

Deep learning

Notes : Perceptron , Activation Function , Gradient Descent , backpropagation , loss Function , Metrics , Optimizer , Learning Rate , CNN , RNN , Structured Data Unstructured Data , Research Papers , History

Convolutional Neural Network

PCA color augmentation

Edge detection example

3	0	1	-1	2	7	4
1	5	8	9	3	1	
2	7	2	5	1	3	
0	1	3	1	7	8	
4	2	1	6	2	8	
2	4	5	2	3	9	

6x6

"convolution"

*

1	0	-1
1	0	-1
1	0	-1

*

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	16

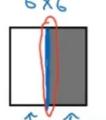
3x3
filter

4x4

$$3 + 1 + 1 \times 1 + \dots = -5 - - - - -$$

Vertical edge detection

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0



$$\begin{matrix} * & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} & = & \begin{matrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{matrix} \\ & 3 \times 3 & & 4 \times 4 \end{matrix}$$

1	1	1
0	0	0
-1	-1	-1

Horizontal

Vertical and Horizontal Edge Detection

1	0	-1
1	0	-1
1	0	-1

Vertical

1	1	1
0	0	0
-1	-1	-1

Horizontal

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

1	1	1
0	0	0
-1	-1	-1

$$\begin{matrix} * & \begin{matrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{matrix} & = & \begin{matrix} 0 & 0 & 0 & 0 \\ 30 & 10 & -10 & -30 \\ 30 & 10 & -10 & -30 \\ 0 & 0 & 0 & 0 \end{matrix} \end{matrix}$$

Andrew Ng

1	0	-1
2	0	-2
1	0	-1

3	0	-3
10	0	-10
3	0	-3

Sobel filter

Sobel filter

$$6 \times 6$$

$n \times n$

*

$$\begin{matrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{matrix}$$

$f \times f$

$= (4,4)$

Padding

$$6 \times 6$$

$n \times n$

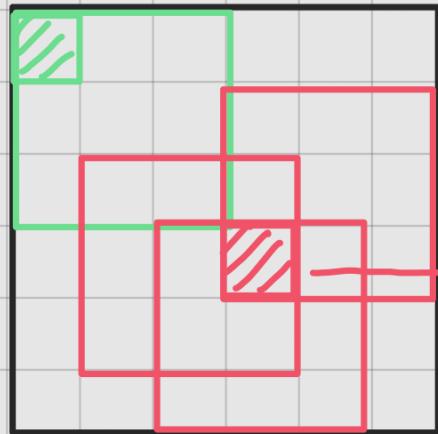
$$3 \times 3$$

$f \times f$

$$= 4 \times 4$$

$$n - f + 1 \times n - f + 1$$

$$6 - 3 + 1 = 4$$



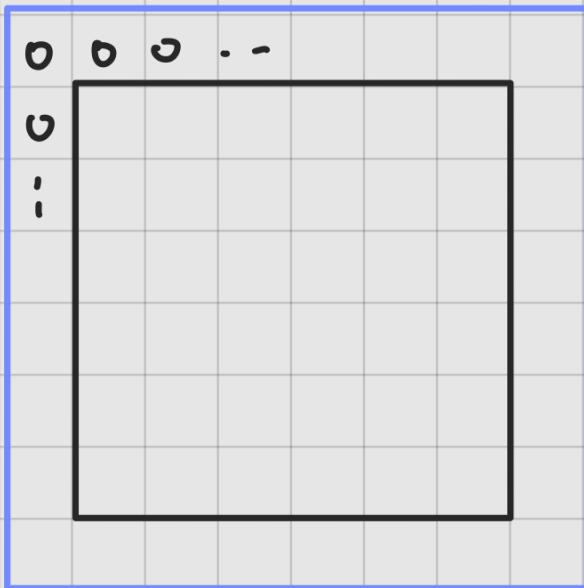
multiple outputs

Only one of the outputs

Problems: . Shrinky output

- throw away info from edges

Padding $p = 1$



$$6 \times 6 \rightarrow 8 \times 8$$

$$8 \times 8 * 3 \times 3 \rightarrow 6 \times 6 \quad (n + 2p - f + 1 = 6)$$

$p = \text{padding amount} = 1$

Valid and Same convolution

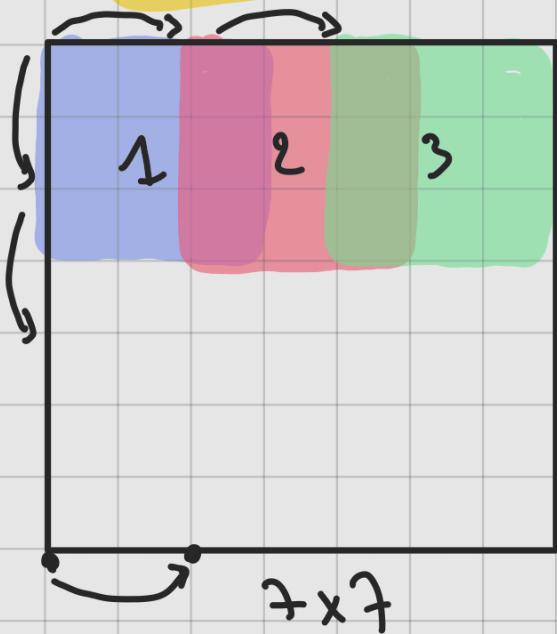
• "Valid": $n \times n * f \times f \rightarrow n - F + 1 \times n - F + 1$
 \hookrightarrow (no padding)

• "Same": Pad so output size same as input

$$\Rightarrow p = \frac{f - 1}{2}$$

f is usually odd

Strided Convolution



$$\star \quad 3 \times 3 \rightarrow 3 \times 3$$

Stride = 2

$$n \times n \quad * \quad f \times f \quad \longrightarrow \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \dots$$

padding p stride s ⌊ floor

Technical note on cross-correlation vs. convolution

Convolution in math textbook:

2	3	7	4	6	2
6	6	9	8	7	4
3	4	8	3	8	9
7	8	3	6	6	3
4	2	1	8	3	4
3	2	4	1	9	8

$$\begin{array}{c}
 * \\
 \begin{array}{|c|c|c|} \hline
 3 & 4 & 5 \\ \hline
 1 & 0 & 2 \\ \hline
 -1 & 9 & 7 \\ \hline
 \end{array}
 \end{array}
 \stackrel{=}{\longrightarrow}
 \begin{array}{|c|c|c|c|c|c|} \hline
 & & & & & \\ \hline
 \end{array}$$

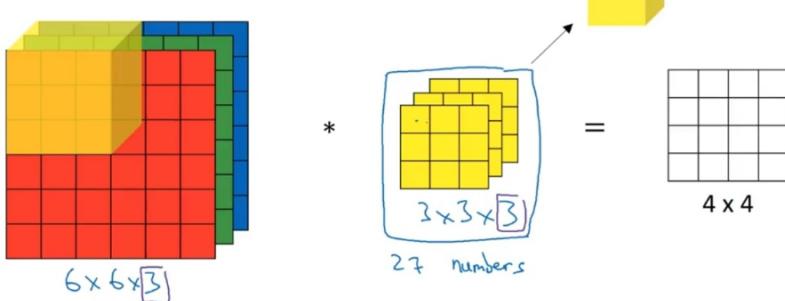
$(A * B) * C = A * (B * C)$

math lit

Andrew Ng

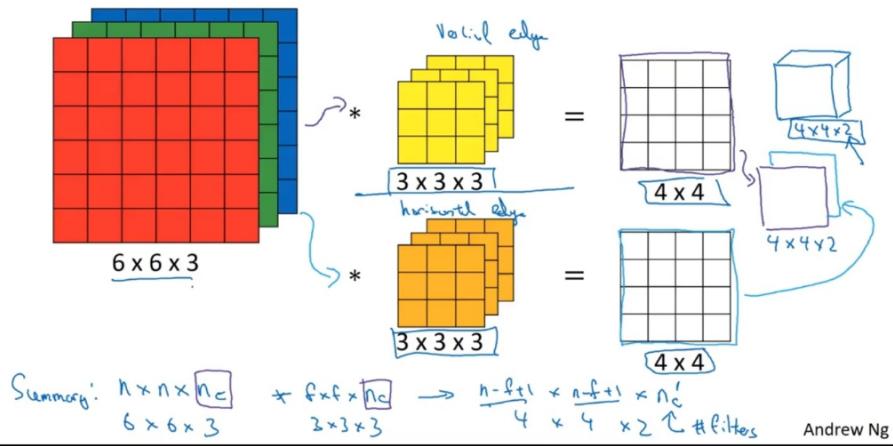
Convolution over volume

Convolutions on RGB image



Andrew Ng

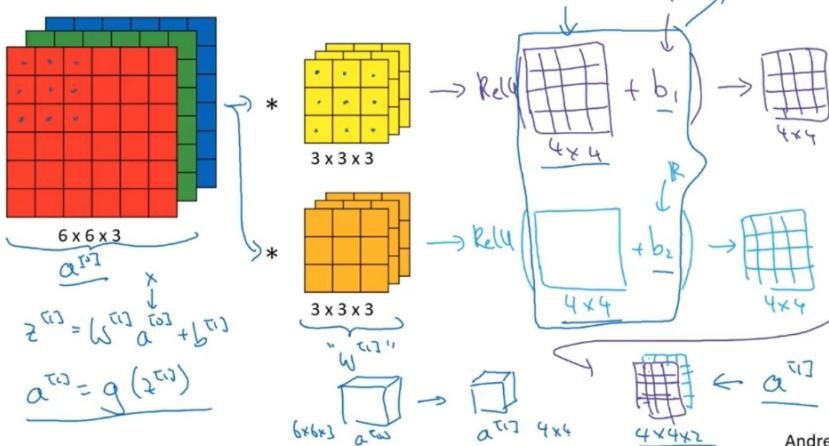
Multiple filters



Andrew Ng

One layer of a Convolutional Network

Example of a layer



Andrew Ng

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

$n_c^{[l]}$ = number of filters

Input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$$n^{[l]} = \left\lfloor \frac{n^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

Each filter is: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

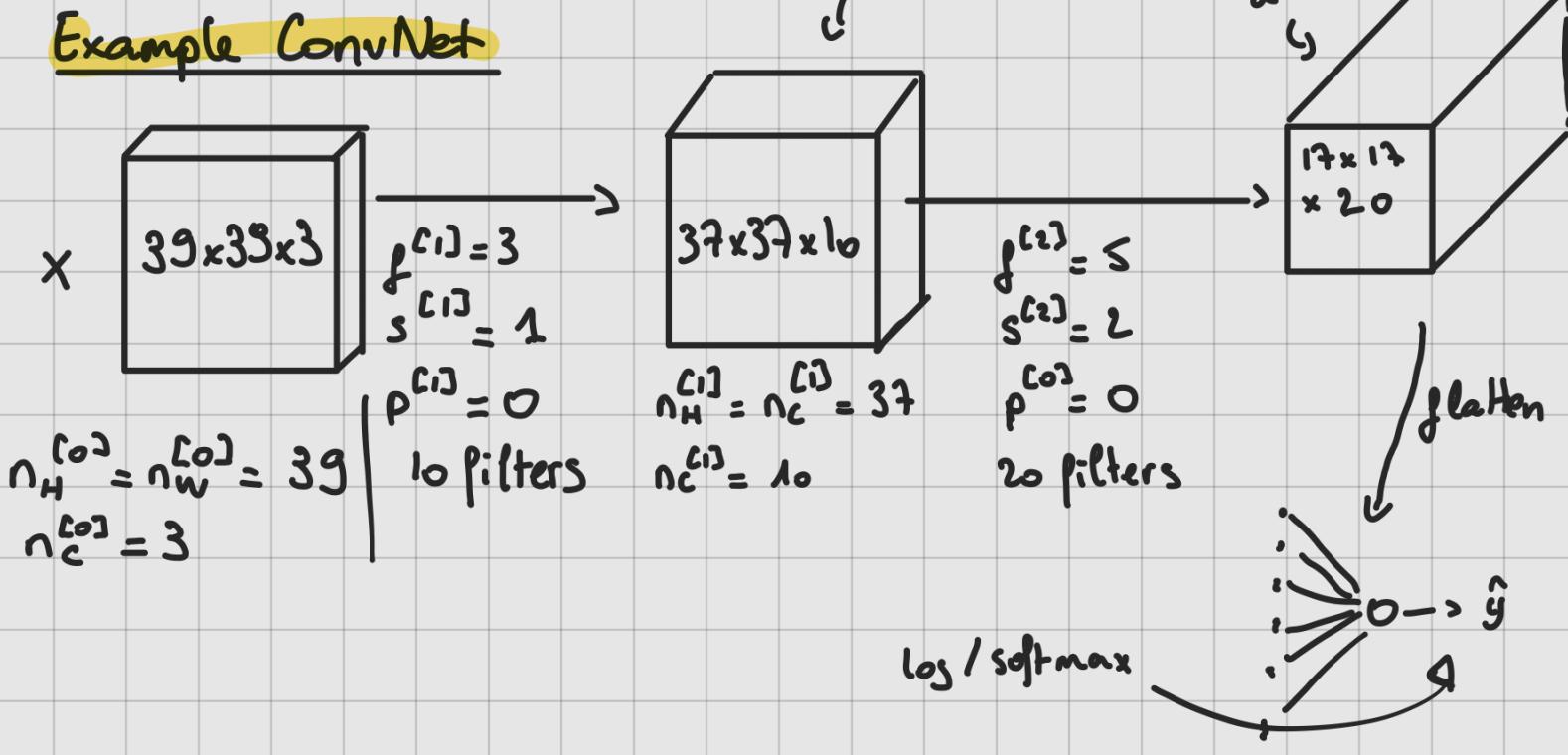
Activations: $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

