Case Studies

Why look at case studies?

deeplearning.ai

# Outline

## Classic networks:

- LeNet-5 $\leftarrow$

- AlexNet $\leftarrow$

- VGG $\leftarrow$

ResNet    (152)

Inception
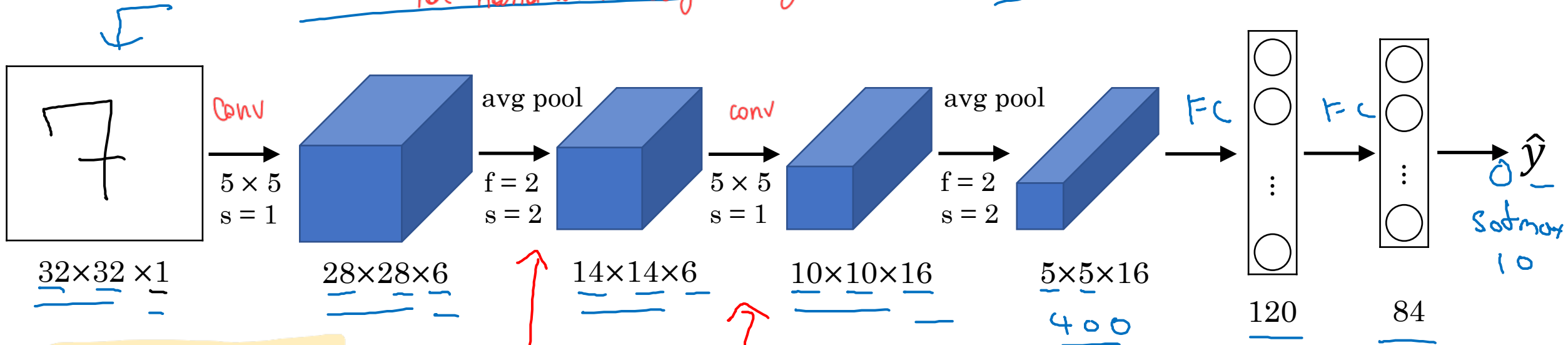
deeplearning.ai

Case Studies

Classic networks

# LeNet - 5

- One of the first ones
- For hand-written digit recognition



32×32×1    28×28×6    14×14×6    10×10×16    5×5×16

Conv 5 × 5, s = 1    avg pool f = 2, s = 2    conv 5 × 5, s = 1    avg pool f = 2, s = 2    FC    FC    $\hat{y}$ softmax 10

400    120    84

60K parameters.

nontinearity afte pooling

$n_H \times n_W \times n_c$    $f \times f \times n_c$

$n_H, n_W \downarrow$    $n_c \uparrow$

conv pool conv pool fc fc output

Advanced: Sigmoid/tanh    ReLU

$\boxed{II}$, $\underline{III}$.

[LeCun et al., 1998. Gradient-based learning applied to document recognition]

Andrew Ng

# AlexNet

- Inspired by leNet-5, with 1k-times more parameters



$11 \times 11$
$s = 4$

$227 \times 227 \times 3$

$55 \times 55 \times 96$

MAX-POOL

$3 \times 3$
$s = 2$

$27 \times 27 \times 96$

$5 \times 5$
same

$27 \times 27 \times 256$

MAX-POOL

$3 \times 3$
$s = 2$

$13 \times 13 \times 256$

$3 \times 3$
same

$13 \times 13 \times 384$

$3 \times 3$
same

$13 \times 13 \times 384$

$3 \times 3$
some

$13 \times 13 \times 256$

MAX-POOL

$3 \times 3$
$s = 2$

$6 \times 6 \times 256$

9216

=

FC          FC

9216        4096        4096

Softmax
1000

- Similary to LeNet, but much bigger.
- ReLU
- Multiple GPUs.
(not used today) - Local Response Normalization (LRN)

13    13    256

160M Parameters

[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]
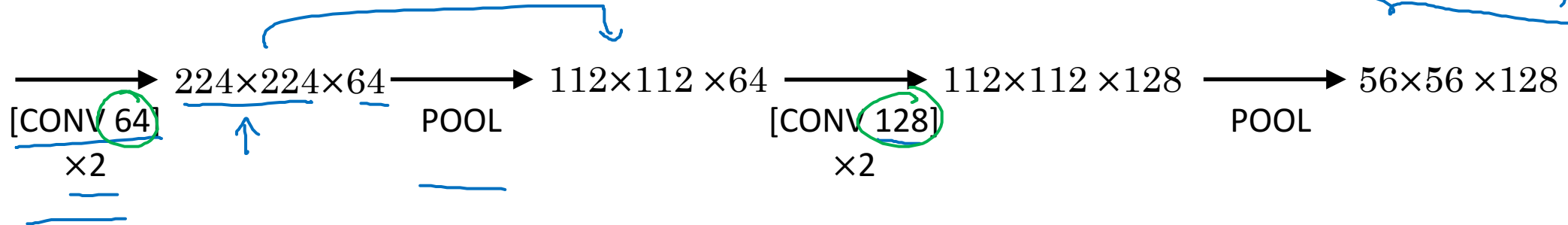
Andrew Ng

