math>N</math></u><br>$d(A, N) = 0.5$              |
| $\frac{\ f(A) - f(P)\ ^2}{d(A, P)} + \alpha \leq 0.2$                                              |                                          | $\frac{\ f(A) - f(N)\ ^2}{d(A, N)} + \alpha \leq 0.7$ |
| Want: $\frac{\ f(A) - f(P)\ ^2}{d(A, P)} + \alpha \leq \frac{\ f(A) - f(N)\ ^2}{d(A, N)} + \alpha$ |                                          |                                                       |

$$\frac{\|f(A) - f(P)\|^2}{\alpha} - \frac{\|f(A) - f(N)\|^2}{\alpha} + \alpha \leq 0$$

Margin

$f(\text{img}) = \vec{0}$

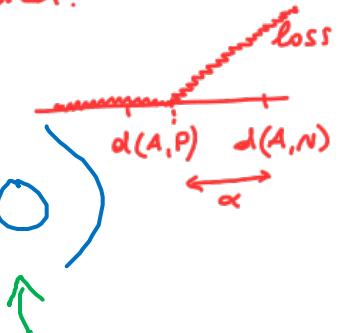
# Loss function

Given 3 images

$A, P, N$ :

The triplet loss uses a max operation to only "consider" cases where  $d(A, P)$  is less than a margin away from  $d(A, N)$  and ignore cases when the margin is already respected.

$$L(A, P, N) = \max \left( \max_{\epsilon > 0} \left( \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \epsilon \right), 0 \right)$$



$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

$A, P$   
↑ ↑

Training set:  $\frac{10k}{\text{C}}$  pictures of  $\frac{1k}{\text{C}}$  persons

NOTE: The SIAMESE NETWORK HAS 3 INPUT BRANCHES.

# Choosing the triplets A,P,N



During training, if A,P,N are chosen randomly,

$d(A, P) + \alpha \leq d(A, N)$  is easily satisfied.

$$\underbrace{\|f(A) - f(P)\|^2}_{\|f(A) - f(N)\|^2} + \alpha \leq \underbrace{\|f(A) - f(N)\|^2}$$

because  $d(A, N)$  is more likely  
to be much greater than  $d(A, P)$

Choose triplets that're “hard” to train on.

for which  $d(A, N) \approx d(A, P)$

Face Net

Deep Face

$$\underbrace{d(A, P)}_{\downarrow} + \alpha \approx \underbrace{d(A, N)}_{\uparrow}$$

details in this paper



# Training set using triplet loss

Anchor



Positive



Negative



:

:

:



J

$d(x^{(i)}, x^{(j)})$

Alternatively to using a triplet loss, one can add a final logistic layer to the network that treats the outputs of two FCs as features for a binary classification problem (same person / different person).



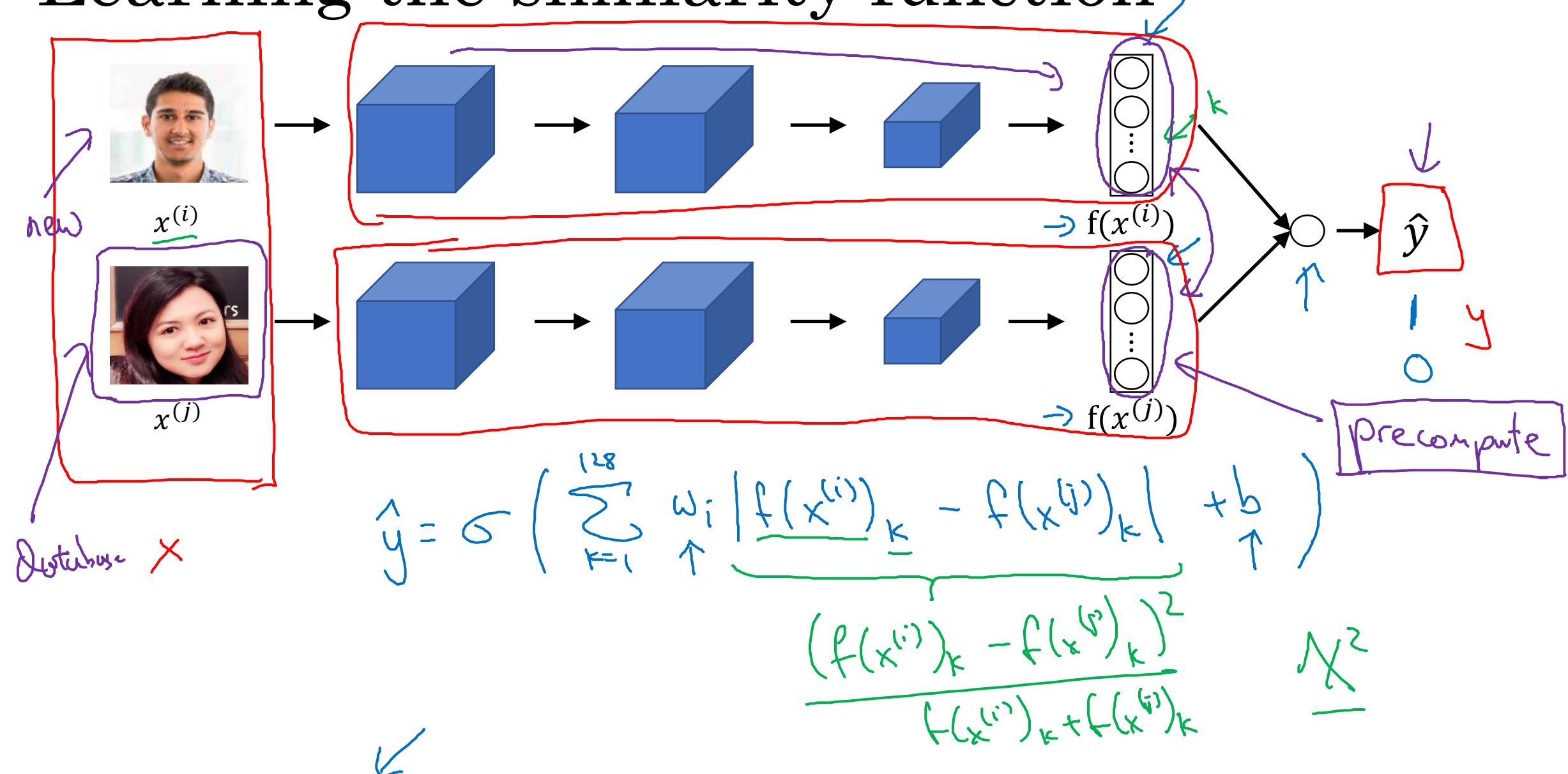
deeplearning.ai

## Face recognition

---

Face verification and  
binary classification

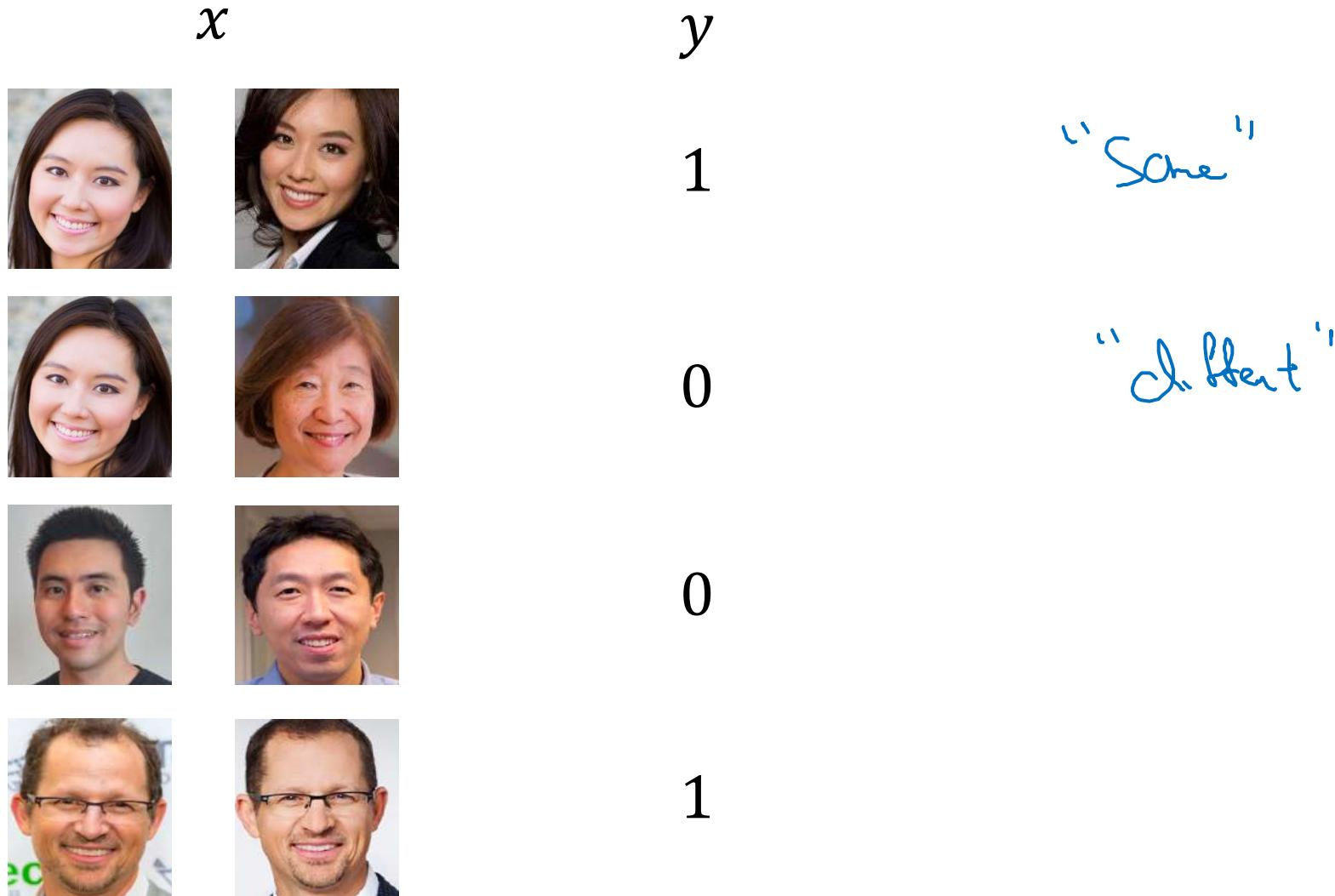
# Learning the similarity function



[Taigman et. al., 2014. DeepFace closing the gap to human level performance]