bias: $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$

Example ConvNet

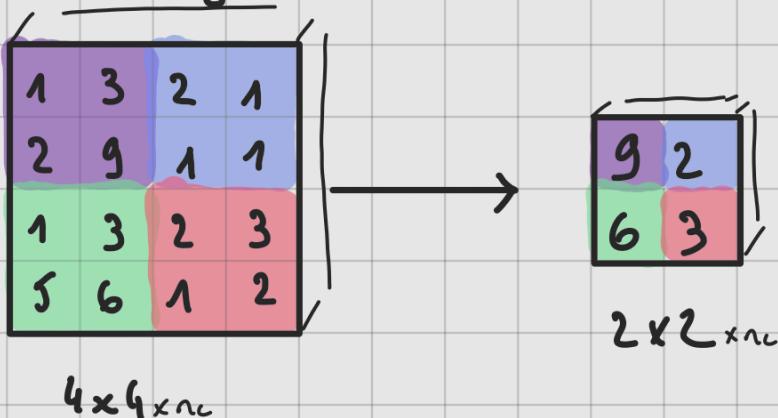


Layers in ConvNet :

- Convolution (Conv)
- Pooling (POOL)
- Fully Connected (FC)

Pooling layers

Max pooling



Hyperparams :

$$f = 2$$

$$s = 2$$

No learning params

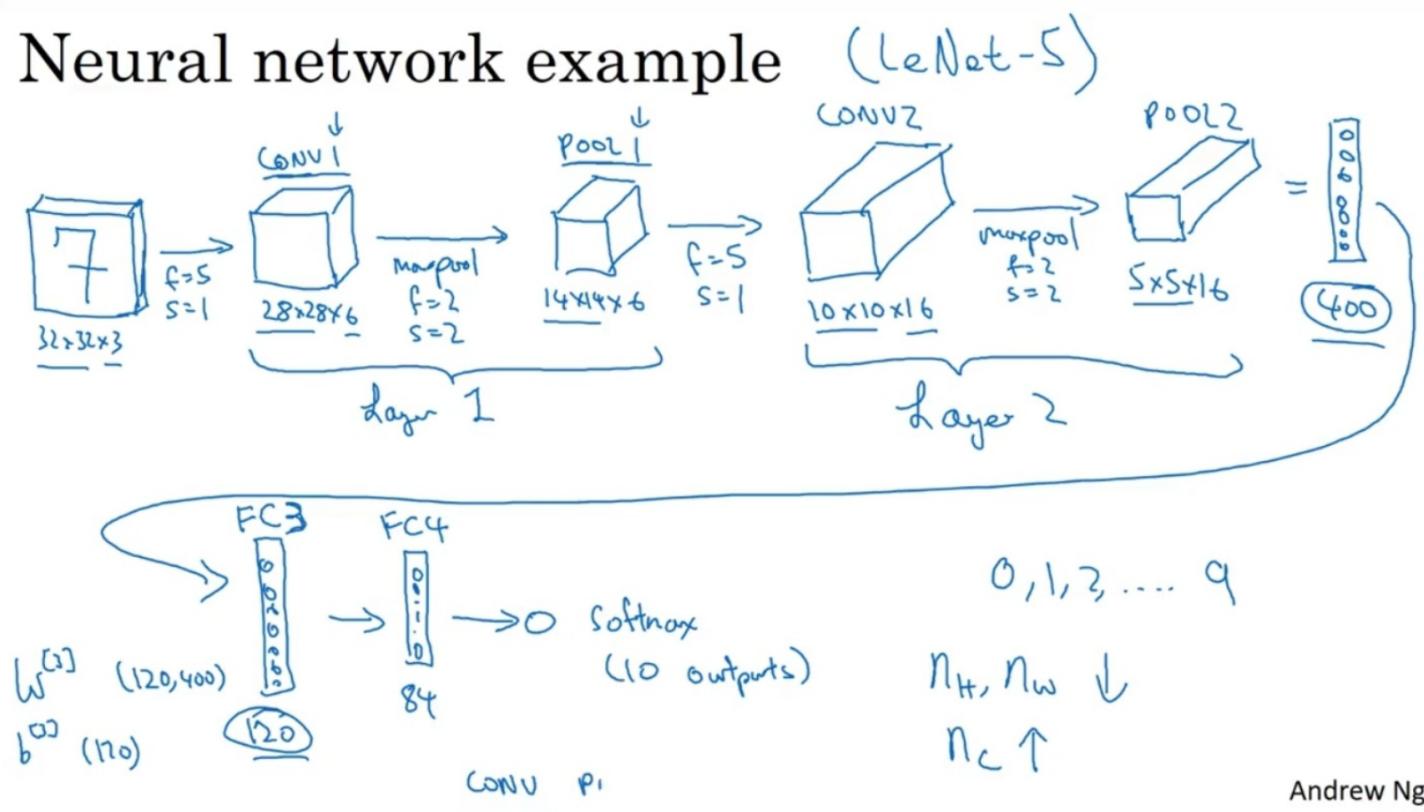
Average pooling : take average



usually not use padding

Example

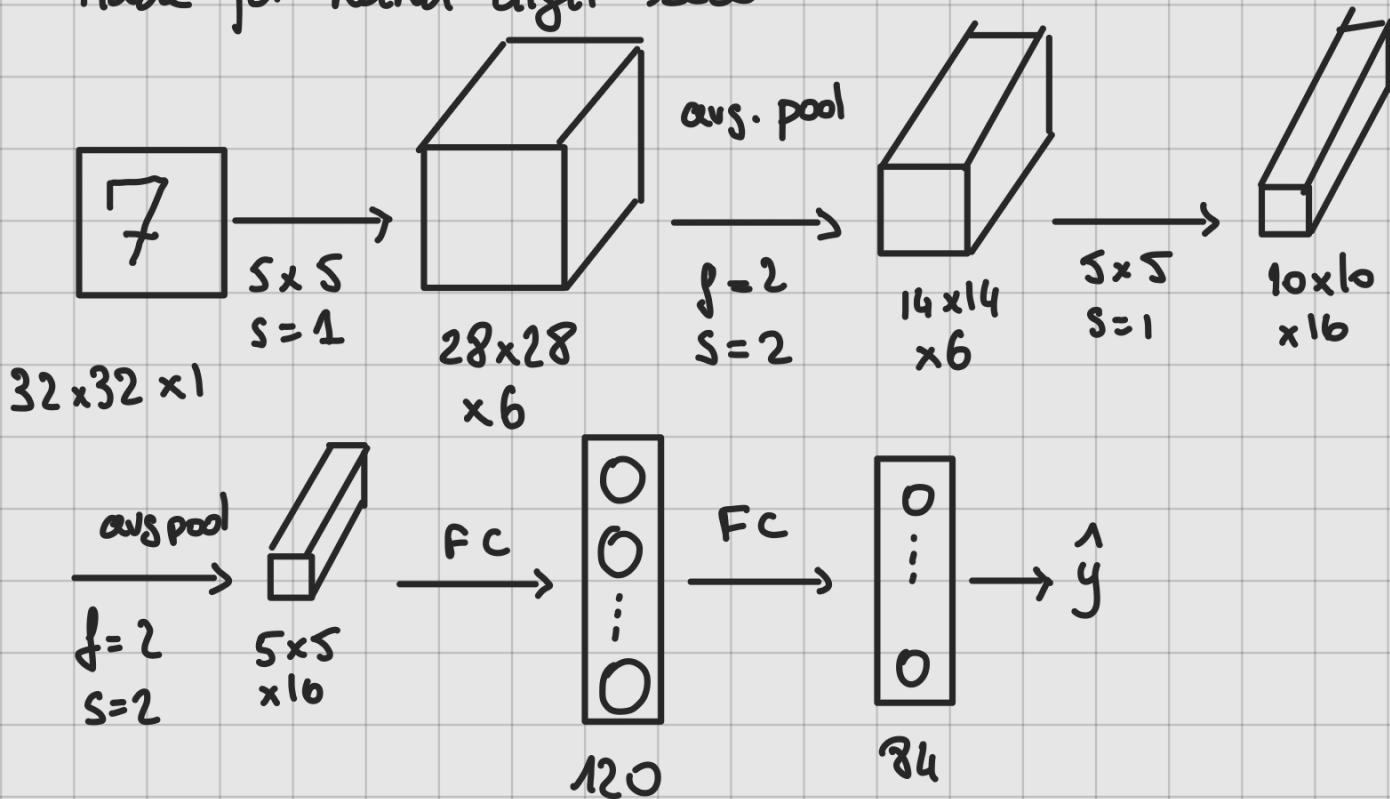
Neural network example



Classic network

LeNet-5 (1998 LeCun)

Made for hand digit reco.