# VGG - 16
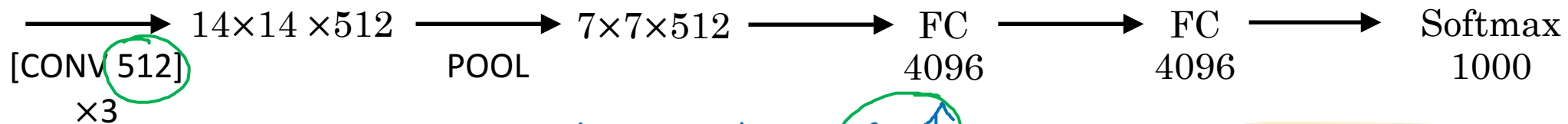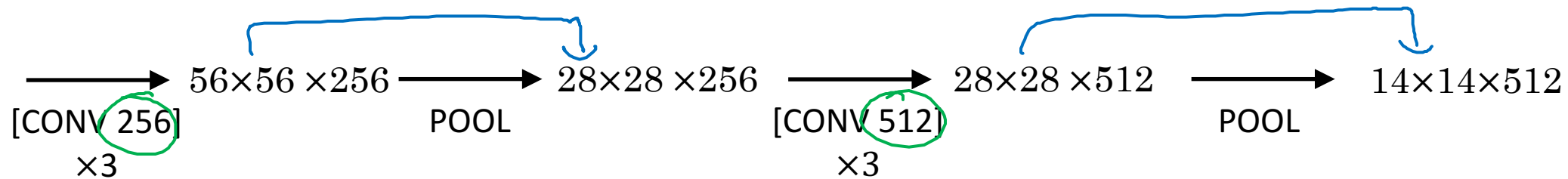
CONV = 3×3 filter, s = 1, same          MAX-POOL = 2×2 , s = 2

VGG-19

- Relatively uniform architecture compared to the previous ones. Only uses CONV(3x3) and MAXPOOL (2x2) operations. But they do it many times.
- Activations increase in depth by a factor of 2 due to conv and shrink in size by a factor of 2 due to pooling

224×224×3

224×224×64 → 112×112 ×64 → 112×112 ×128 → 56×56 ×128

[CONV 64]    POOL         [CONV 128]    POOL
×2                        ×2

224×224 ×3

56×56 ×256 → 28×28 ×256 → 28×28 ×512 → 14×14×512

[CONV 256]   POOL         [CONV 512]    POOL
×3                        ×3

14×14 ×512 → 7×7×512 → FC → FC → Softmax
                        4096   4096   1000

[CONV 512]   POOL
×3

$n_H, n_W \downarrow$         $n_c \uparrow$         ~138M

[Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition]          Andrew Ng

• Motivation: adding too many layers can worsen the training error. This is because redundant layers should learn an identity function, which is hard in general.

• This problem can be solved by adding "skip-connections" that forward the activation of one layer to the linear operation of a later layer. This solves the problem because redundant blocks can easily learn the identity function by setting their weights to zero.



$a^{[l]} \rightarrow a^{[l+1]} \rightarrow a^{[l+2]}$

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

RESIDUAL BLOCK

The residual block can implement identity by setting $W^{[l+2]}, b^{[l+2]}$ to zero.

The group of layers included within the skip connection is defined "residual block."

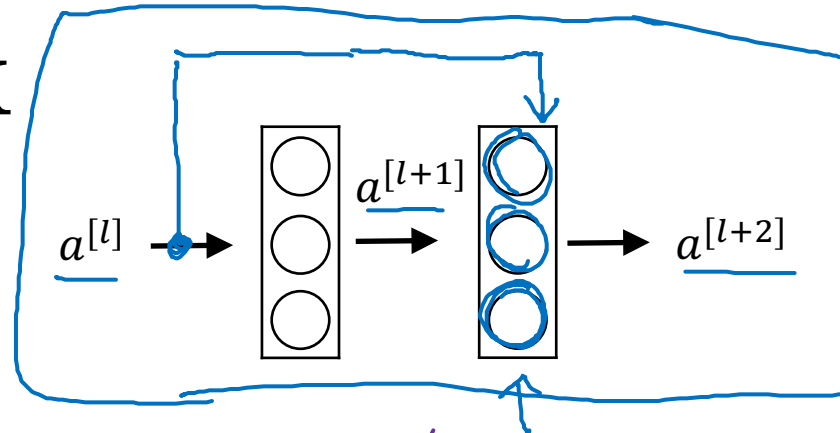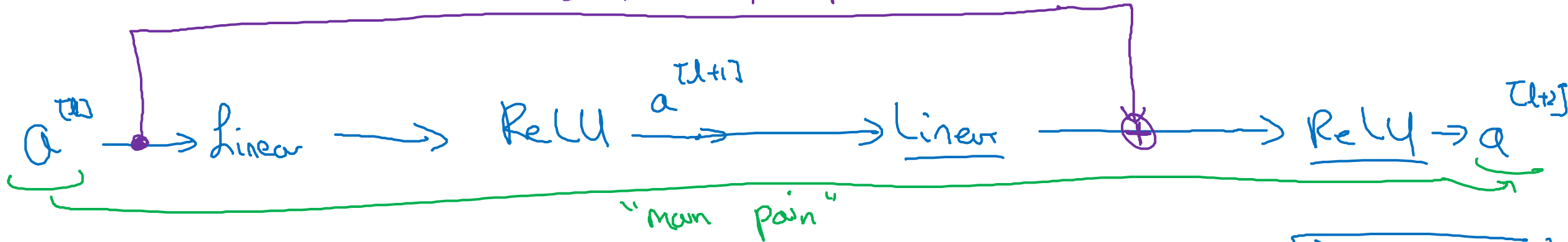• ResNets are very deep networks containing residual blocks.