Andrew Ng

# Face verification supervised learning





deeplearning.ai

# Neural Style Transfer

---

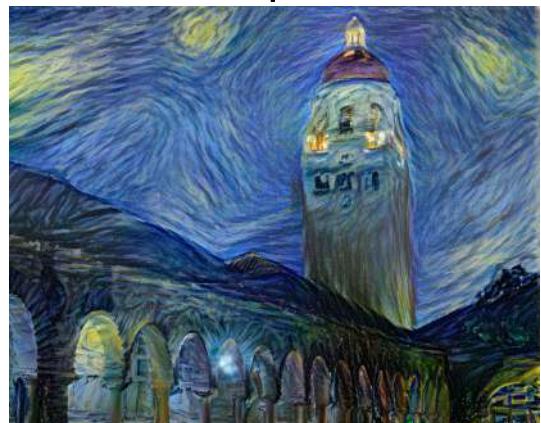
What is neural style  
transfer?

# Neural style transfer



Content ( $c$ )

Style ( $s$ )



Generated image ( $g$ )



Content ( $c$ )

Style ( $s$ )



Generated image ( $g$ )



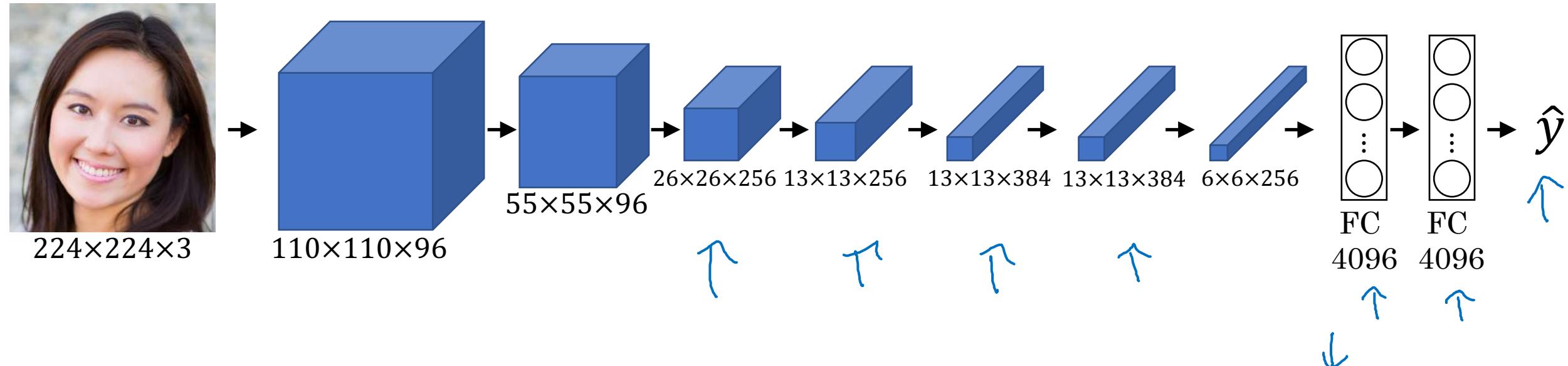
deeplearning.ai

# Neural Style Transfer

---

What are deep  
ConvNets learning?

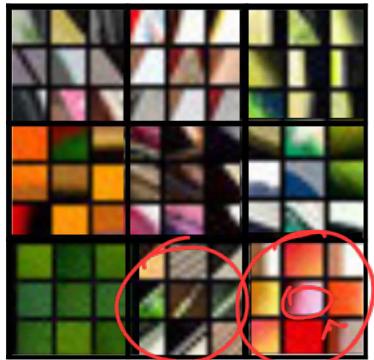
# Visualizing what a deep network is learning



Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation. → visualize all 9 image patches .

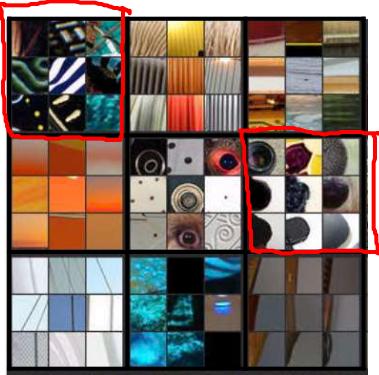
Repeat for other units.

# Visualizing deep layers



Layer 1  
2 neurons of Layer 1

the corresponding patch  
in one of the 9 images  
of the dataset that mostly  
activated that neuron



Layer 2



Layer 3

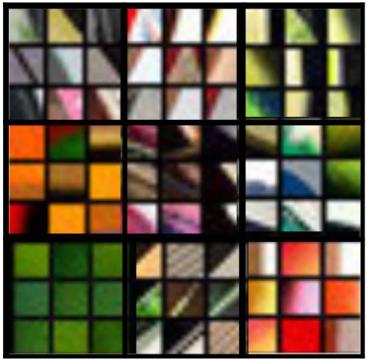


Layer 4



Layer 5

# Visualizing deep layers: Layer 1



Layer 1



Layer 2



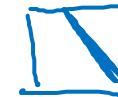
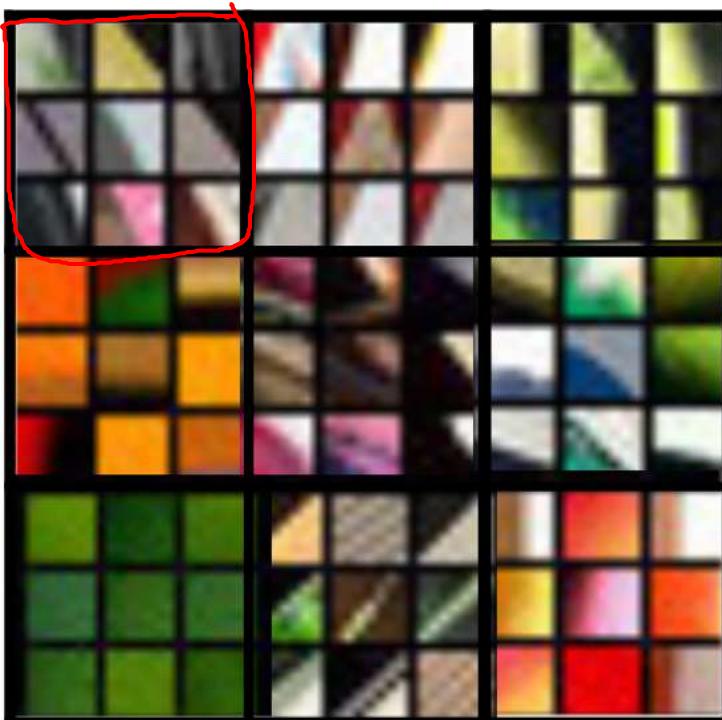
Layer 3