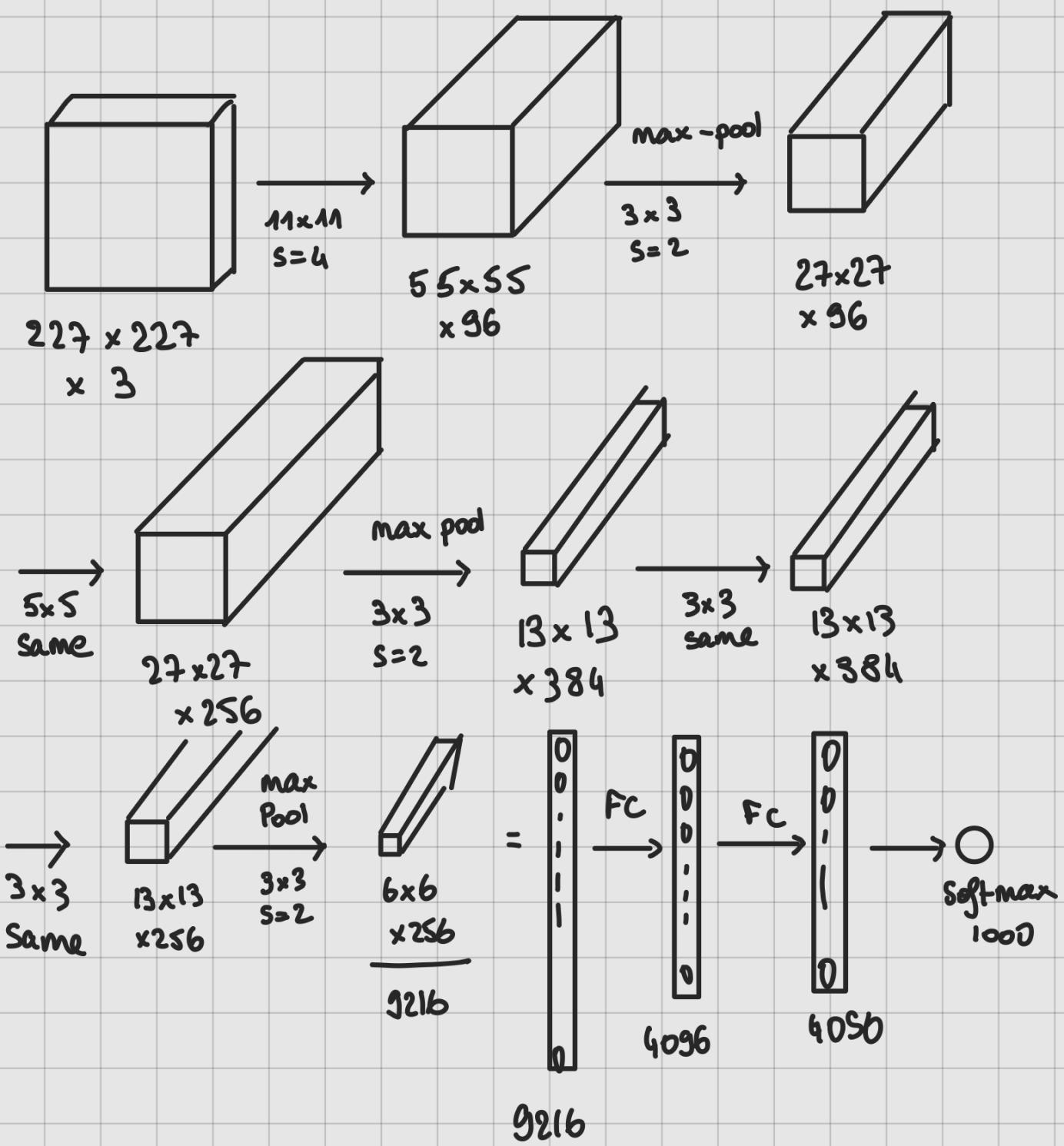
~ 60K parameters

$n_H \downarrow, n_W \downarrow, n_C \uparrow$

[Conv Pool] - [Conv Pool] - FC - FC - Output

Use sigmoid/tanh

AlexNet (2012 Krizhevsky)



$\sim 60M$ parameters

(multiple GPUs
Local Response Normalized (LRN))

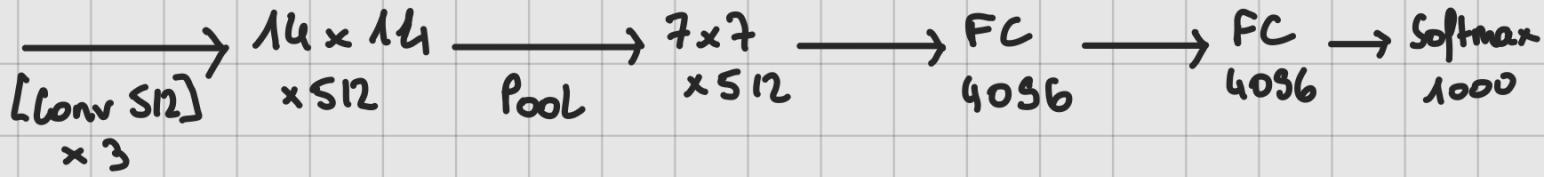
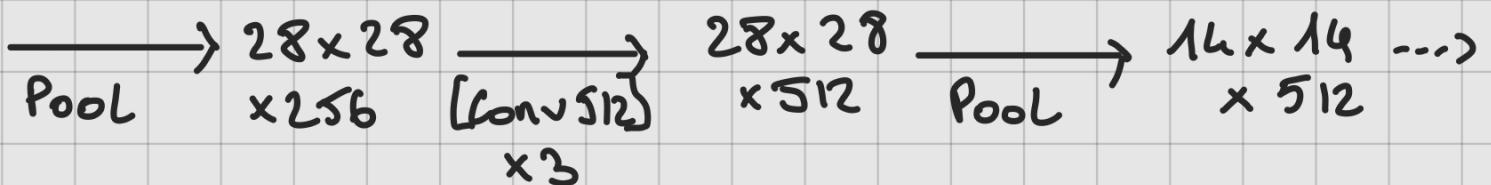
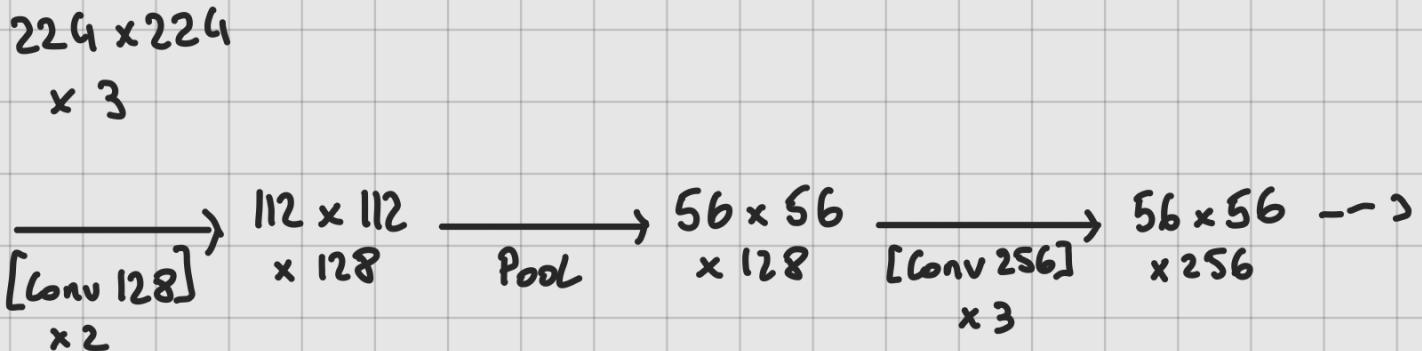
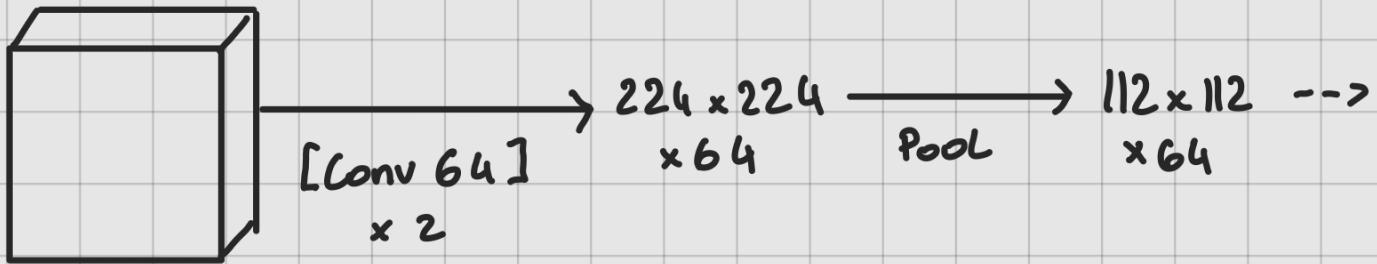
Use ReLU

↑
not use

VGG-16 (2015 Simonyan, Zisserman)

Conv = 3×3 filter, $s=1$, same

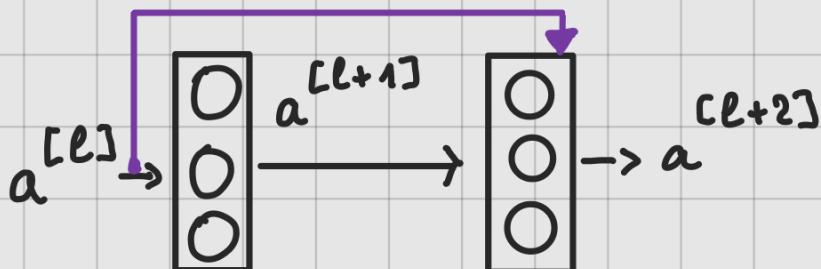
Max Pool = 2×2 , $s=2$



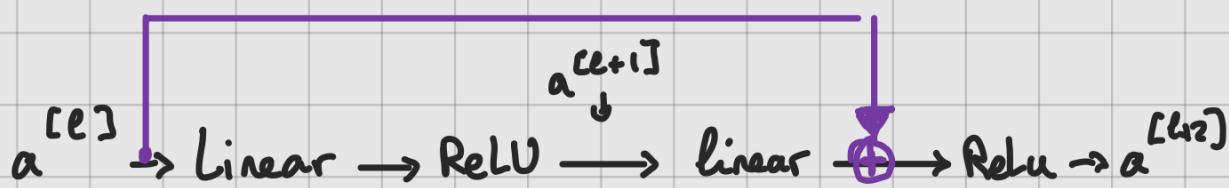
~ 138 M params

ResNets

Residual Block



"short cut" / Skip connection



$$z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]} \quad a^{[l+1]} = g(z^{[l+1]})$$

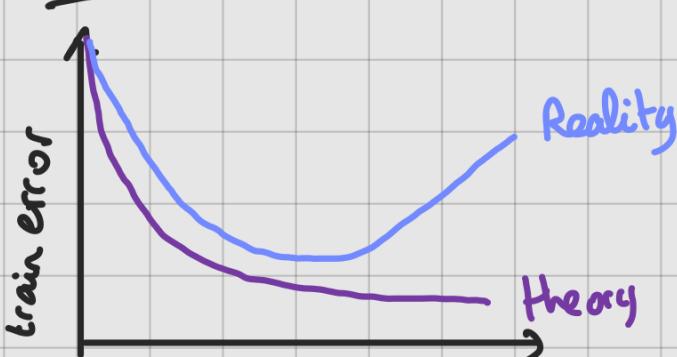
$$z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]} \quad a^{[l+2]} = g(z^{[l+2]})$$

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

"Plain Network" \times $\square \rightarrow \square \rightarrow \square$

"Residual Net" \times $\square \rightarrow \square \xrightarrow{\text{skip}} \square \rightarrow \square$

Plain



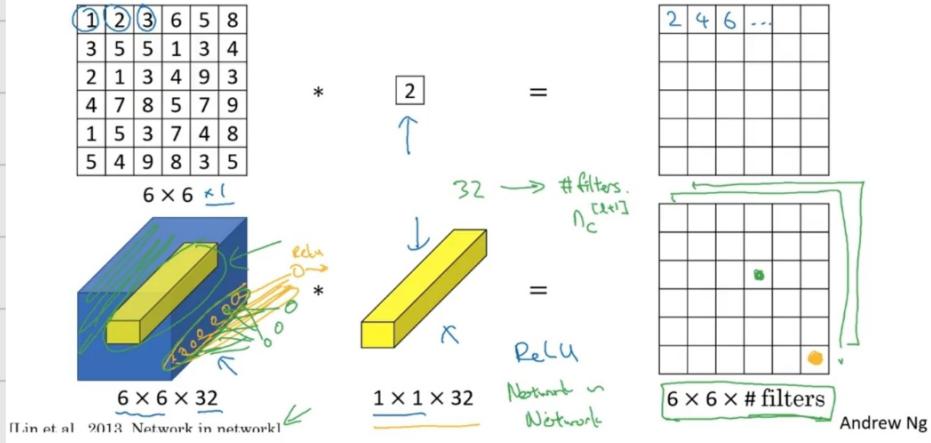
Resnet



Networks in Networks and 1×1 convolution

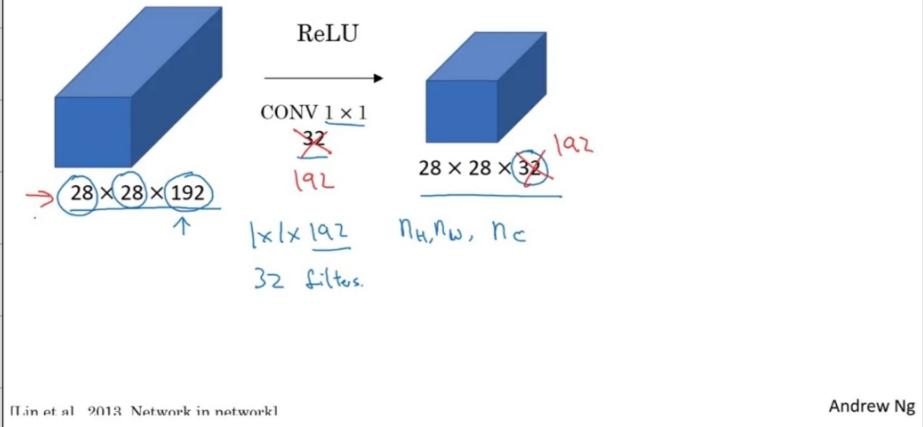
(2013, Lin)

Why does a 1×1 convolution do?