# Case Studies

# Residual Networks (ResNets)

deeplearning.ai

In general the activations of all layers within a residual block have the same dimension. However, it is possible to have different dimensions if an intermediate "conjunction" matrix is used to sum $a^{[l]}$ to $z^{[l+2]}$.

# Residual block



"short cut" / skip connection
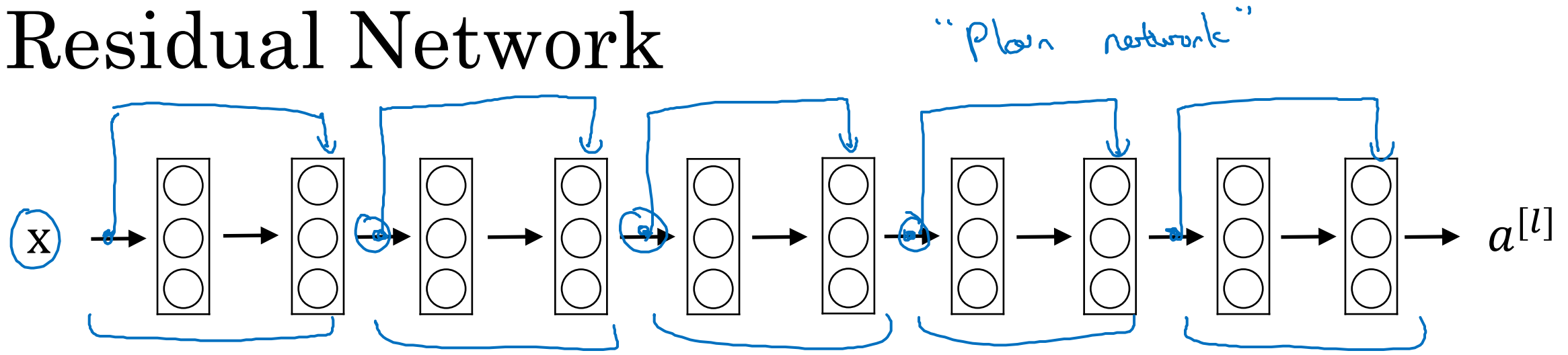
"main pain"

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

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

[He et al., 2015. Deep residual networks for image recognition]

Andrew Ng

# Residual Network



"Plain network"

$x$ → ... → $a^{[l]}$

Plain

residual blocks

ResNet

redundant layers struggle to learn the simple identity function.

"reality"

theory

training error vs # layers (Plain)

training error vs # layers (ResNet)

[He et al., 2015. Deep residual networks for image recognition]

Andrew Ng

Case Studies

Why ResNets work

deeplearning.ai

# Why do residual networks work?



$X \rightarrow \boxed{\text{Big NN}} \rightarrow a^{[l]}$

$X \rightarrow \boxed{\text{Big NN}} \rightarrow a^{[l]}$

$a^{[l+2]}$

"Same"

Identity function is easy for Residual block to learn!

to join activations with different sizes.

ReLU.  $a \geq 0$

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

$$= g\left( W^{[l+2]} a^{[l+1]} + b^{[l+2]} + W_s a^{[l]} \right) = g\left( a^{[l]} \right)$$

$$= a^{[l]}$$

256

$\mathbb{R}^{256 \times 128}$

128

If $W^{[l+2]} = 0$, $b^{[l+2]} = 0$

# ResNet

## Plain



## ResNet



same size activations

different size activations

$z^{[l+2]} + a^{[l]}$

3x3 same

Pool

Pool

$W_s$

[He et al., 2015. Deep residual networks for image recognition]
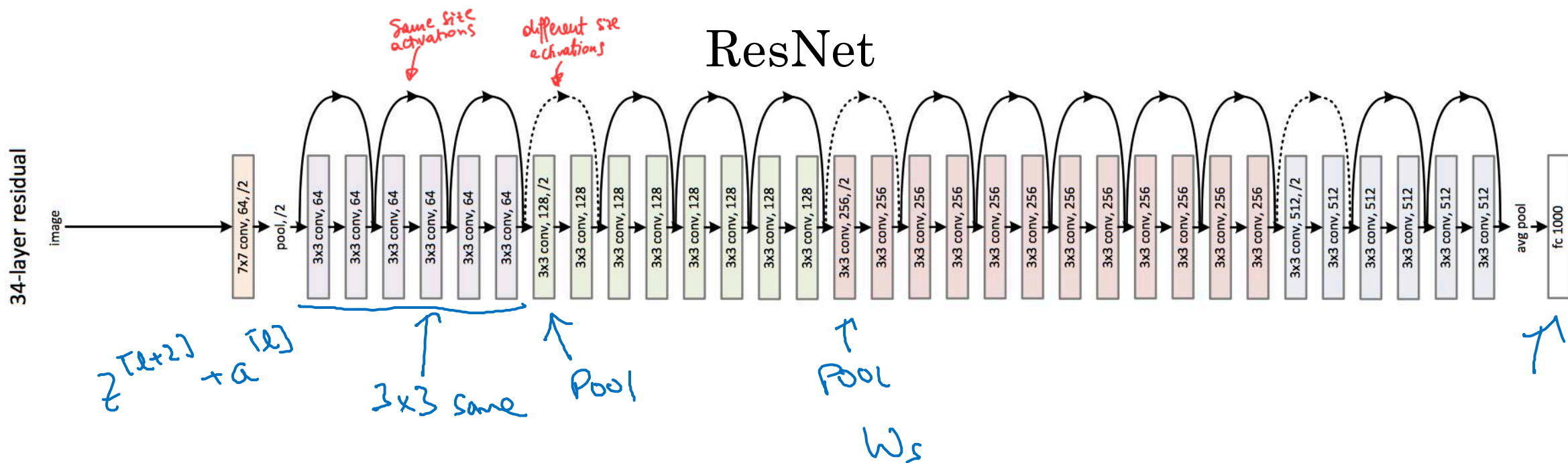
Andrew Ng