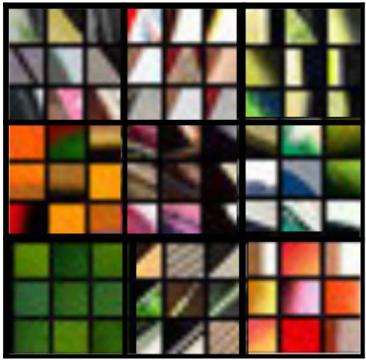
Layer 4



Layer 5



# Visualizing deep layers: Layer 2



Layer 1



Layer 2



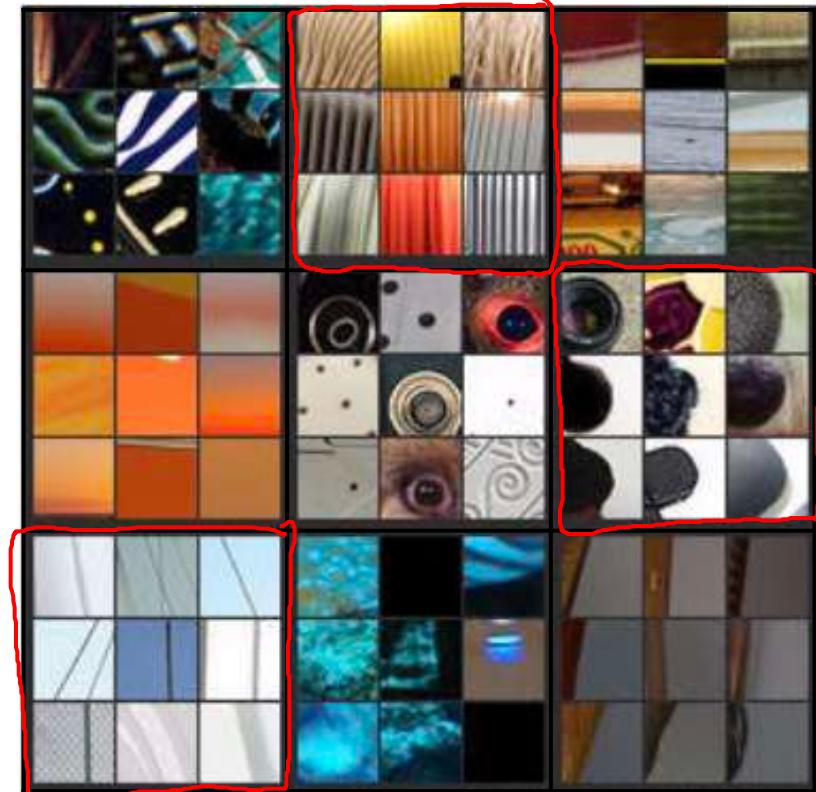
Layer 3



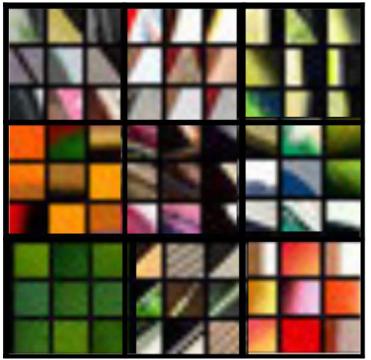
Layer 4



Layer 5



# Visualizing deep layers: Layer 3



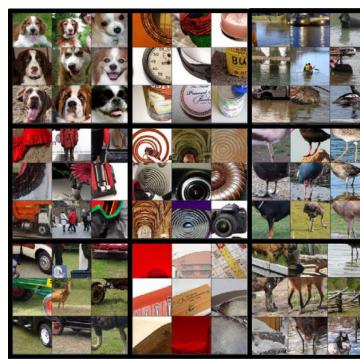
Layer 1



Layer 2



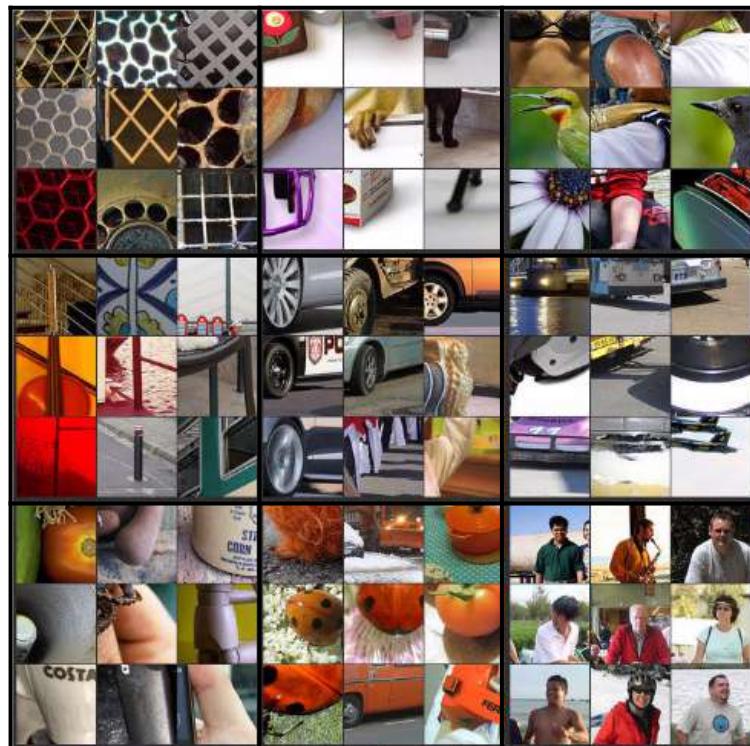
Layer 3



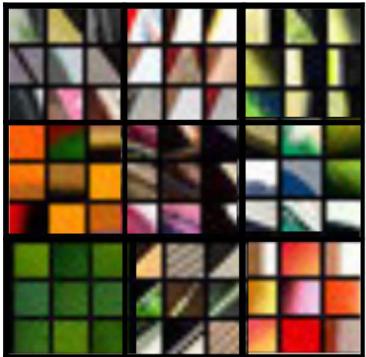
Layer 4



Layer 5



# Visualizing deep layers: Layer 3



Layer 1

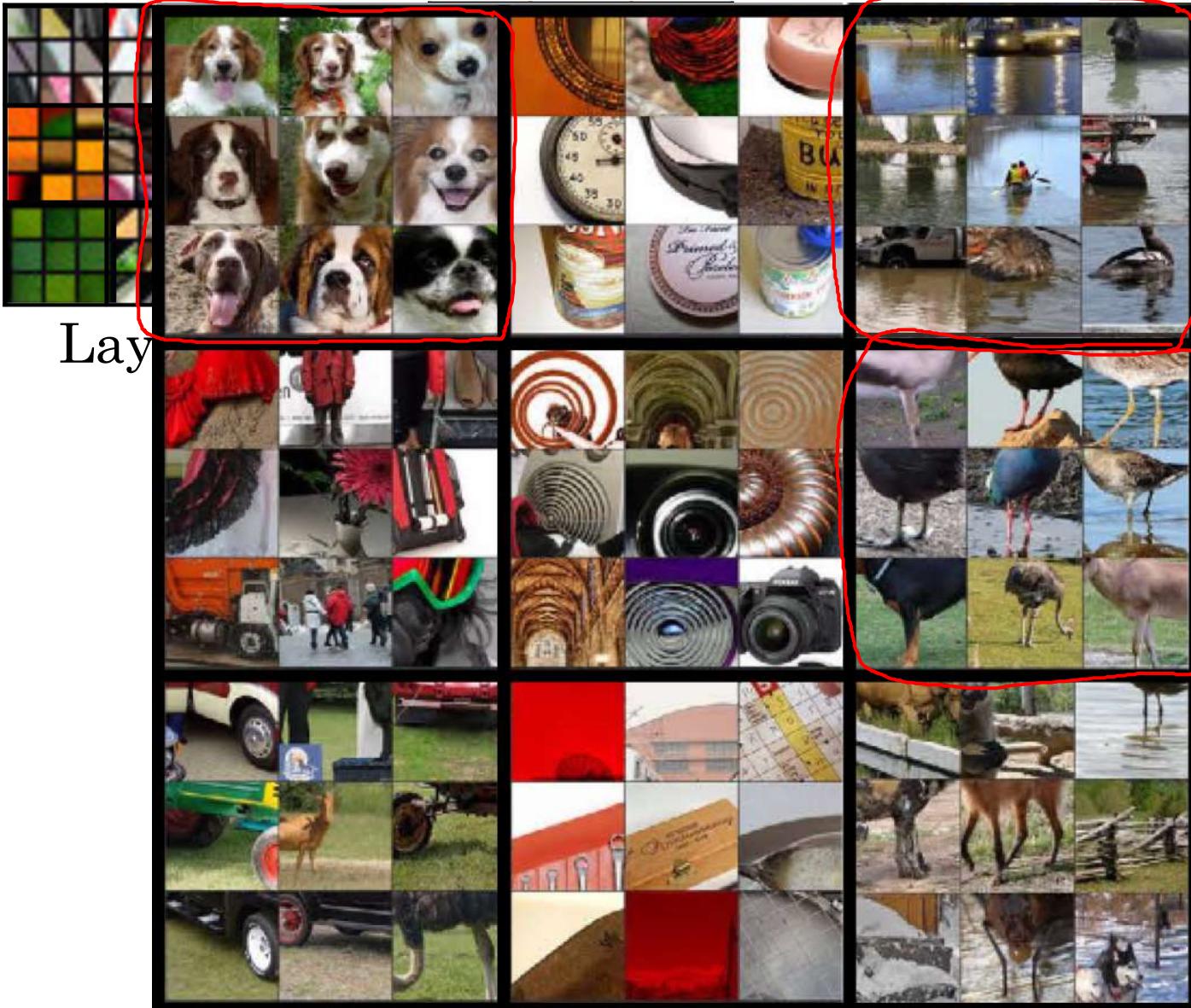


L



Layer 5

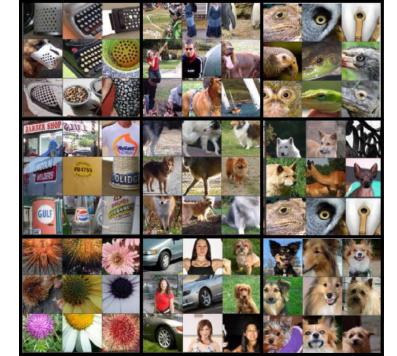
# Visualizing deep layers: Layer 4



Layer 4

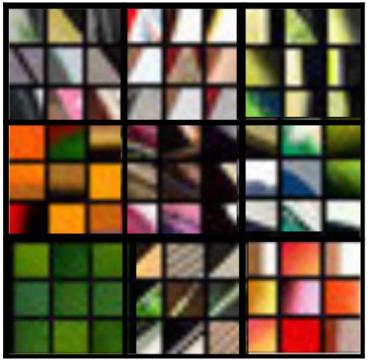


Layer 4



Layer 5

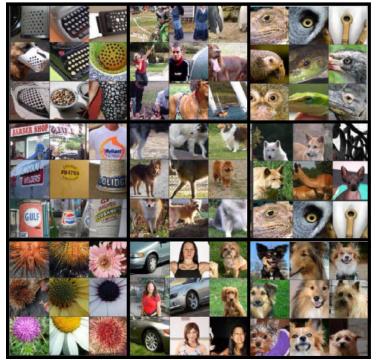
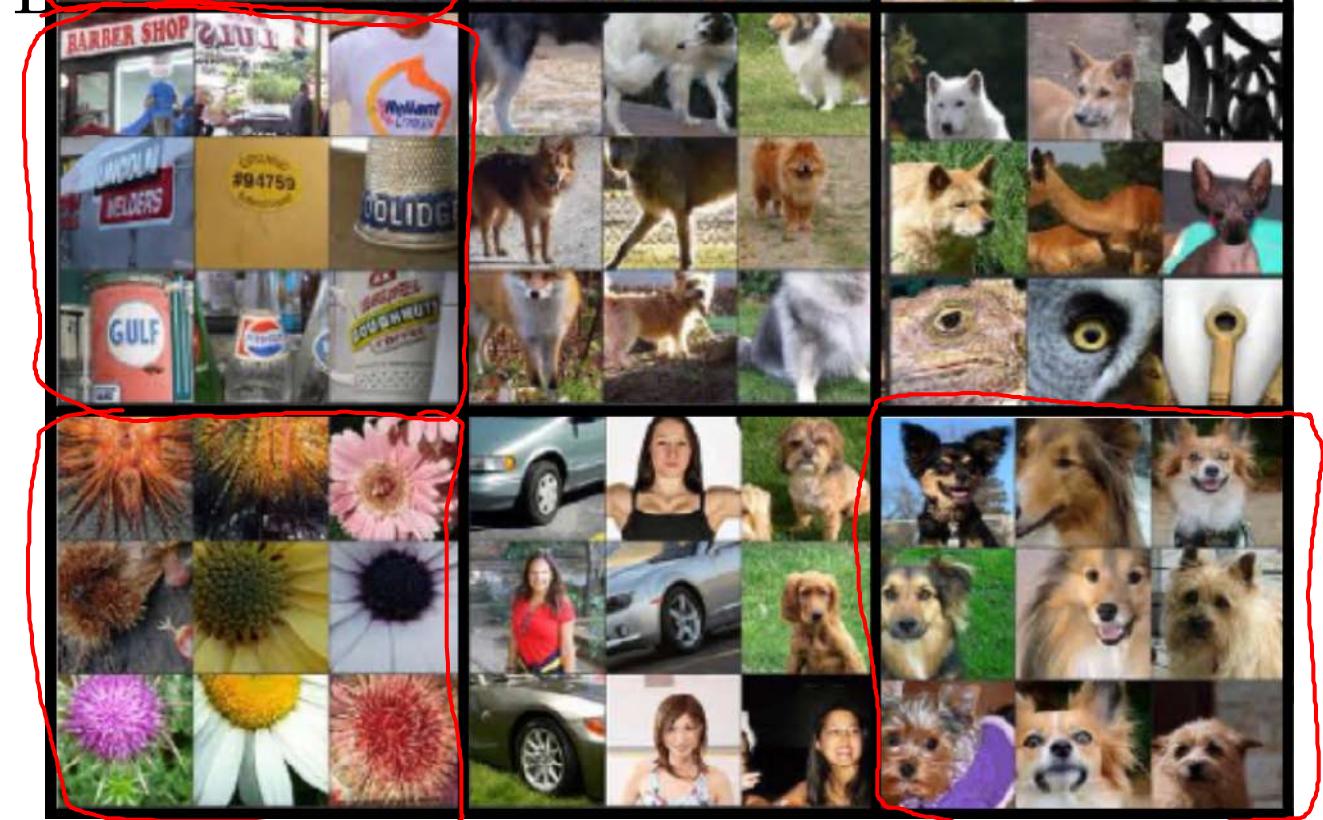
# Visualizing deep layers: Layer 5



Layer 1



L2



Layer 5

These are the steps for neural style transfer:

- Take content ( $C$ ), style ( $S$ ) images and generate a stylized image randomly ( $G$ ).
- Pass  $C, S, G$  into a pretrained CNN up to a certain intermediate layer  $L$  and compute the difference in contents between the representations of  $C$  and  $G$  and the difference in style between  $S$  and  $G$ . The difference in style

requires the preliminary computation of the style matrices of  $S$  and  $G$ .

- Based on the content and style differences, compute the loss associated with the current  $G$ .

- At this point



deeplearning.ai

# Neural Style Transfer

## Cost function

idea: you want similarity between:

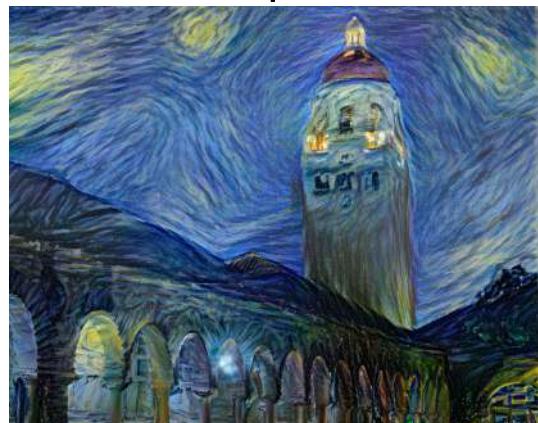
- generated image and content image
- style of the generated image and content image

# Neural style transfer cost function



Content C

Style S



Generated image G

cost function

$$J(G) = \alpha J_{\text{Content}}(C, G)$$

Similarity between C and G

$$+ \beta J_{\text{Style}}(S, G)$$

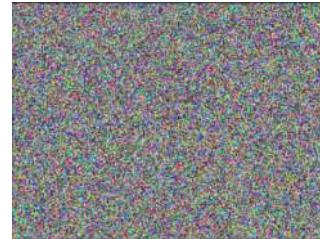
Similarity between the style of S and G