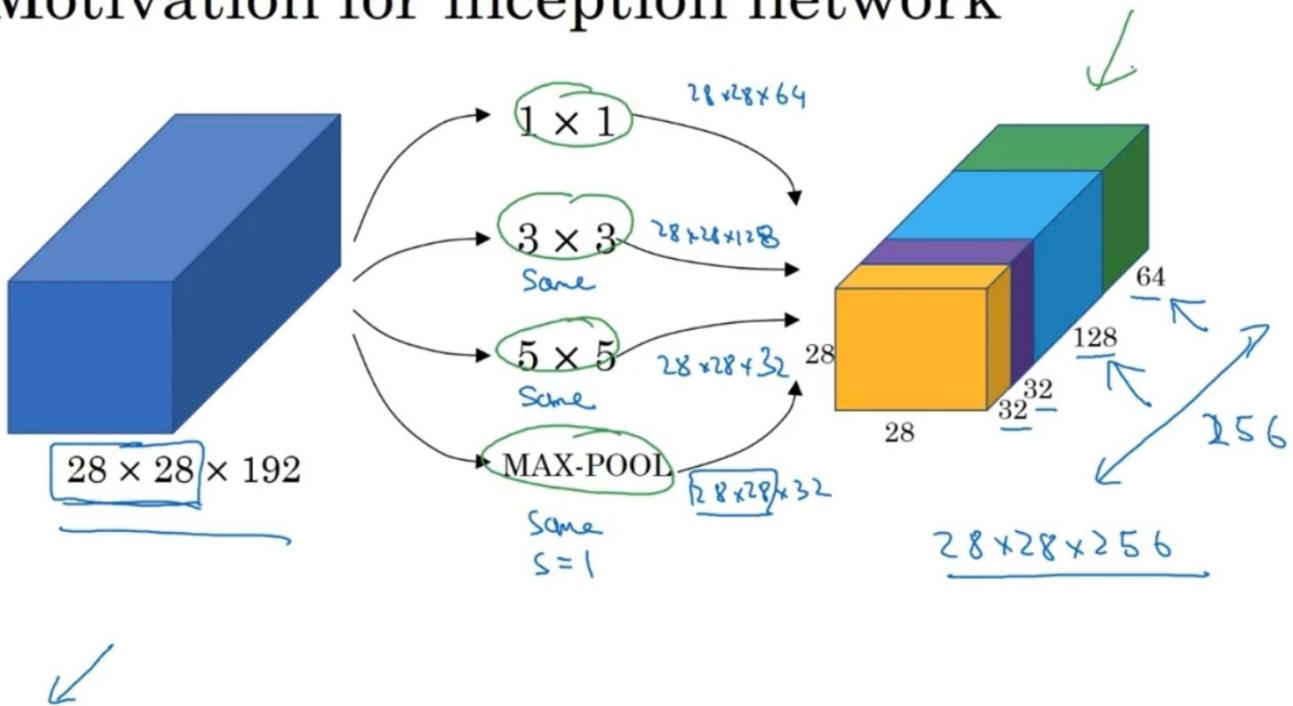
Shrink nb of filters

Using 1×1 convolutions



Inception Network (2014 Szegedy)

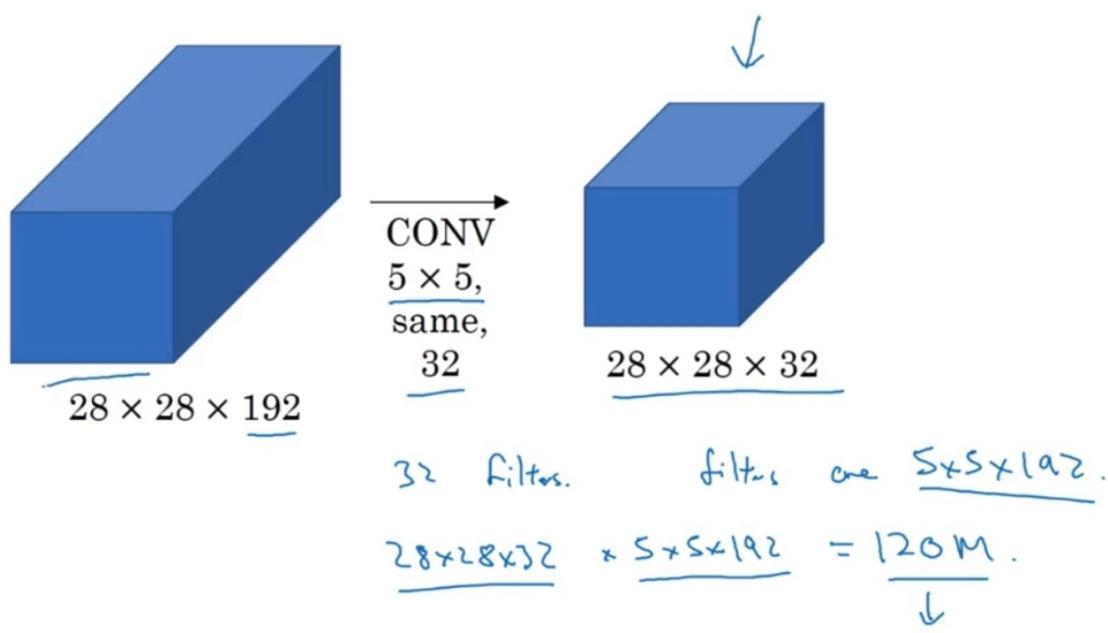
Motivation for inception network



[Szegedy et al. 2014. Going deeper with convolutions]