Applies a nontrivial function of the channels*

→ Used to shrink the number of channels in a nontrivial way (it's like a fancy version of pooling). This is especially useful for inception networks that tend to have extremely high number of channels.

# Case Studies

# Network in Network and 1×1 convolutions

deeplearning.ai
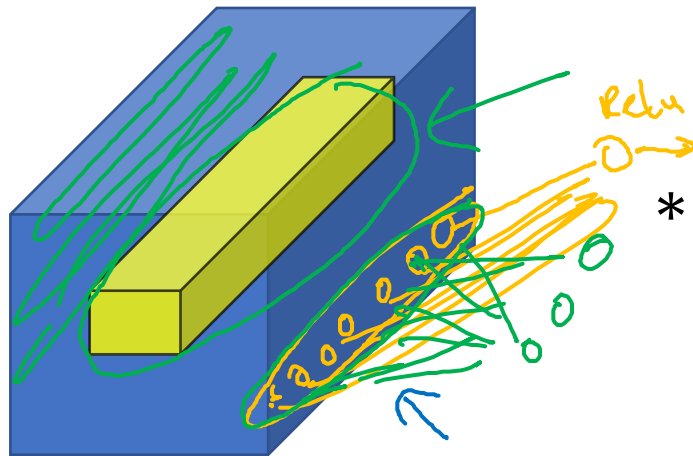
*"AKA:"Network in network" because convolving all the C channels of 1 neuron with a 1×1×C filter corresponds to a FC layer taking the C channels as inputs.

# Why does a 1 × 1 convolution do?

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

$6 \times 6$  × 1

$*$  $\boxed{2}$  $=$

| 2 | 4 | 6 | .. | | |
|---|---|---|----|-|-|
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

$6 \times 6 \times 32$

$32 \longrightarrow$ #filters.  $n_c^{[l+1]}$

Relu

$0 \rightarrow$

$*$  $1 \times 1 \times 32$  $=$

ReLU

Network in Network

$6 \times 6 \times$ # filters

[Lin et al., 2013. Network in network]

Andrew Ng

# Using 1×1 convolutions

Shrinking channels
(fancy pooling).



ReLU

CONV $1 \times 1$

~~32~~

192

$\rightarrow$ 28 × 28 × 32 ~~192~~

$1 \times 1 \times 192$

32 filters.

$n_H, n_W, n_c$

[Lin et al., 2013. Network in network]

Andrew Ng

- Idea: instead of choosing which filters to apply, apply them all (in parallel) and stack them.
- however this creates layers with many channels which increases enormously the complexity of convolutions. ○ The inception module is designed to save computation.

The idea is to replace each convolution with a 2 step operation:

1 - "compress the no. channels" by performing n CONV 1x1xno.channels where n is "small within reason".
⇒ RESULTS IN A "BOTTLENECK LAYER" where the information is compressed;

## Case Studies

2 - apply the convolution for the desired output dimension to the bottleneck layer.

## Inception network

This 2 step operation con reduce computations by a factor of 10 while obtaining similar results.

## motivation

Size of the bottleneck layer must be "small within reason."

deeplearning.ai

e.g: instead of   28x28x192 $\xrightarrow{\text{32 CONV 5x5}}$ 28x28x32 , do  28x28x192 $\xrightarrow{\text{16 CONV 1x1}}$ 28x28x16 $\xrightarrow{\text{32 CONV 5x5}}$ 28x28x32