# Find the generated image G

1. Initiate G randomly

$G: \underbrace{100 \times 100}_{\text{---}} \times \underbrace{3}_{\text{---}}$

$\uparrow$   
RGB



2. Use gradient descent to minimize  $\underline{J(G)}$

$$G_t := G - \frac{\partial}{\partial G} J(G)$$





deeplearning.ai

# Neural Style Transfer

---

## Content cost function

# Content cost function

$$\underline{J(G)} = \alpha \underline{J_{content}(C, G)} + \beta J_{style}(S, G)$$

↑      ↓  
mid-level (not too shallow, not too deep)

- Say you use hidden layer  $\underline{l}$  to compute content cost.
- Use pre-trained ConvNet. (E.g., VGG network)
- Let  $\underline{a^{[l](C)}}$  and  $\underline{a^{[l](G)}}$  be the activation of layer  $\underline{l}$  on the images
- If  $a^{[l](C)}$  and  $a^{[l](G)}$  are similar, both images have similar content

$$J_{content}(C, G) = \frac{1}{2} \| \underbrace{a^{[l](C)}}_{\uparrow} - \underbrace{a^{[l](G)}}_{\uparrow} \|_2^2$$

IDEA: The difference in style between two images is measured as the difference of their respective "style matrices". These matrices are computed on the feature representation of an image at an hidden layer of a pretrained CNN (3D feature representation  $W \times h \times c$ ).



deeplearning.ai

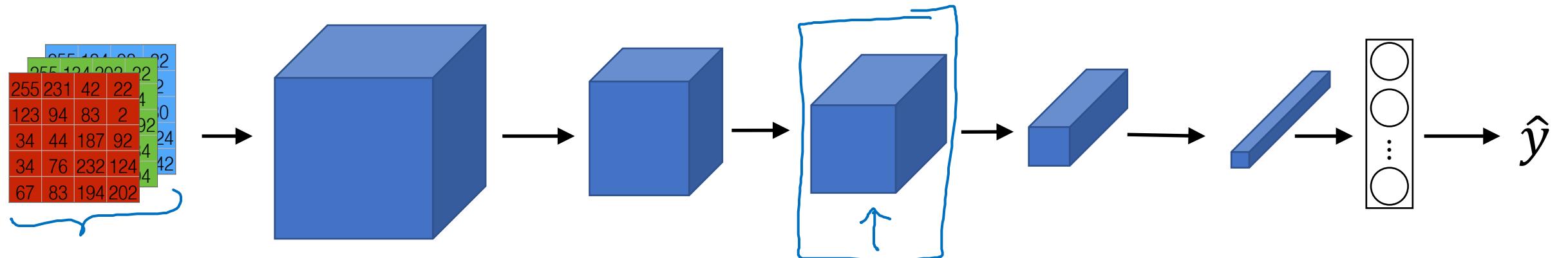
- take images ( $S$ tyle) and ( $G$ )enerated, pass them through a CNN and take their representation at layer  $l$ .
- for each of them compute the style matrix  $G^{(S)}$  and  $G^{(G)}$  ( $G$  stands for GRAM matrix actually). The style matrix is of  $n \times c \times channel \times channel$  and on the element  $kk'$  they have the correlation of channels  $k$  and  $k'$  (that is  $G_k$  linearized  $\cdot G_{k'}^T$  linearized).
- The loss is based on the Frobenius norm between the style and generated images.
- The loss is backpropagated through the network (with the  $G$  as input), keeping all the weights fixed except for those feeding back to the input pixels  $\Rightarrow$  change  $G$ .