Andrew Ng

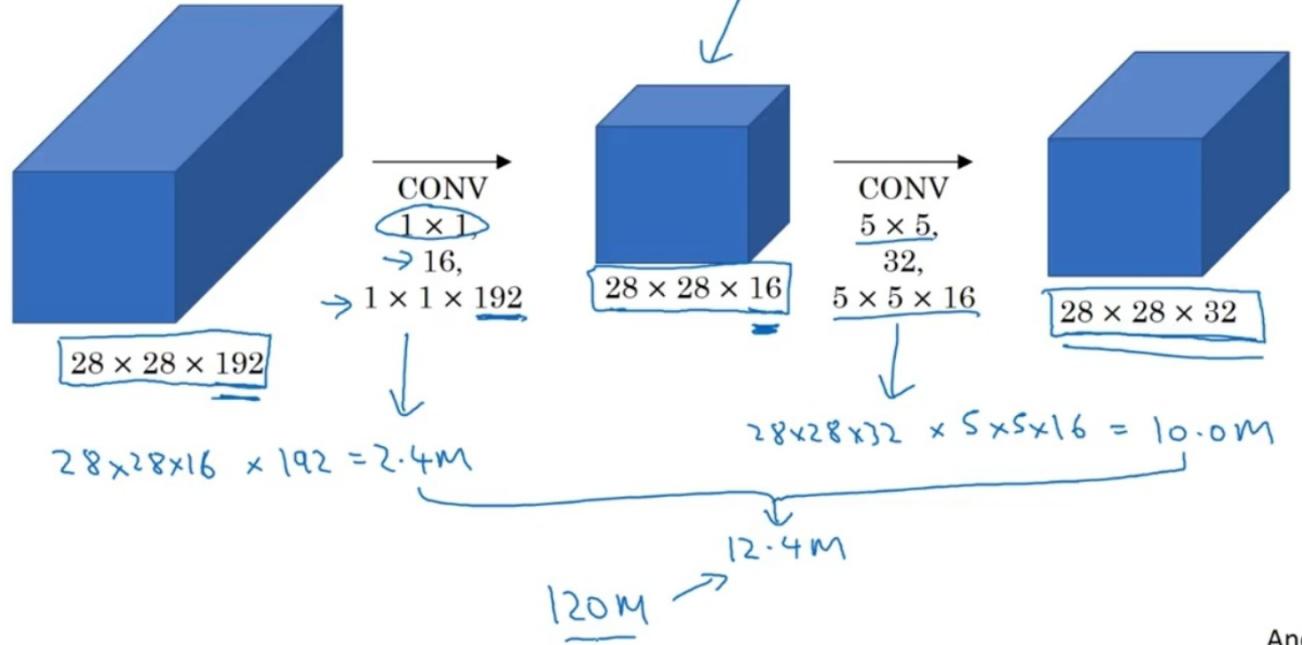
The problem of computational cost



Andrew Ng

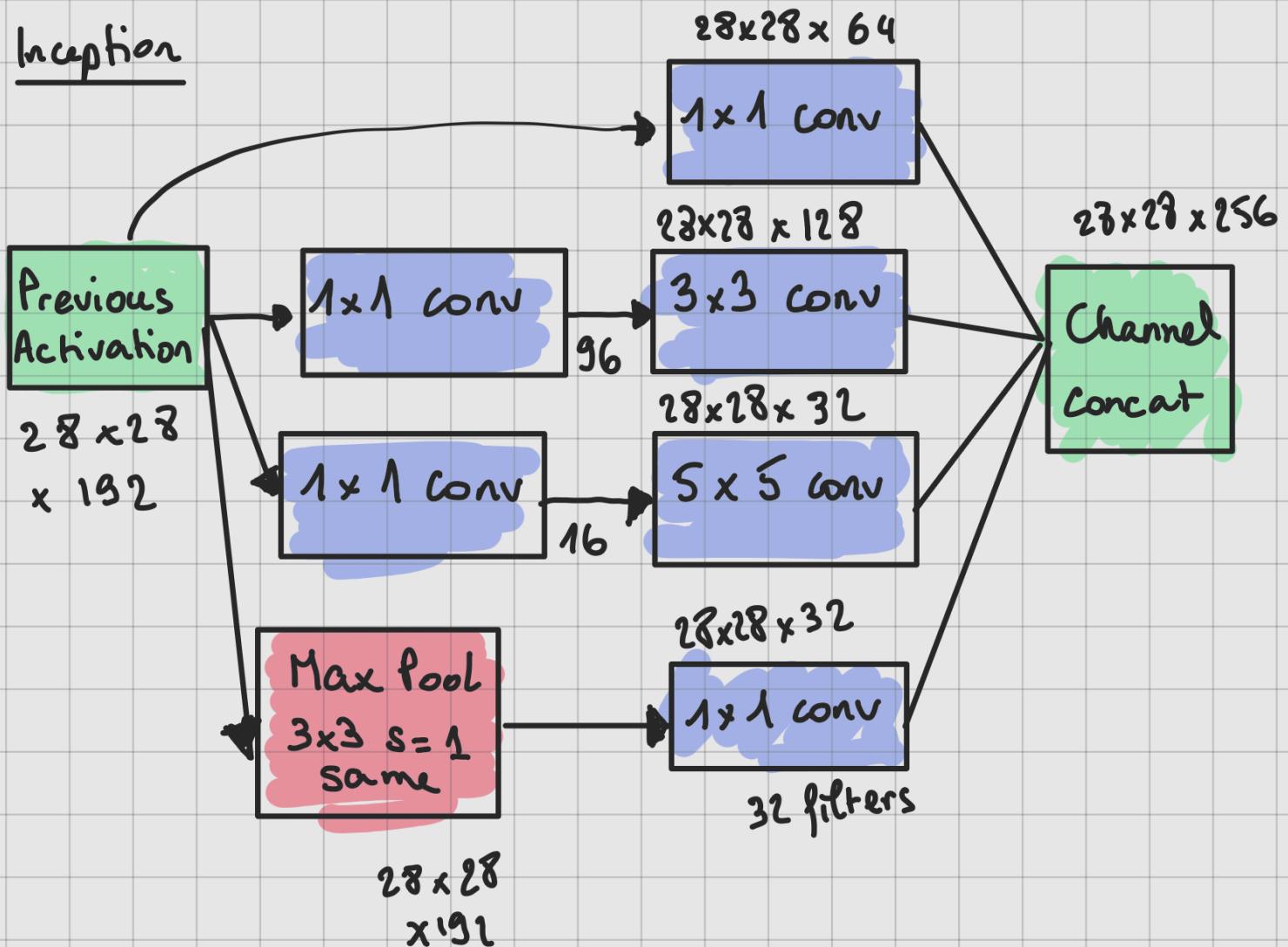
Solution ↴

Using 1×1 convolution



Andrew Ng

Inception

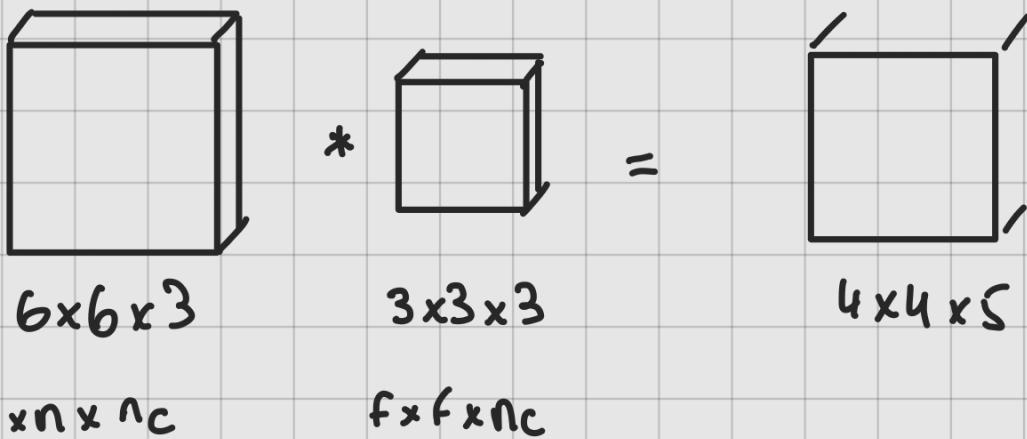


googleNet

MobileNet (2017 Howard)

- Low computational cost at deployment
- Useful for mobile and embedded vision apps
- Key idea: Normal vs depthwise-separable conv

Normal conv:

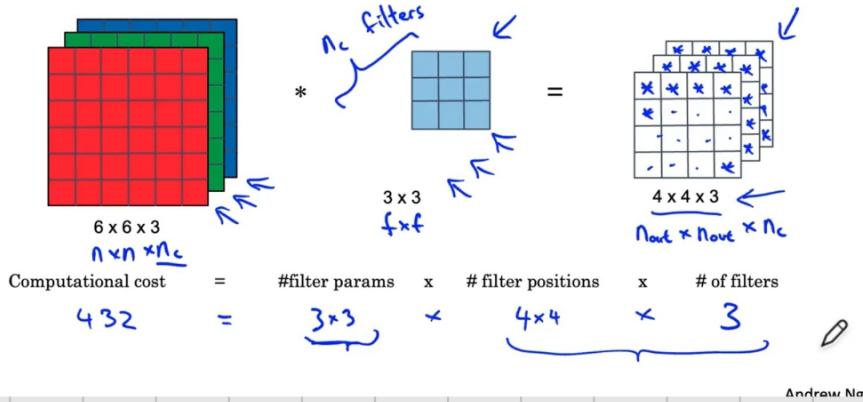


$$\text{Computational cost} = \cancel{\text{filter params}} \times \cancel{\text{filter position}} \times \cancel{\text{filters}} \\ 3 \times 3 \times 3 \quad \times \quad 4 \times 4 \quad \times \quad 5 \\ = 2160$$

Depthwise Separable Convolution

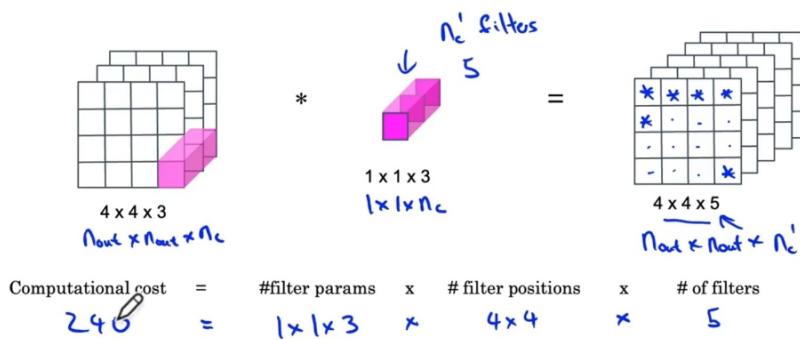
2 steps :

Depthwise Convolution



432

Pointwise Convolution



240

Total : 672 vs Normal : 240

$$\frac{672}{240} = 0,31$$

General rule :

$$= \frac{1}{n_c'} + \frac{1}{f^2}$$

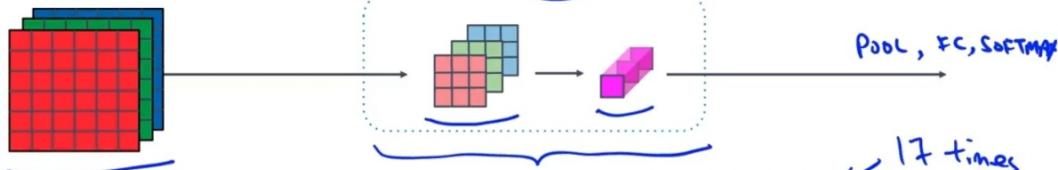
$$\exp = \frac{1}{512} + \frac{1}{3^2}$$

~ 10 times cheaper

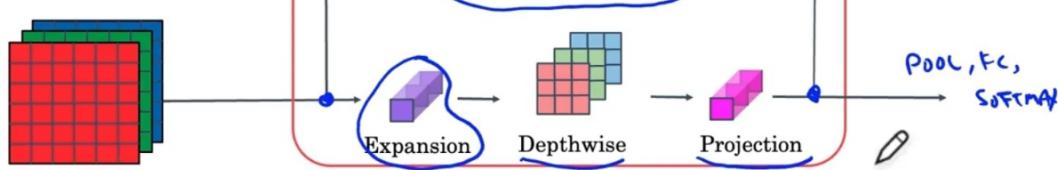
Archi

MobileNet

MobileNet v1

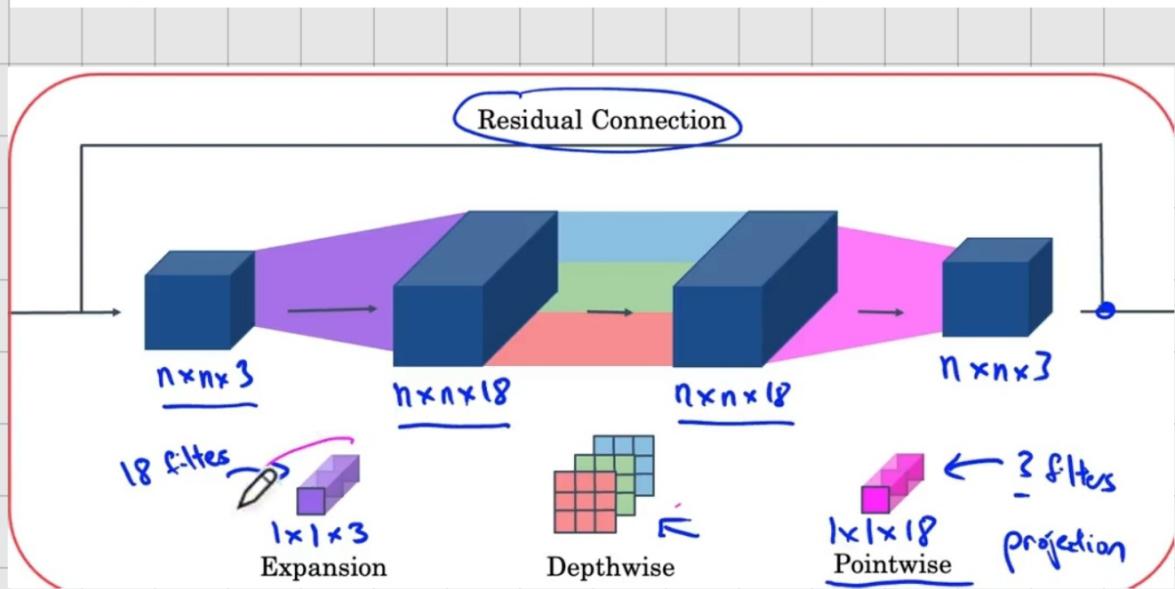


MobileNet v2



[Sandler et al. 2019, MobileNetV2: Inverted Residuals and Linear Bottlenecks]