about 120M operations     about 12M operations

# Motivation for inception network



$28 \times 28 \times 192$

$1 \times 1$    28×28×64

$3 \times 3$    28×28×128
Same

$5 \times 5$    28×28×32
Same

MAX-POOL    28×28×32
Same
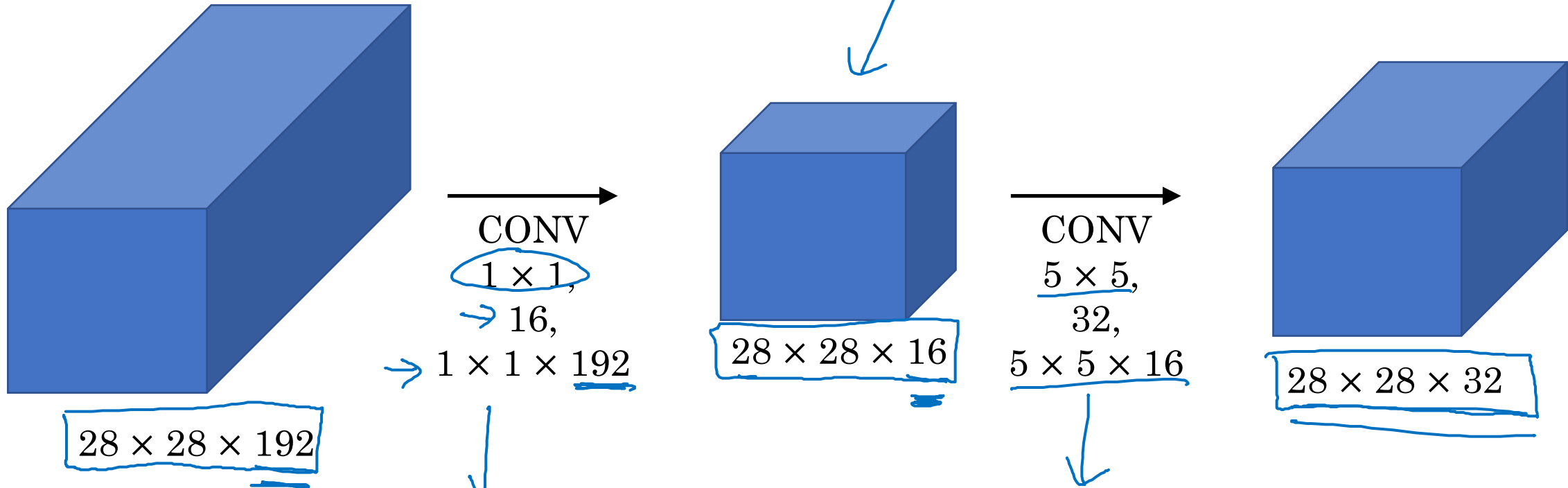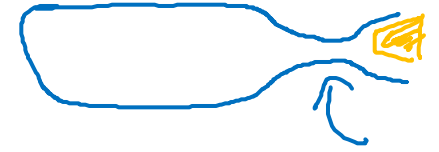S=1

28

28

32 32

128

64

256

$28 \times 28 \times 256$

cryptic: actually
maxpool returns same no. channels (192)
but typically this is followed by a 1×1 conv
to shrink the no. channel (in the inception architecture).

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

Andrew Ng

# The problem of computational cost



CONV
5 × 5,
same,
32

28 × 28 × 192

28 × 28 × 32

32 filters.     filters are $5 \times 5 \times 192$.

$28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120M.$

Andrew Ng

# Using 1×1 convolution



"bottleneck layer"