# Neural Style Transfer

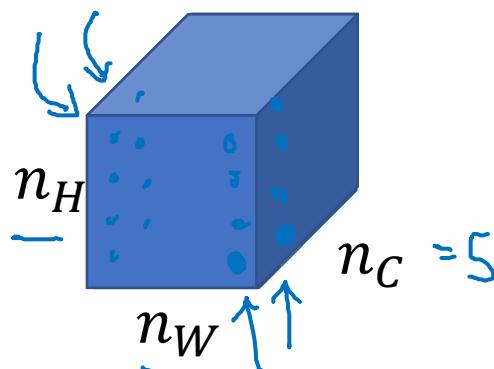
## Style cost function

The original image . In essence :

# Meaning of the “style” of an image

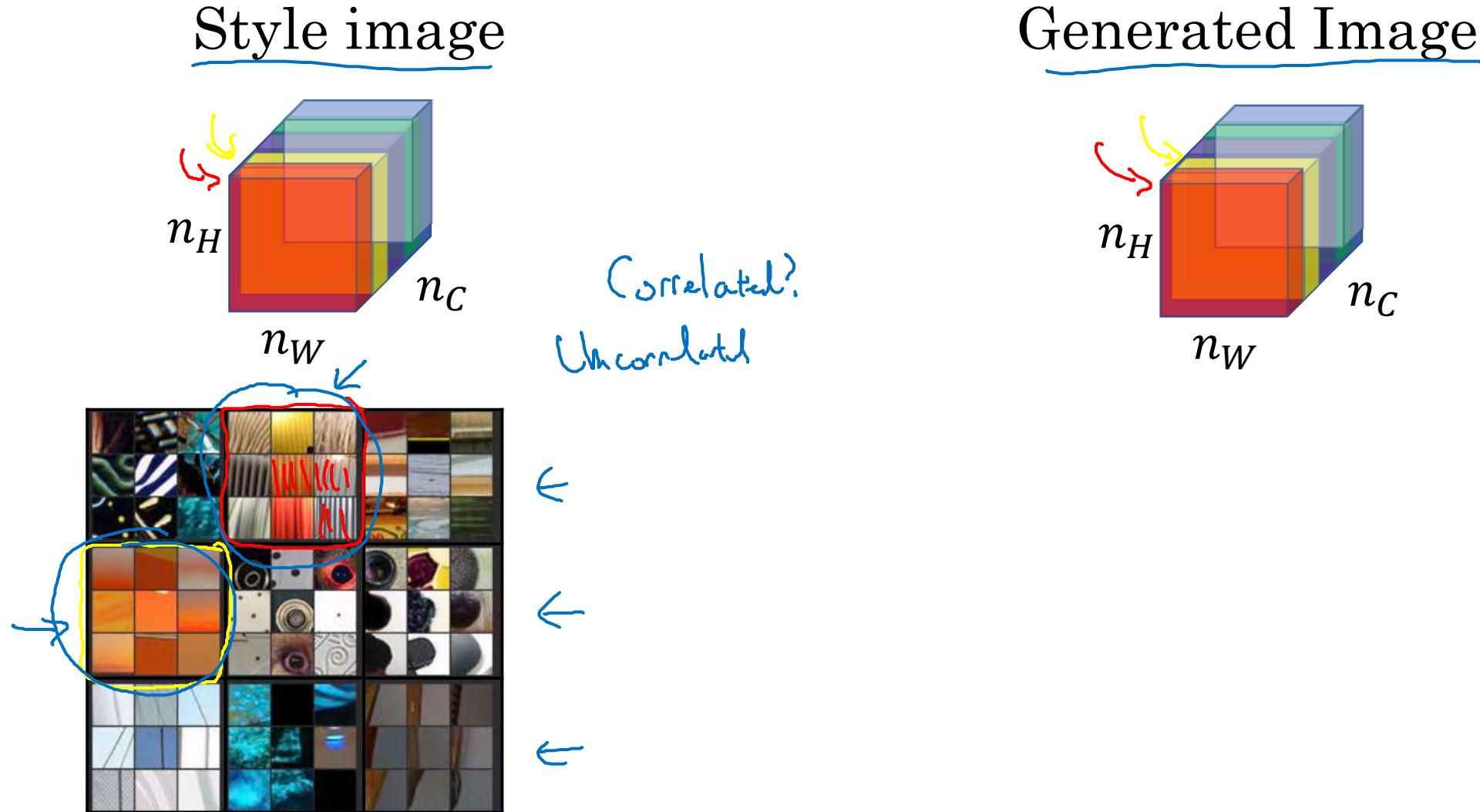


Say you are using layer  $l$ 's activation to measure “style.”  
Define style as correlation between activations across channels.



How correlated are the activations  
across different channels?

# Intuition about style of an image



# Style matrix

= GRAM MATRIX

H    W    C  
↓    ↓    ↗

Let  $a_{i,j,k}^{[l]}$  = activation at  $(i, j, k)$ .  $G^{[l]}$  is  $n_c^{[l]} \times n_c^{[l]}$

$$\rightarrow G_{kk'}^{[l](S)} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l](S)} a_{ijk'}^{[l](S)}$$

$$\rightarrow G_{kk'}^{[l](G_{\text{content}})} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l](G)} a_{ijk'}^{[l](G)}$$

"Gram matrix"

$n_c$

$$G_{kk'}^{[l]}$$

$$k = 1, \dots, n_c$$

$$\begin{aligned} J_{\text{style}}^{[l]}(S, G) &= \frac{1}{(\dots)} \| G^{[l](S)} - G^{[l](G)} \|_F^2 \\ &= \frac{1}{(2n_h^{[l]} n_w^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)})^2 \end{aligned}$$

# Style cost function

$$\left\| G^{[l](s)} - G^{[l](G)} \right\|_F^2$$

$$J_{style}^{[l]}(S, G) = \frac{1}{\left(2n_H^{[l]} n_W^{[l]} n_C^{[l]}\right)^2} \sum_k \sum_{k'} (G_{kk'}^{[l](s)} - G_{kk'}^{[l](G)})$$

$$J_{style}(S, G) = \sum_l \lambda^{[l]} J_{style}^{[l]}(S, G)$$

$$\underline{J(G)} = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$



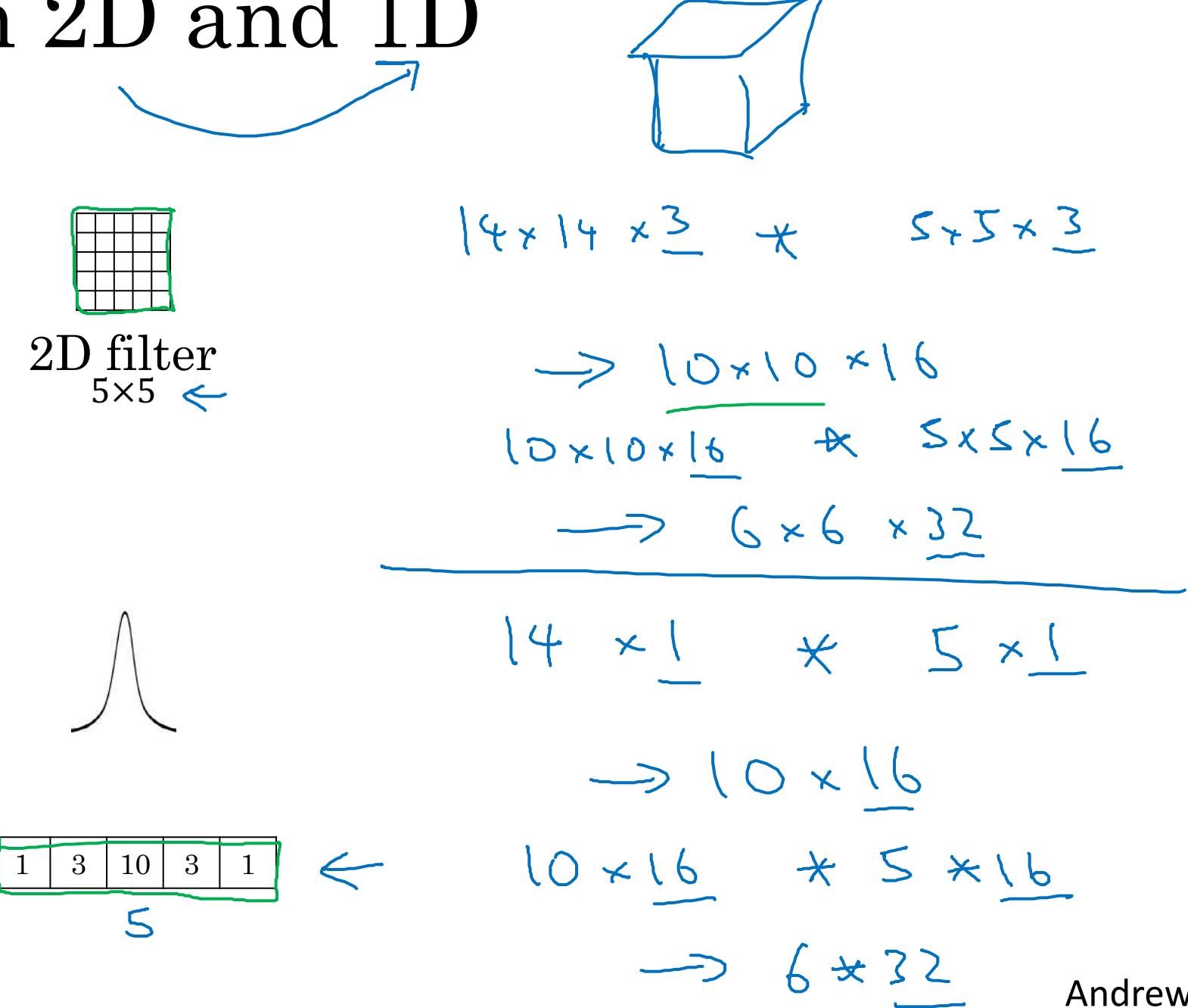
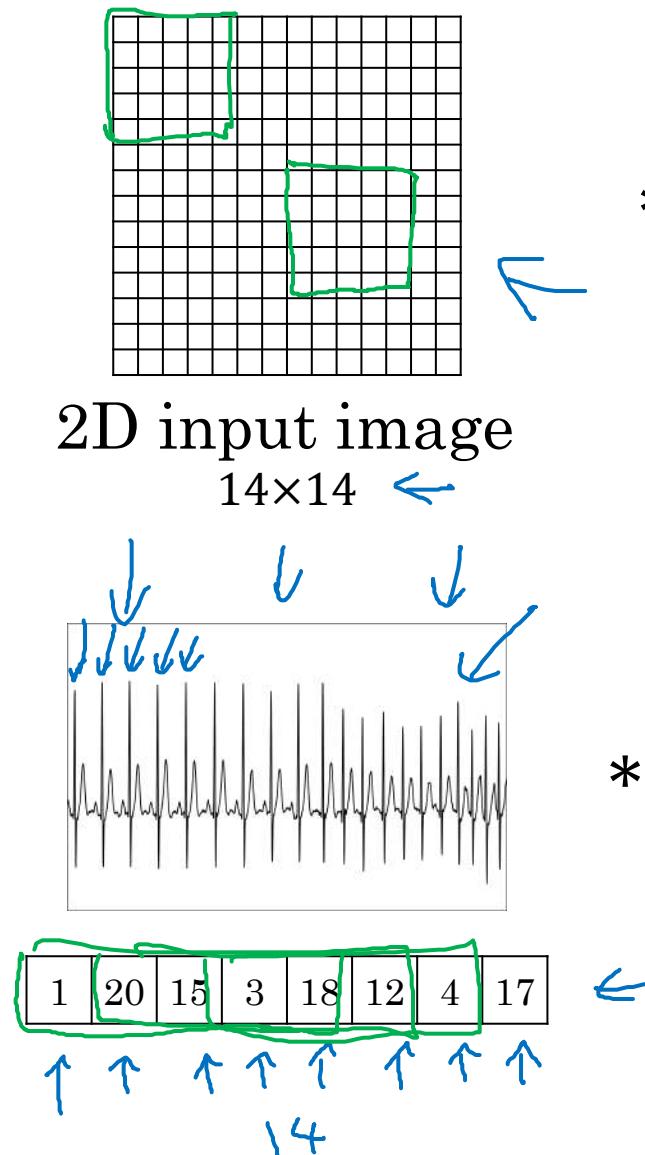
deeplearning.ai