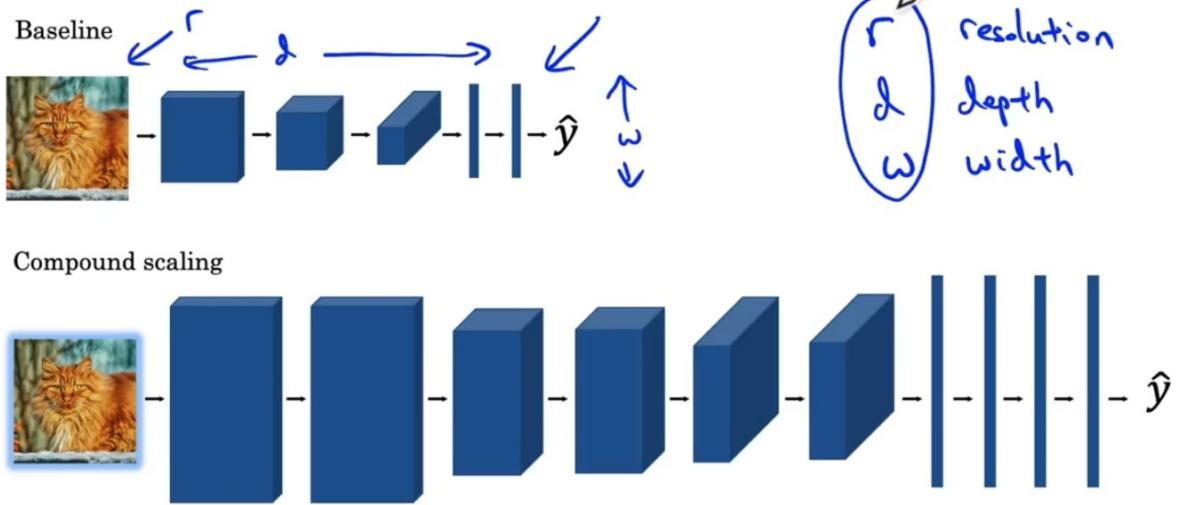
Andrew Ng



$$30 \times 5 + 20 \times 30$$

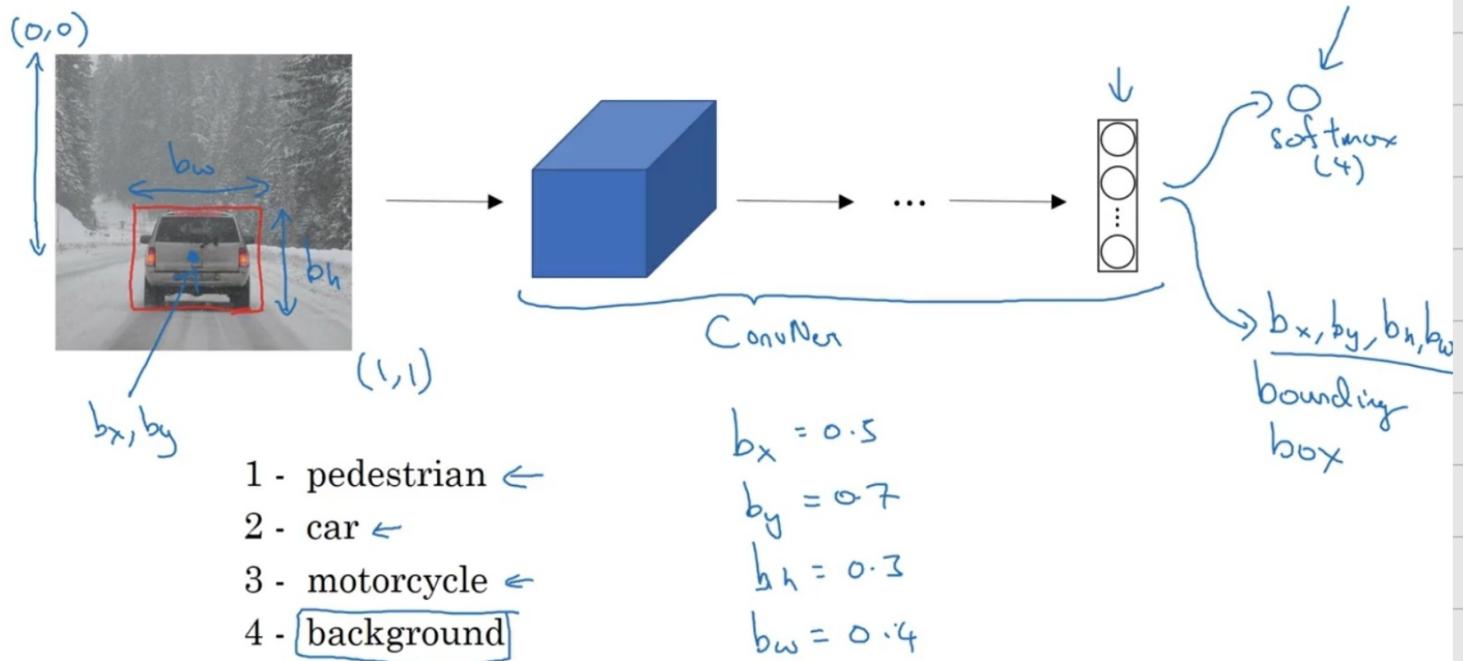
EfficientNet (2019 Tan and Le)

EfficientNet



Object Localization

Classification with localization



1 object

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \leftarrow \text{is there any object?}$$

No obj :

$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \leftarrow \text{"don't care"}$$

$$\mathcal{L} = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \\ \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

↑
for squared error

Can use | log like feature loss for C_1, C_2, C_3
 | squared error for bounding box
 | logistic reg loss for P_c

Landmark Detection

←

Landmark Detection

Landmark detection



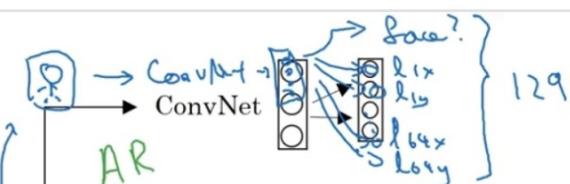
b_x, b_y, b_h, b_w



$\begin{matrix} l_{1x}, l_{1y}, \\ l_{2x}, l_{2y}, \\ l_{3x}, l_{3y}, \\ l_{4x}, l_{4y}, \\ \vdots \\ l_{64x}, l_{64y} \end{matrix} \} X, Y$



$\begin{matrix} l_{1x}, l_{1y}, \\ \vdots \\ l_{l}x, l_{l}y \end{matrix}$



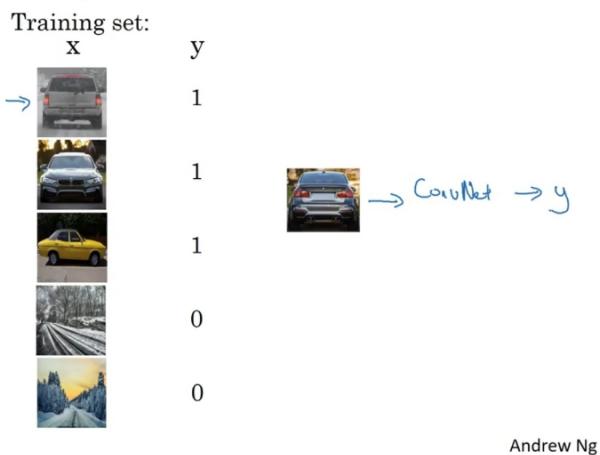
Object Detection

Sliding window detection:

←

Object Detection

Car detection example



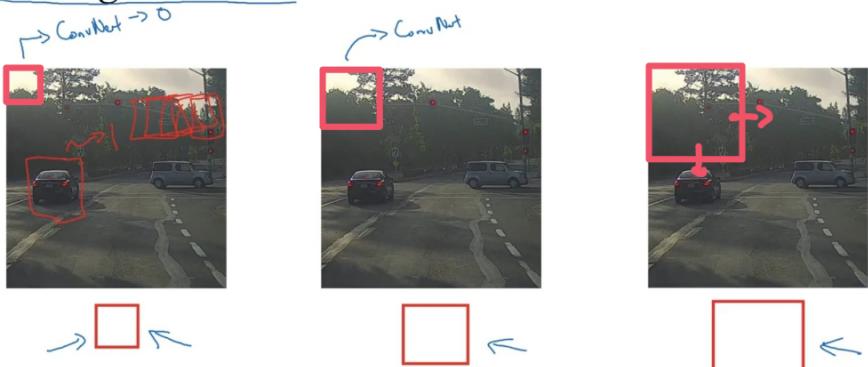
1) Train

with car that
take entire image

←

Object Detection

Sliding windows detection



2) Crop part of

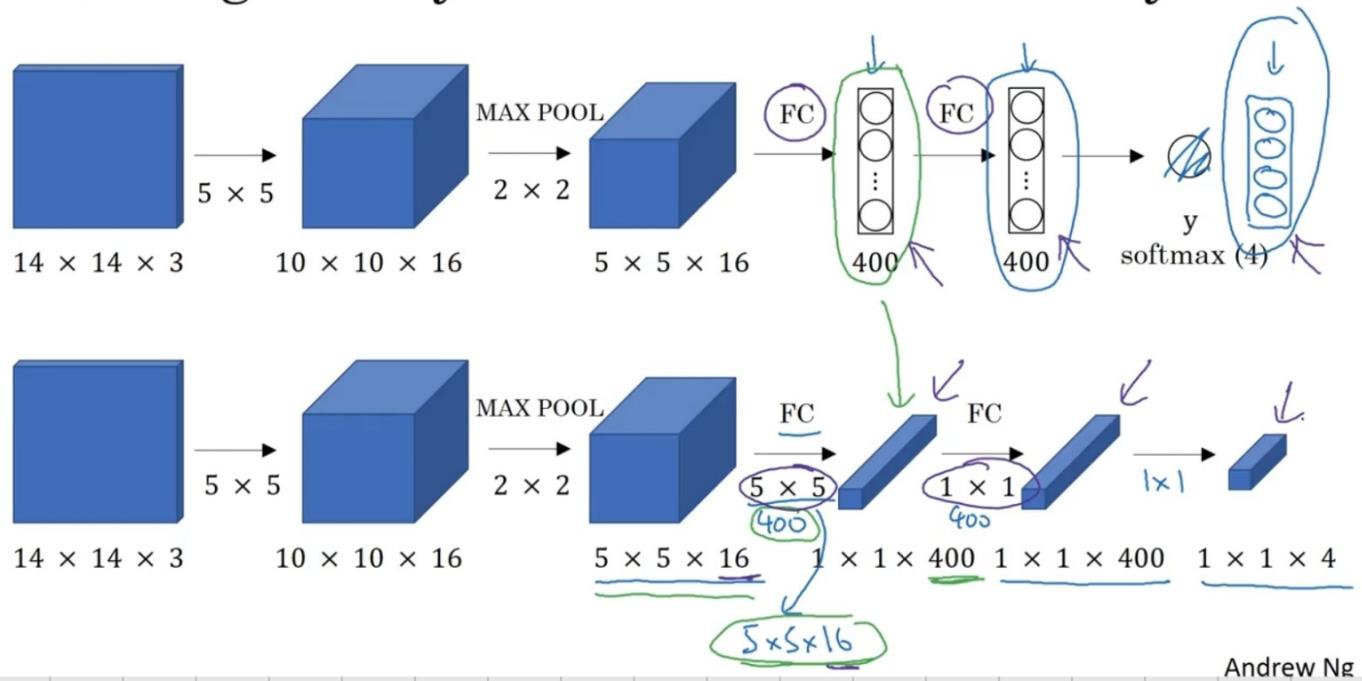
the image and
predict car or not

disadvantage: Computational cost

Convolutional Implementation of Sliding Windows

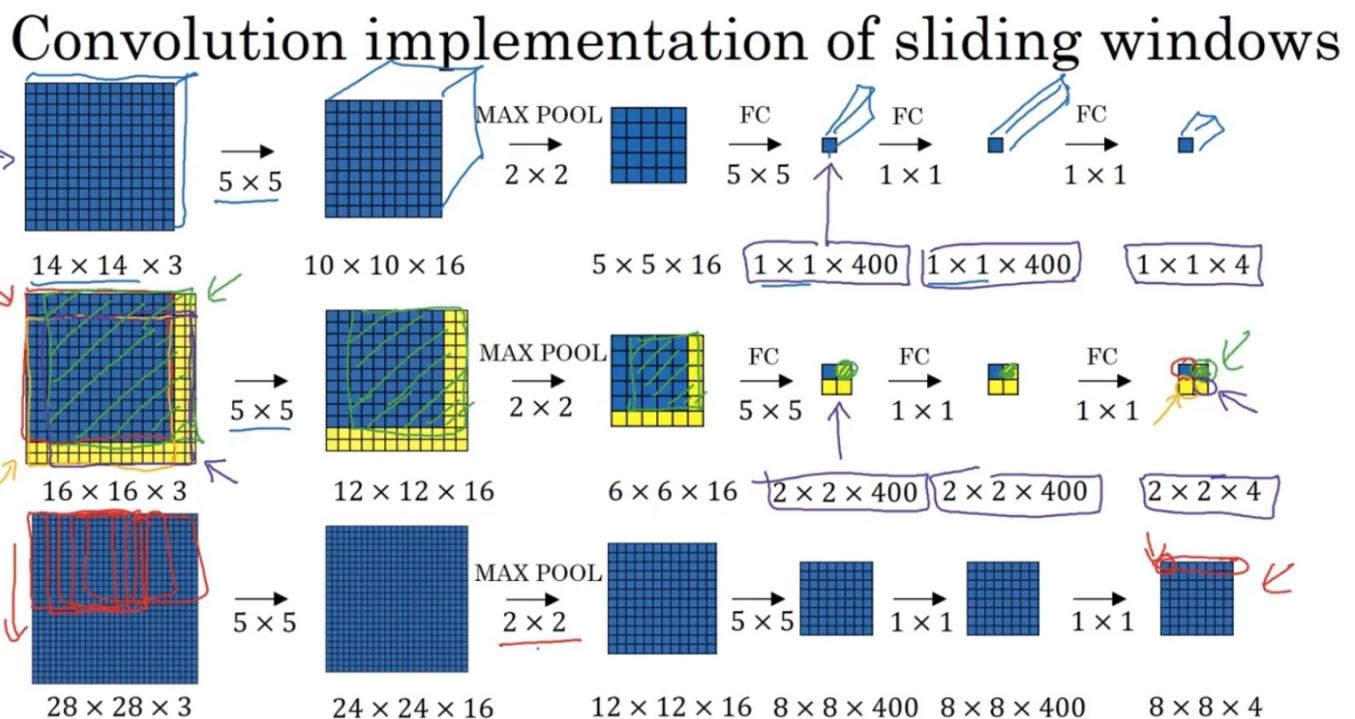
Turning FC layer into conv layer

Turning FC layer into convolutional layers



←

Convolutional Implementation of Sliding Windows

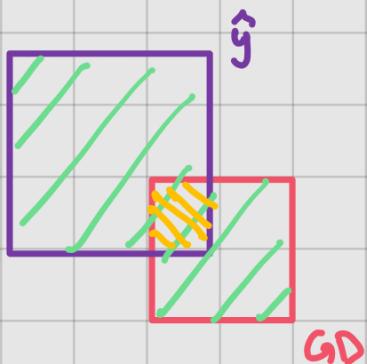


[Sermanet et al., 2014, OverFeat: Integrated recognition, localization and detection using convolutional networks]