CONV
1 × 1,
→ 16,
→ 1 × 1 × 192

28 × 28 × 192

28 × 28 × 16

CONV
5 × 5,
32,
5 × 5 × 16

28 × 28 × 32

$28 \times 28 \times 16 \times 192 = 2.4M$

$28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10.0M$

12.4M

$120M \rightarrow$

Andrew Ng

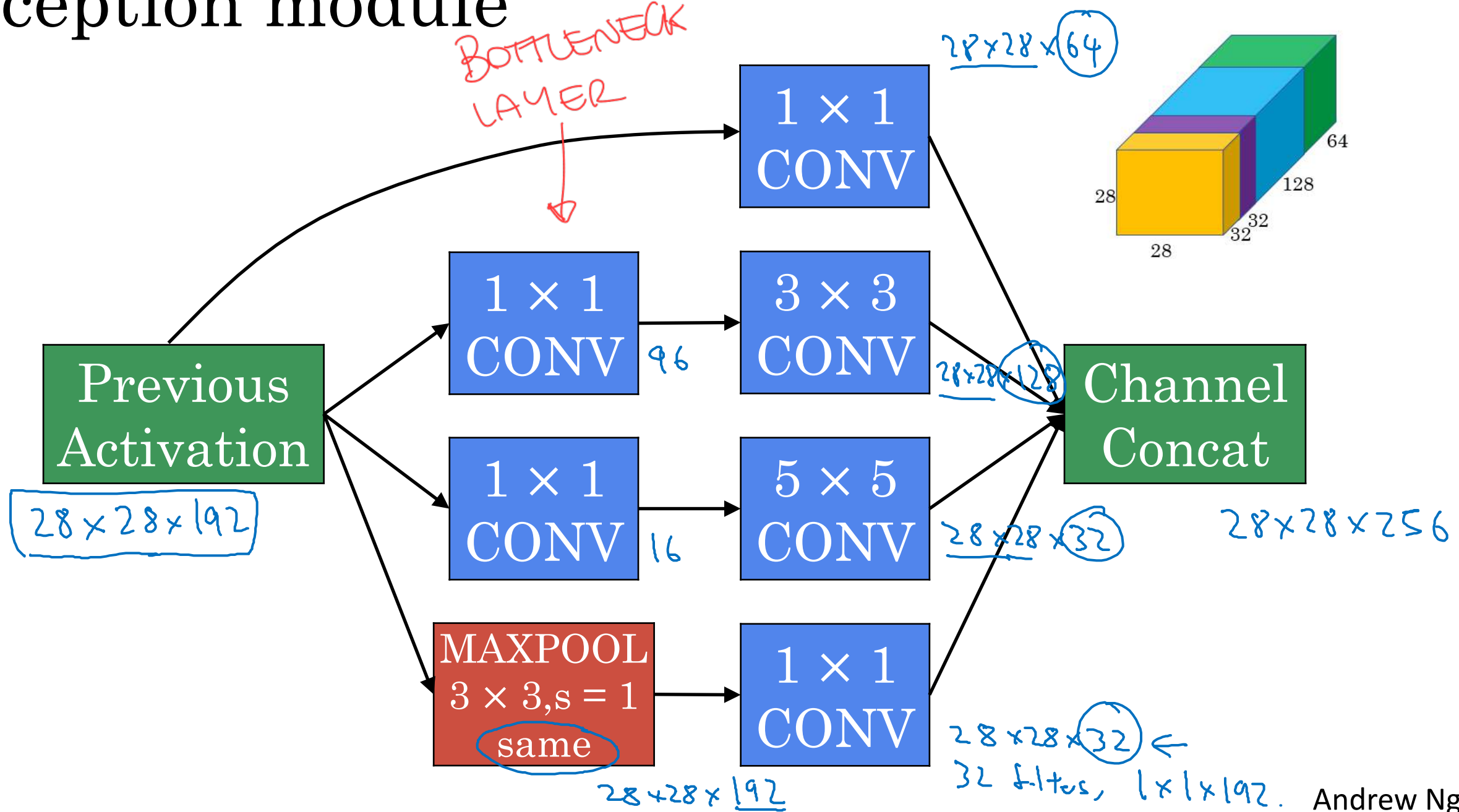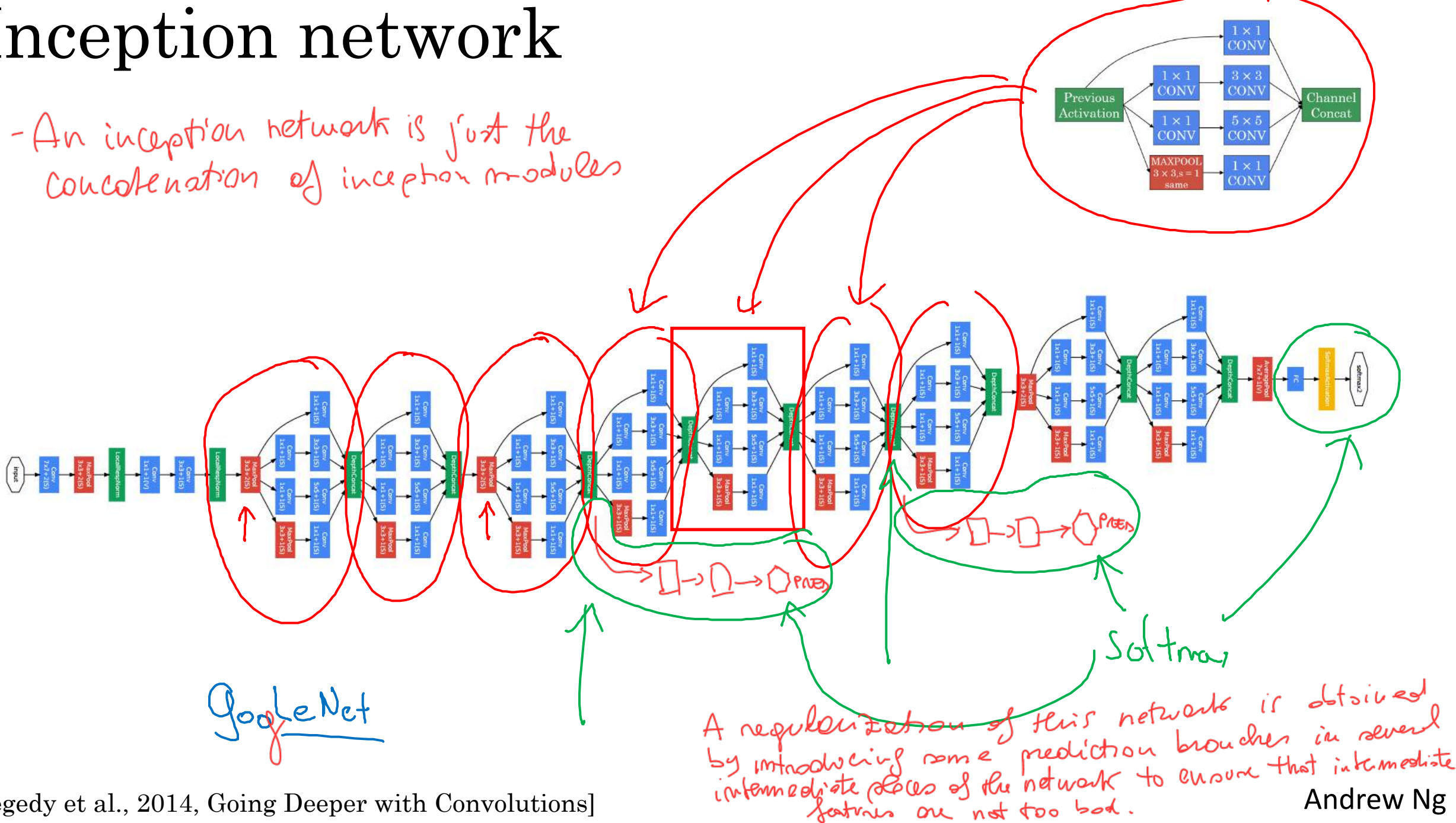Case Studies