# Convolutional Networks in 1D or 3D

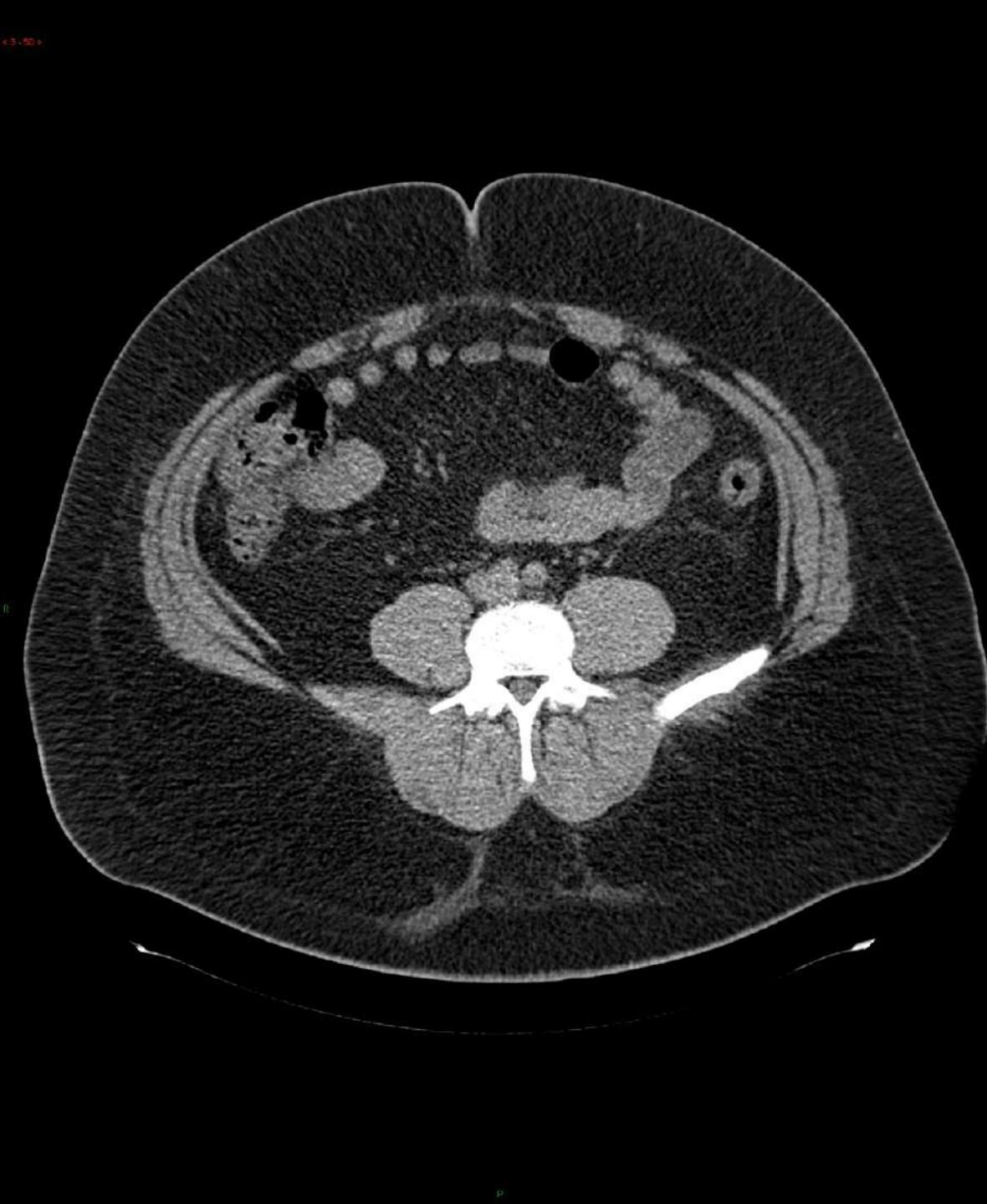
---

1D and 3D  
generalizations of  
models

# Convolutions in 2D and 1D



# 3D data



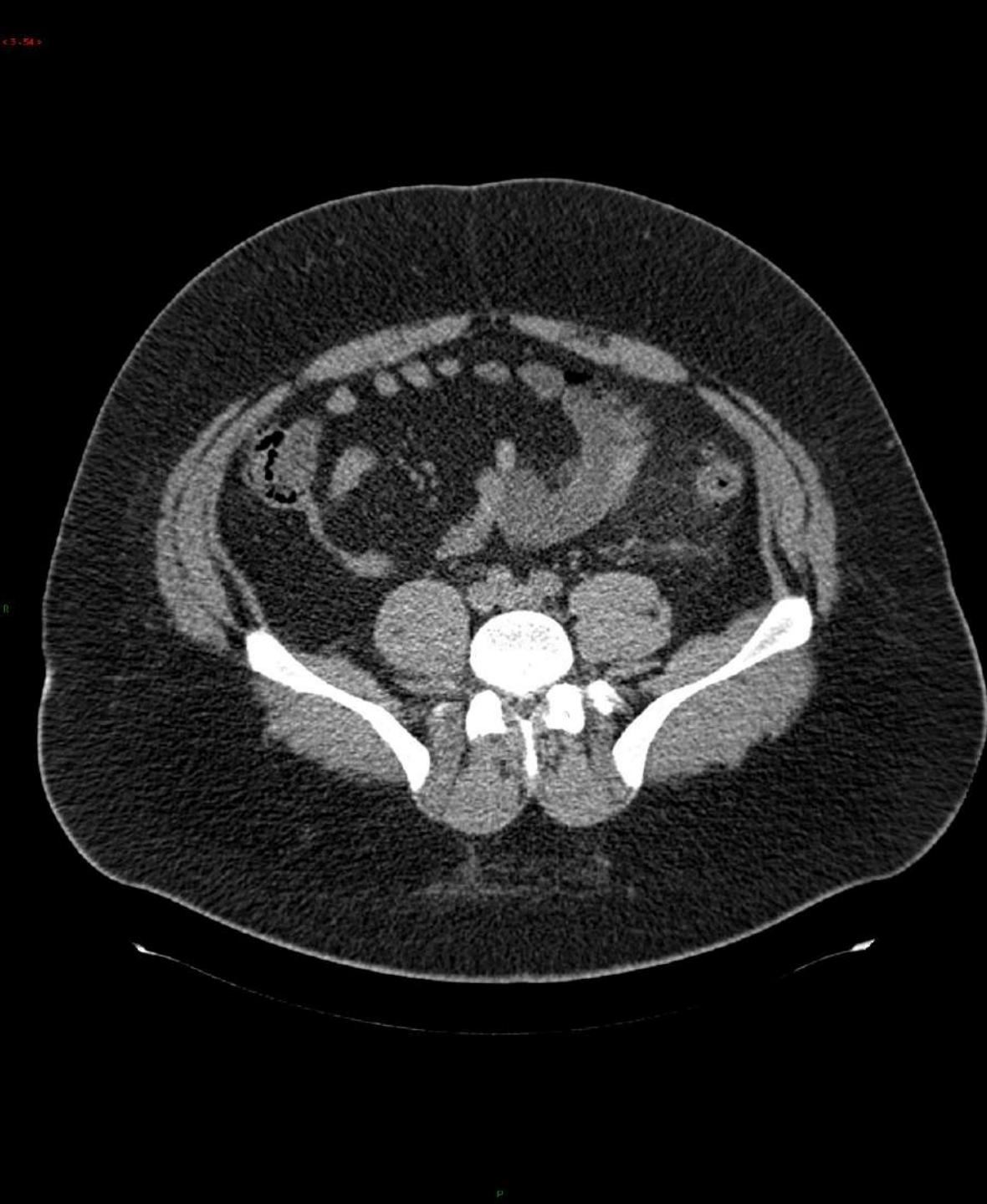
Andrew Ng

3D data



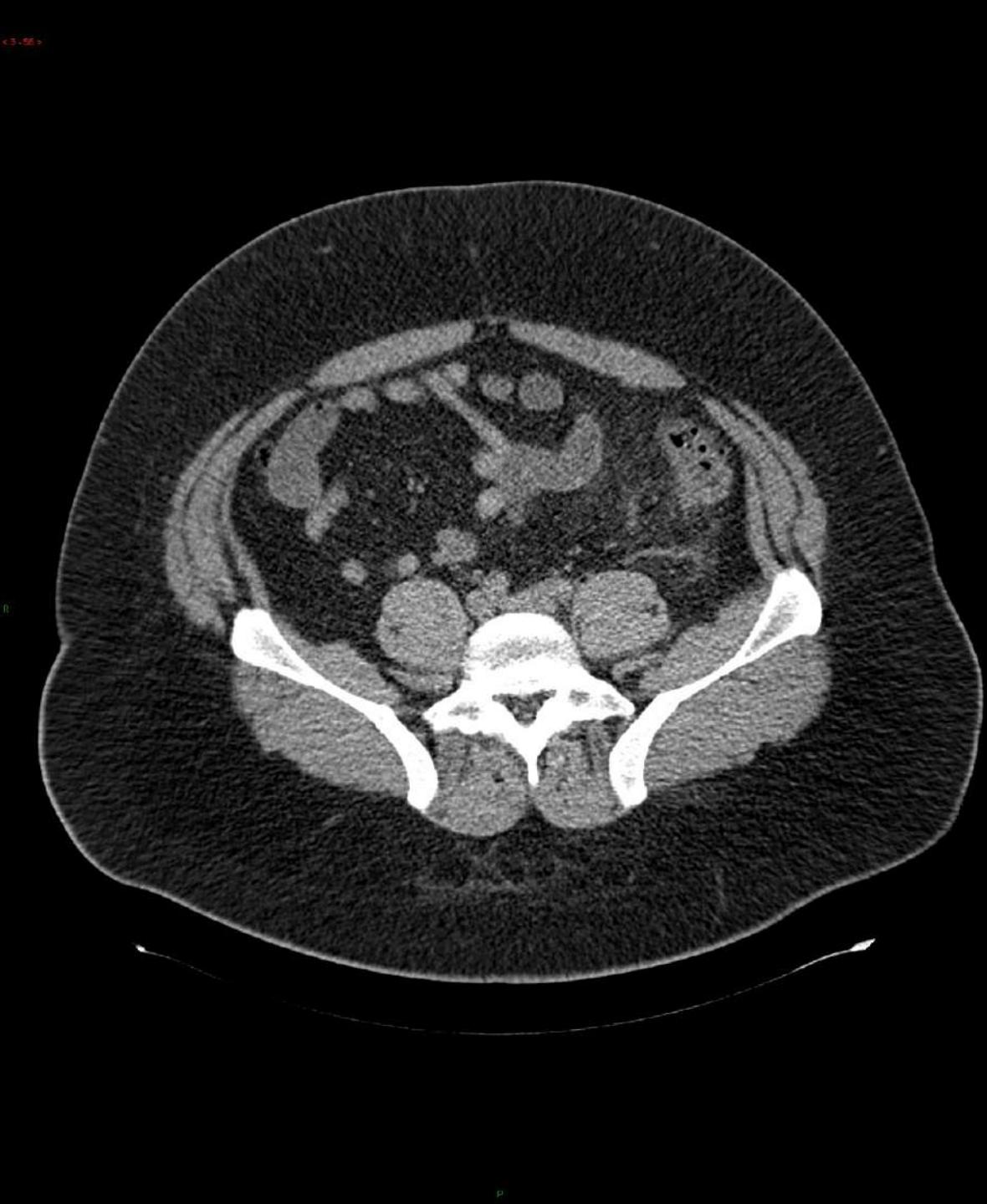
Andrew Ng

# 3D data

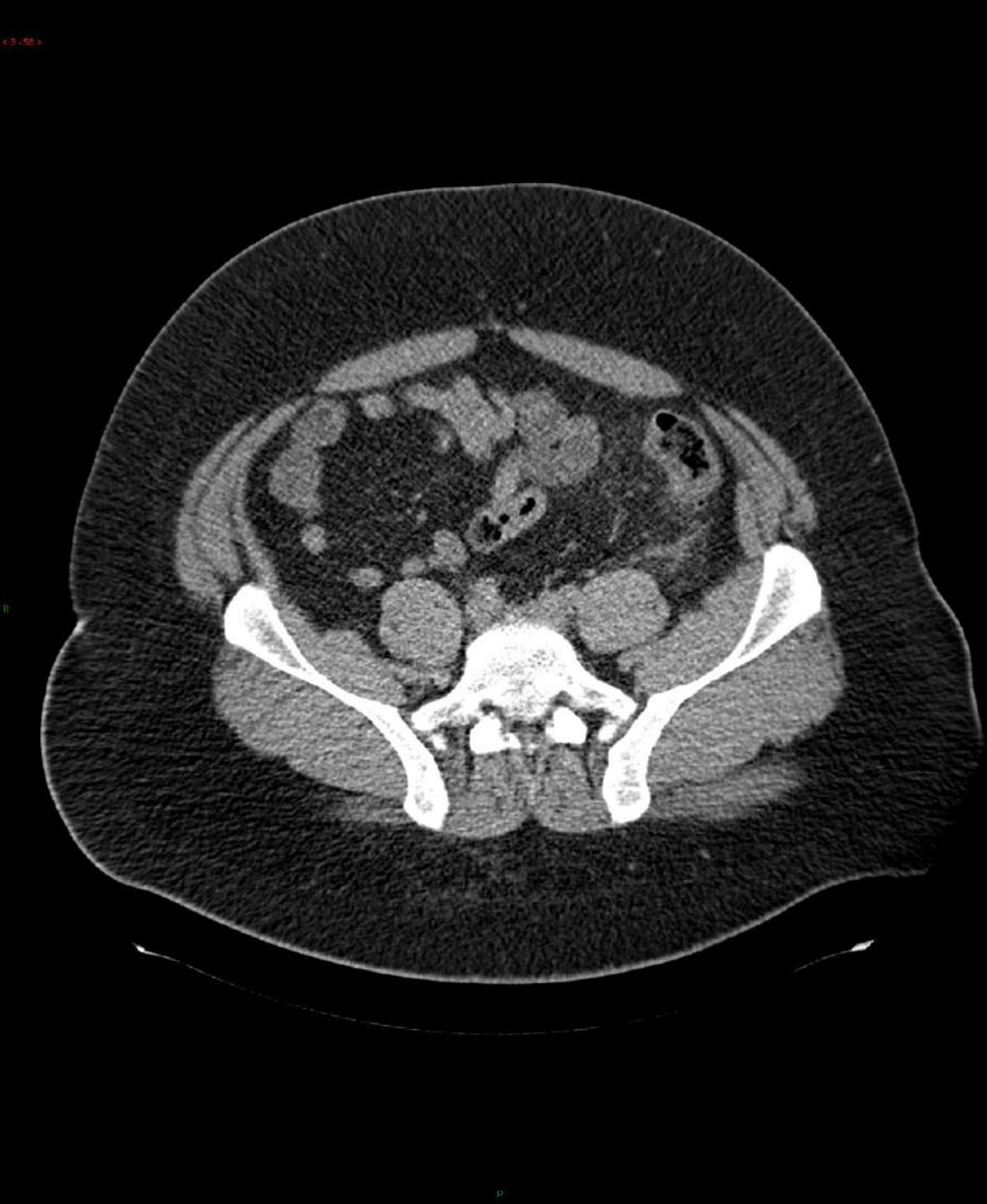


Andrew Ng

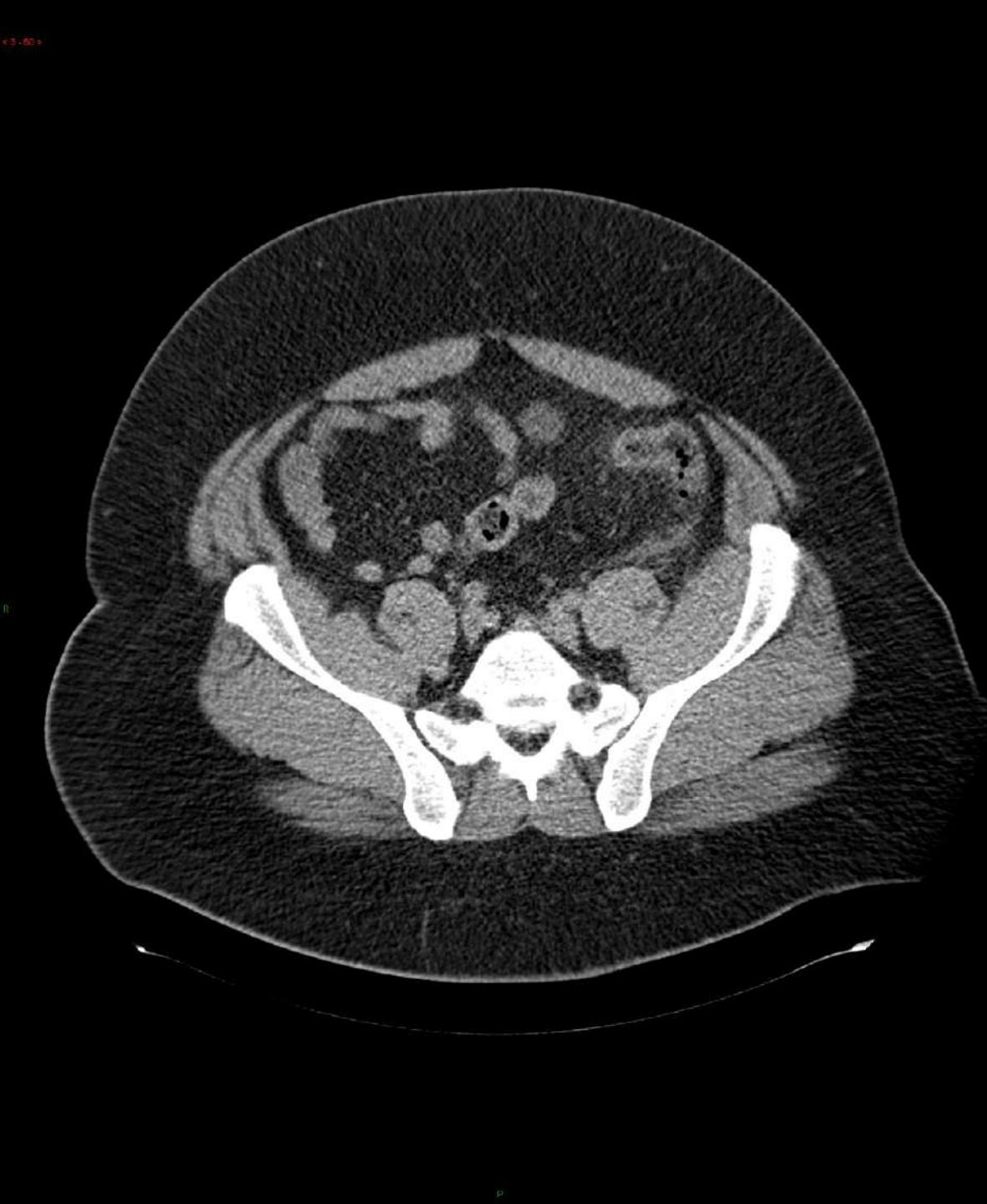
# 3D data



# 3D data



# 3D data



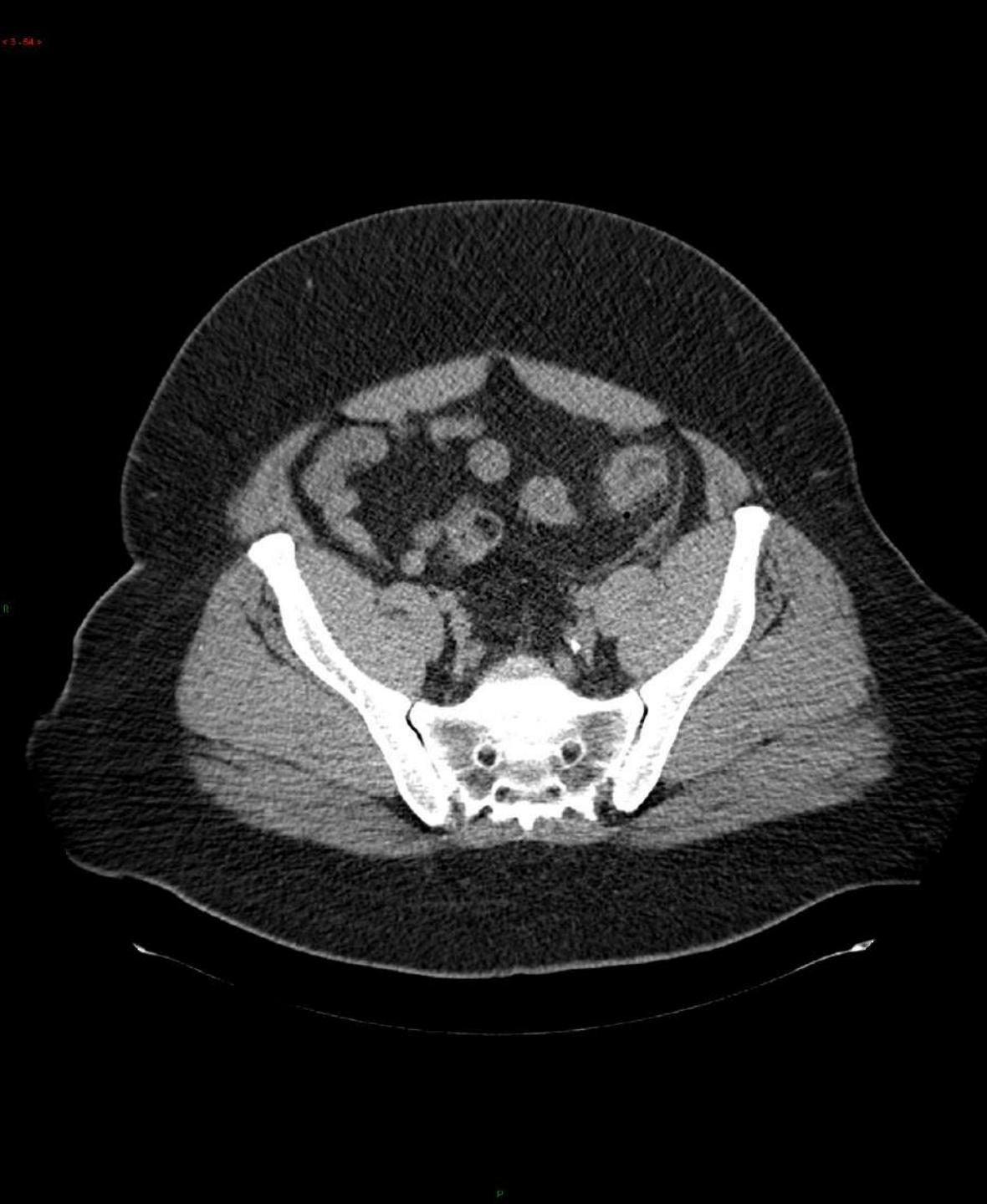
Andrew Ng

# 3D data



Andrew Ng