Intersection vs Union

(IoU)

$$= \frac{\text{size of } I}{\text{size of } U}$$



"correct" if $\text{IoU} > 0.5$

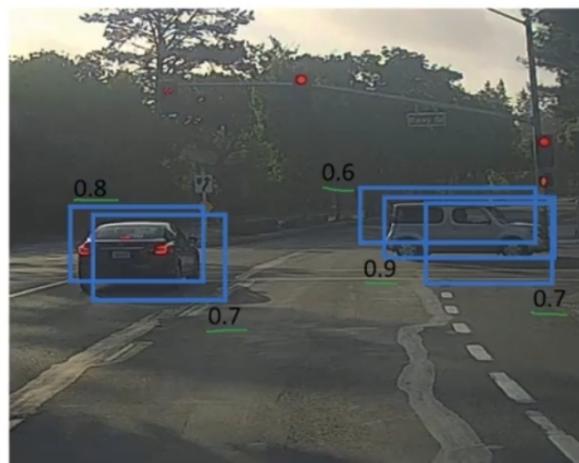
0.6

..

Non-max Suppression

Detect each object only once

Non-max suppression example

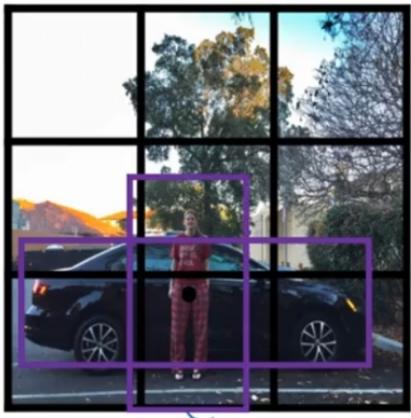


Andrew Ng

Anchor Boxes

for overlapping obj

Anchor box example



Anchor box 1: Anchor box 2:



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Andrew Ng

YOLO Algorithm

Outputting the non-max suppressed outputs



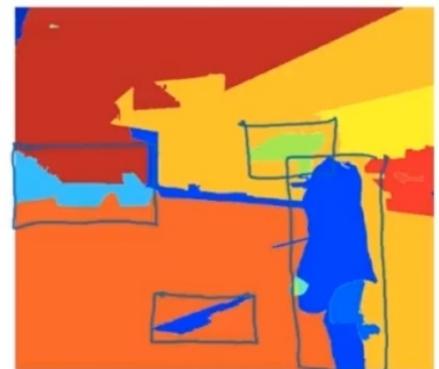
- For each grid cell, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

Andrew Ng

R-CNN

Region CNN

Region proposal: R-CNN



Segmentation algorithm

~2000

[Girshik et. al, 2013, Rich feature hierarchies for accurate object detection and semantic segmentation] Andrew Ng

Faster algorithms

→ R-CNN: Propose regions. Classify proposed regions one at a time. Output label + bounding box. ←

Fast R-CNN: Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions. ←

Faster R-CNN: Use convolutional network to propose regions.

[Girshik et. al, 2013. Rich feature hierarchies for accurate object detection and semantic segmentation]

[Girshik, 2015. Fast R-CNN]

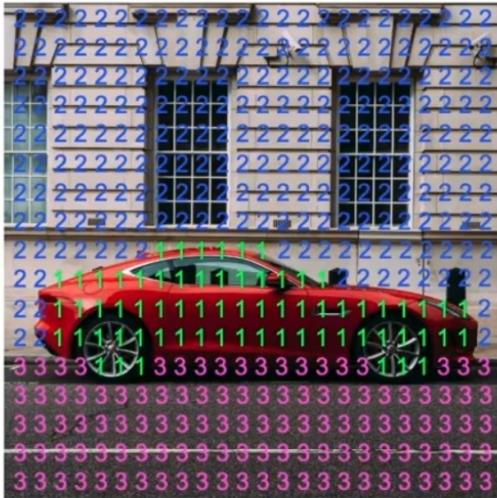
[Ren et. al, 2016. Faster R-CNN: Towards real-time object detection with region proposal networks]

Andrew Ng

Semantic Segmentation with U-Net

Per-pixel class labels

↓



1. Car
 2. Building
 3. Road

Segmentation Map

Andrew Ng

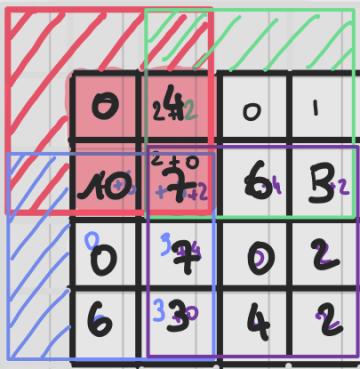
Transpose Convolution

Normal conv :

$$\begin{array}{c} \boxed{6 \times 6} \\ \times 3 \end{array} * \begin{array}{c} \boxed{3 \times 3} \\ \times 3 \\ \times 5 \end{array} = \begin{array}{c} \boxed{4 \times 4} \\ \times 5 \end{array}$$

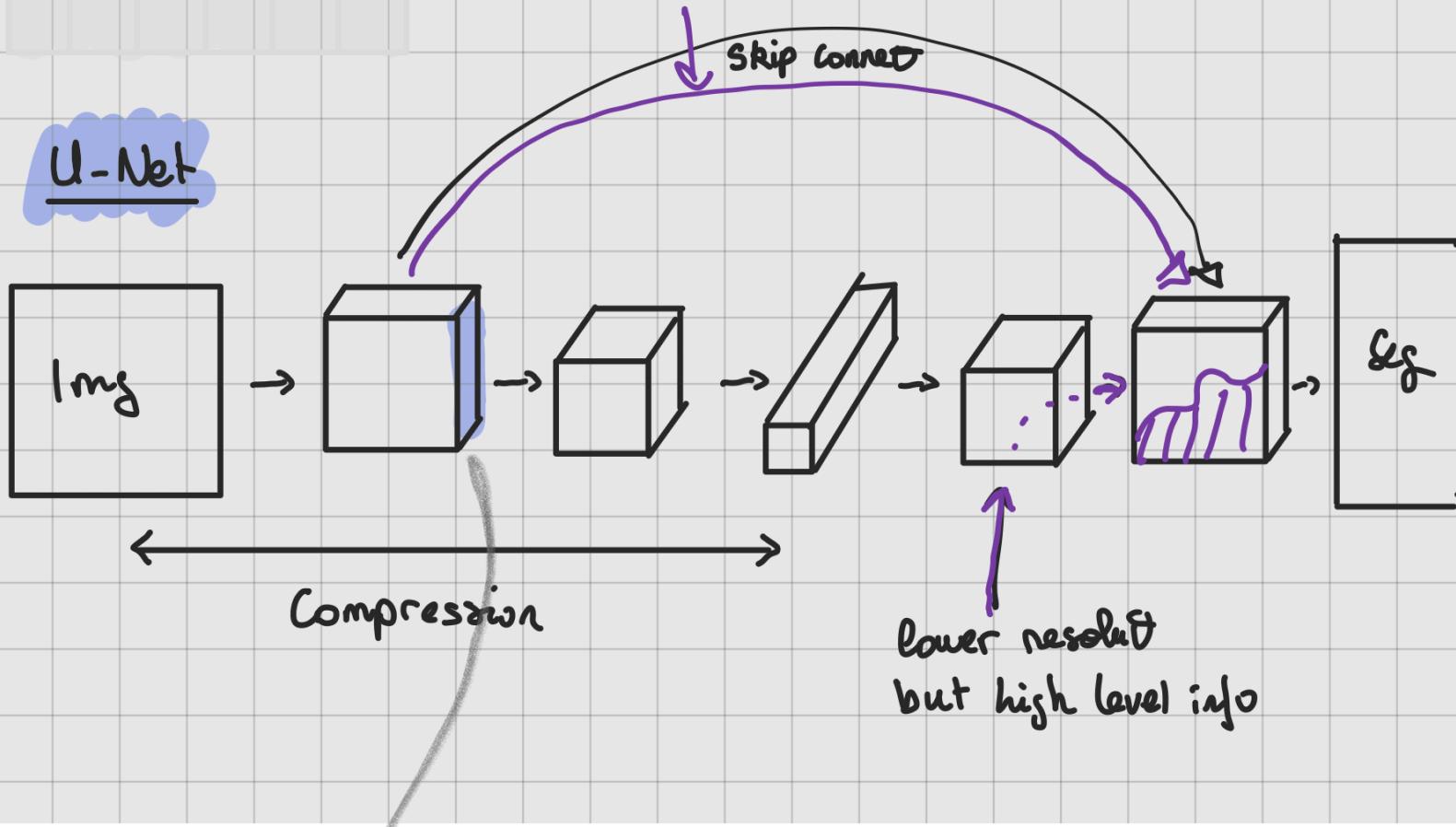
Transpose conv:

$$\begin{array}{|c|c|} \hline
 2 & 1 \\ \hline
 3 & 2 \\ \hline
 \end{array}
 \ast
 \begin{array}{|c|c|} \hline
 1 & 2 \\ \hline
 3 & 1 \\ \hline
 2 & 2 \\ \hline
 2 & 0 \\ \hline
 3 & 1 \\ \hline
 2 & 2 \\ \hline
 0 & 2 \\ \hline
 3 & 1 \\ \hline
 \end{array}$$

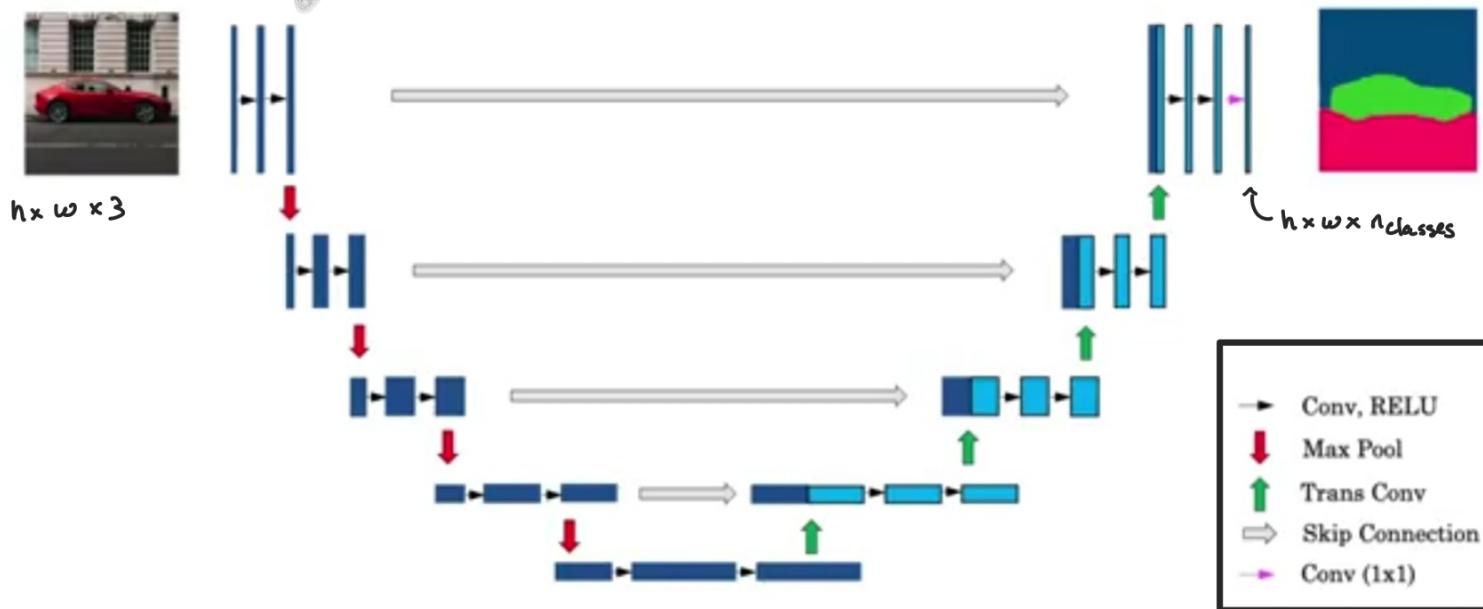


low level info
but high resolution

U-Net



U-Net



Face Recognition

Face Verification :