Inception network

deeplearning.ai

# Inception module



BOTTLENECK LAYER

$1 \times 1$ CONV

$28 \times 28 \times 64$

$1 \times 1$ CONV 96

$3 \times 3$ CONV

$28 \times 28 \times 128$

Previous Activation

$28 \times 28 \times 192$

$1 \times 1$ CONV 16

$5 \times 5$ CONV

$28 \times 28 \times 32$

Channel Concat

$28 \times 28 \times 256$

MAXPOOL $3 \times 3, s = 1$ same

$1 \times 1$ CONV

$28 \times 28 \times 192$

$28 \times 28 \times 32$ ← 32 filters, $1 \times 1 \times 192$.

Andrew Ng

# Inception network

- An inception network is just the concatenation of inception modules



GoogLeNet

Softmax

A regularization of this network is obtained by introducing some prediction branches in several intermediate places of the network to ensure that intermediate features are not too bad.

[Szegedy et al., 2014, Going Deeper with Convolutions]

Andrew Ng

http://knowyourmeme.com/memes/we-need-to-go-deeper

Andrew Ng

# Practical advice for using ConvNets

## Transfer Learning

deeplearning.ai

# Transfer Learning

Tigger    Misty    Neither

T
M
N

softmax

x → [network layers] → softmax / logit → ŷ

freeze    trainableParams = 0    freeze = 1

save to disk

x → [network layers] → O → ŷ (T, M, N)

freeze

train

x → [network layers] → ŷ

train

Andrew Ng

deeplearning.ai

Practical advice for using ConvNets

Data augmentation

# Common augmentation method

Mirroring



Random Cropping



Rotation

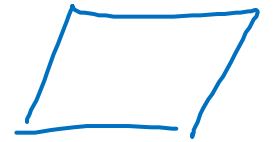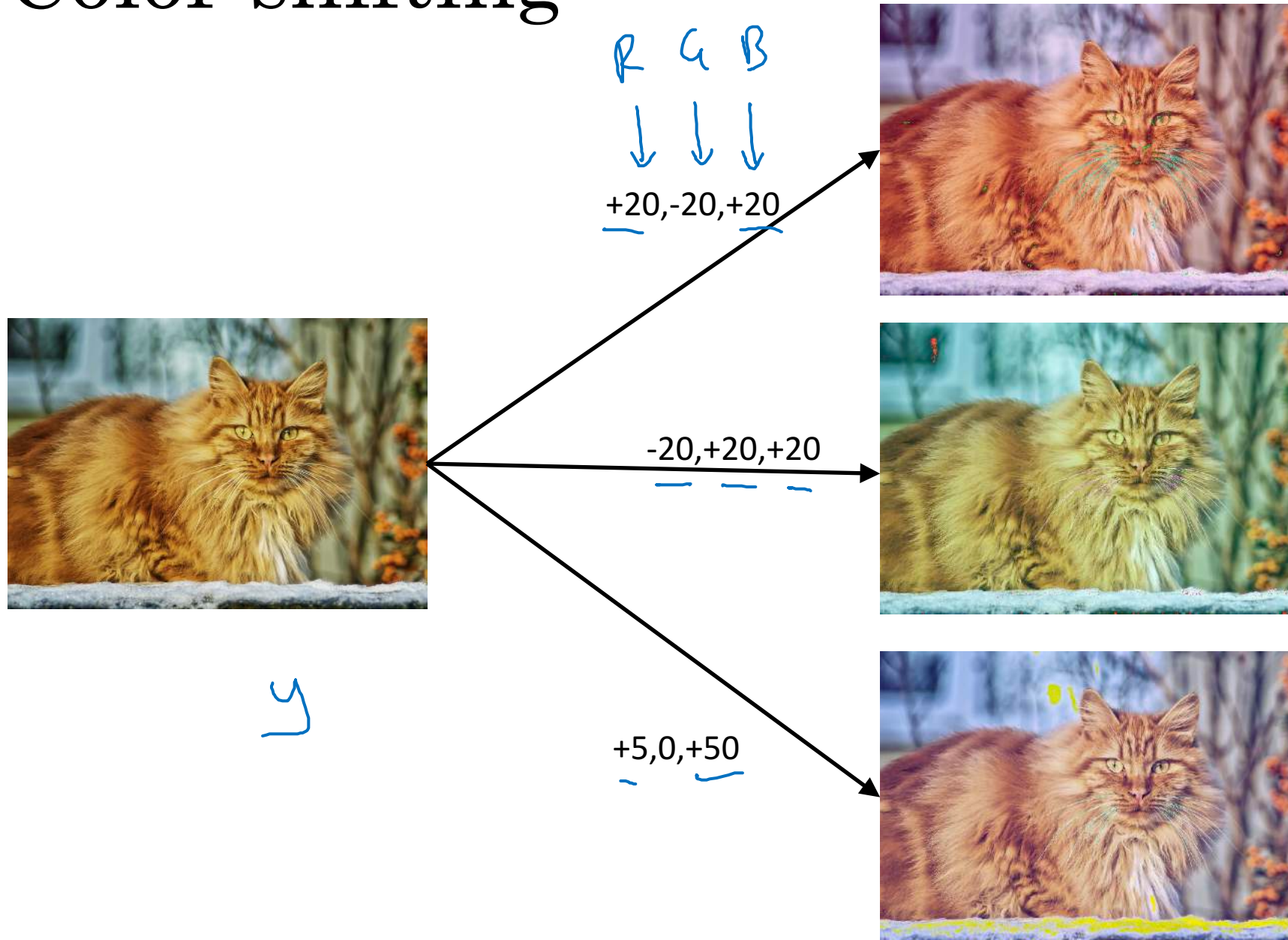Shearing