# 3D data



Andrew Ng

# 3D data

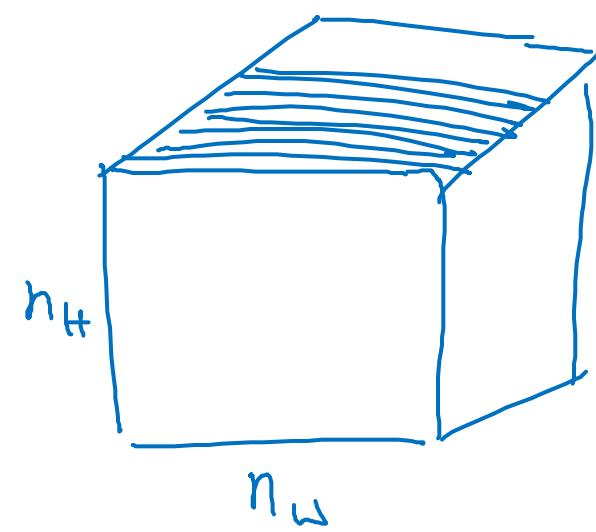


# 3D data



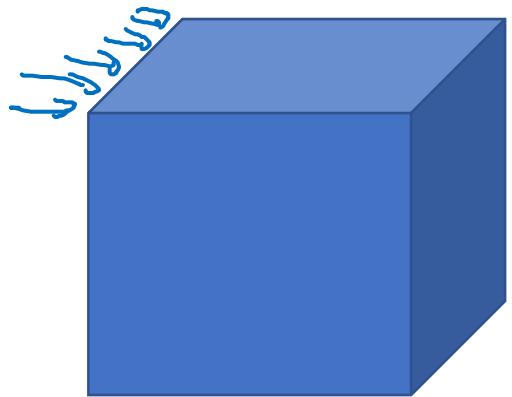
Andrew Ng

# 3D data

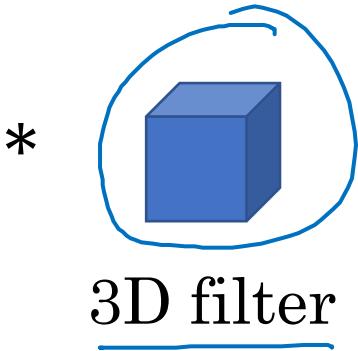


Andrew Ng

# 3D convolution



3D volume



$$\begin{array}{c} \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ \underbrace{4 \times 4 \times 4}_{\text{Input}} \times 1 \\ * \quad \underbrace{5 \times 5 \times 5}_{\text{Filter}} \times 1 \quad 16 \text{ filters.} \\ \rightarrow 10 \times 10 \times 10 \times 16 \\ * \quad \underbrace{5 \times 5 \times 5}_{\text{Stride}} \times 16 \\ \rightarrow 6 \times 6 \times 6 \times 32 \end{array}$$

Diagram illustrating the 3D convolution process:

- The input is a 3D volume of size  $4 \times 4 \times 4$ .
- The input is processed by a 3D filter of size  $5 \times 5 \times 5$ , resulting in 16 filters of size  $10 \times 10 \times 10$ .
- The output is then processed by another 3D filter of size  $5 \times 5 \times 5$  with stride 2, resulting in 32 filters of size  $6 \times 6 \times 6$ .