- Input image, name / ID 1:1
 - Output whether the input = claimed person

Face reco :

- Has a database of K persons $1:K$
 - Output ID if the image is any of the K persons

Learning a “similarity” function

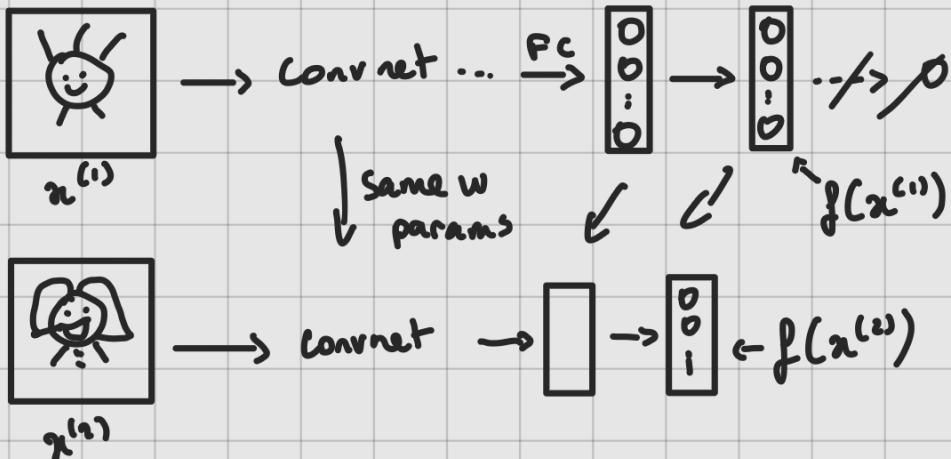
→ $d(\text{img1}, \text{img2})$ = degree of difference between images

If $d(\text{img1}, \text{img2}) \leq \tau$ "same"
 $> \tau$ "different" } Verification.



Andrew Ng

Siamese Network



$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

Learn parameters so that:

If $x^{(i)}, x^{(j)}$ same person $\rightarrow d(x^{(i)}, x^{(j)})$ is small

If $x^{(i)}, x^{(j)}$ diff person $\rightarrow d(x^{(i)}, x^{(j)})$ is large

Triplet Loss

Learning Objective



$$\text{Want: } \frac{\|f(A) - f(P)\|^2}{d(A, P)} + \alpha \leq 0.2$$

$$\frac{\|f(A) - f(N)\|^2}{d(A, N)} + \alpha \leq 0.7$$

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

[Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering]

Andrew Ng

Given 3 images A, P, N :

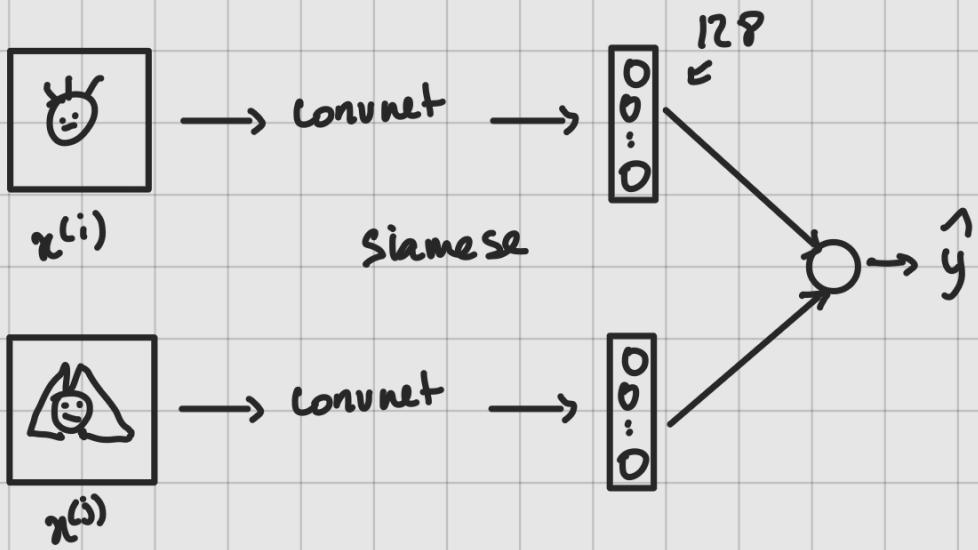
$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

During training, if A, P, N are chosen randomly, its easy to satisfy

→ Choose triplets that're "hard" to train on

$$\hookrightarrow d(A, P) \approx d(A, N)$$

Other way: Binary Classification For face verif



$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

Neural Style Transfer

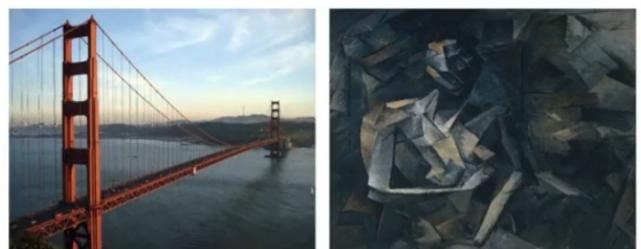
Neural style transfer



Content (C) Style (S)



Generated image (G)



Content (C) Style (S)



Generated image (G)

[Images generated by Justin Johnson]

Andrew Ng

Cost Function:

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

Find G :

1. Initiate G randomly ($G: 100 \times 100 \times 3$)
2. Use Gradient Descent to min $J(G)$

Content Cost Function

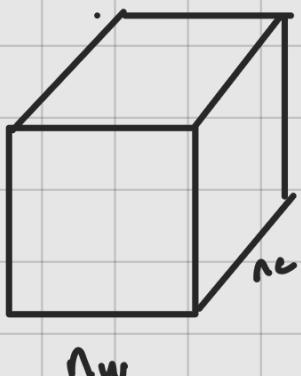
- Say you use hidden layer l to compute content cost
(l neither too shallow nor too deep)
- Use pre-trained ConvNet (E.g. VGG)
- Let $a^{[l](c)}$ and $a^{[l](g)}$ activation of layer l
- If $a^{[l](c)}$ and $a^{[l](g)}$ are similar \rightarrow img similar

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[l](c)} - a^{[l](g)} \|_2^2$$

normalize

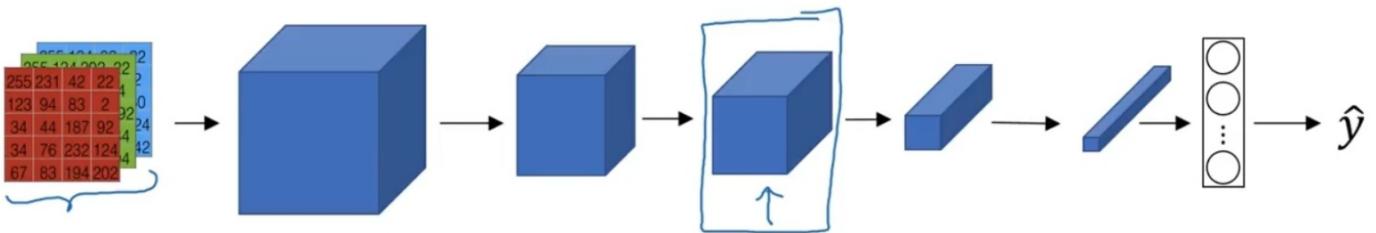
Style Cost Function

Say you are using layer l 's activation to measure "style." Define style as correlation between activation across channels.

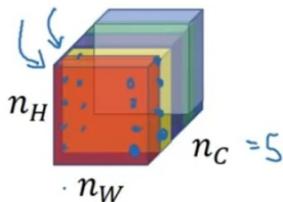


How correlated are the activations across channels?

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.

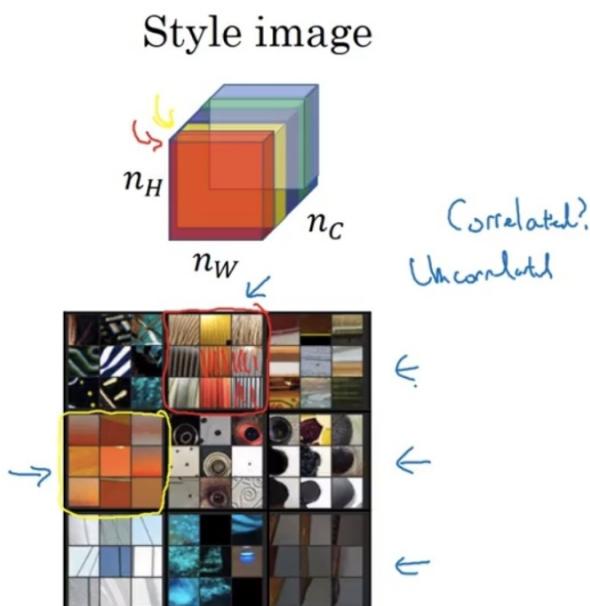


How correlated are the activations across different channels?

[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

Intuition about style of an image



Generated Image

[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

Style matrix

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

$$\begin{matrix} h & w & c \\ \downarrow & \downarrow & \downarrow \\ l & l & l \end{matrix}$$

$$G_{k,k'}^{[l]}$$

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

$$G_{kk'}^{[e](s)} = \sum_i^{n_h} \sum_j^{n_w} a_{ijk}^{[e](s)} a_{ijk'}^{[e](s)}$$

"Gram matrix"

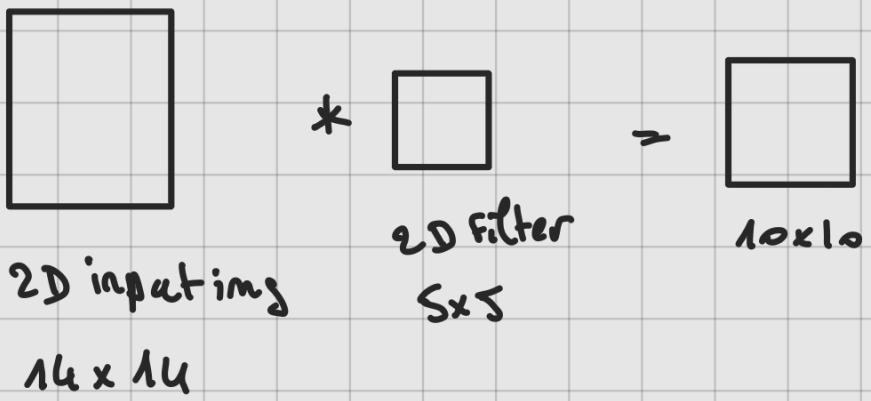
$$G_{kk'}^{[e]G} = \sum_i^{n_h} \sum_j^{n_w} a_{ijk}^{[e](G)} a_{ijk'}^{[e](G)}$$

$$J_{\text{Style}}^{[e]}(s, G) = \frac{1}{2n_h^{[e]} n_w^{[e]} n_c^{[e]}} \| G^{[e](s)} - G^{[e](G)} \|_F^2$$

normalize

$$J_{\text{Style}}(s, G) = \sum_e \lambda^{[e]} J_{\text{Style}}^{[e]}(s, G)$$

1D and 3D Generalized



Convolutions in 2D and 1D