Local warping

...

Andrew Ng

# Color shifting



R  G  B
↓  ↓  ↓
+20,-20,+20

-20,+20,+20

+5,0,+50

y

Advanced:
PCA
ml-class.org
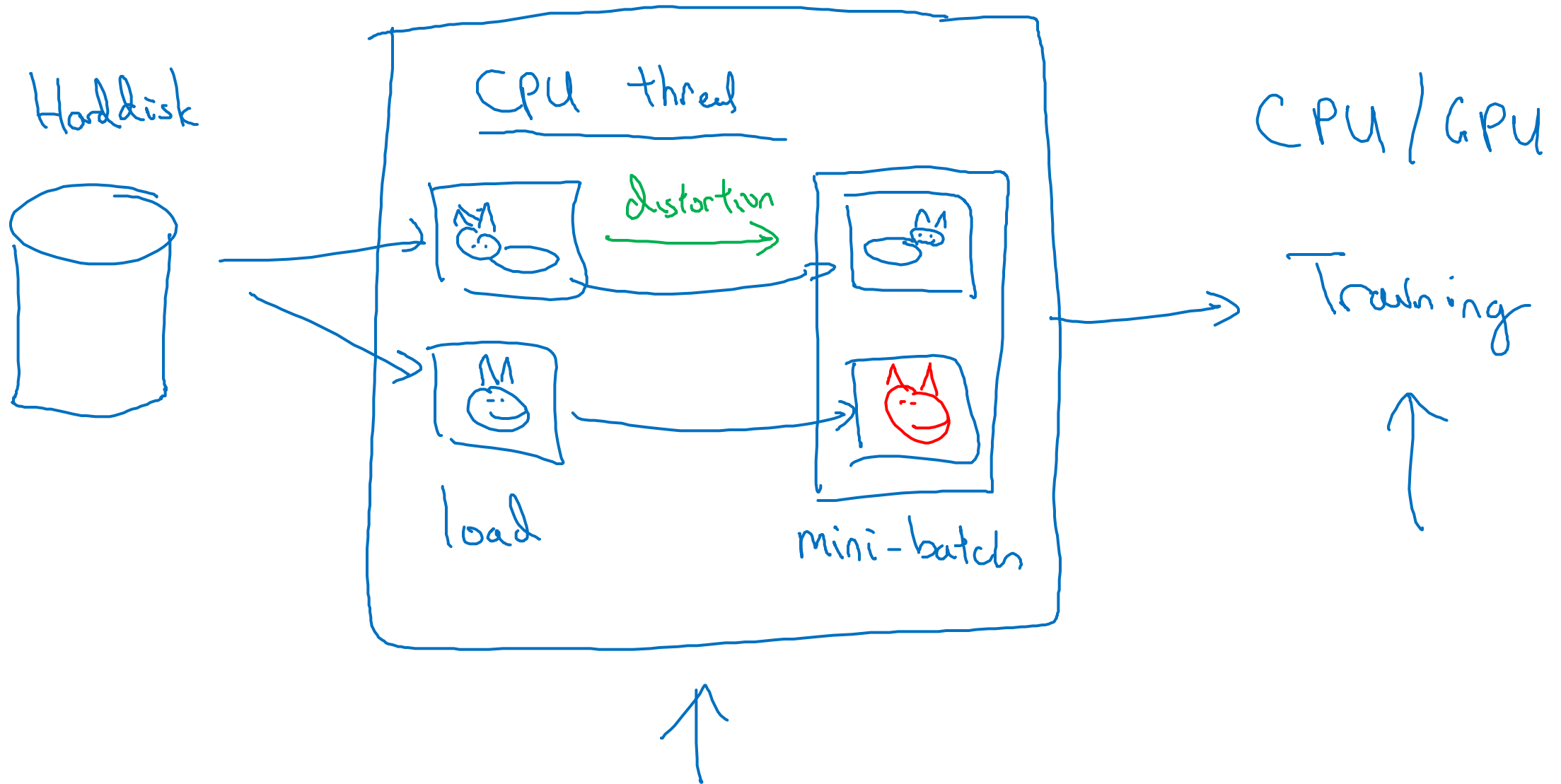[AlexNet paper
["PCA color
augmentation."

R B       G

Andrew Ng

# Implementing distortions during training



Andrew Ng

# Practical advice for using ConvNets

---

# The state of computer vision

deeplearning.ai

# Data vs. hand-engineering

Little data

Lots of data

More hand-engineering ("hacks")

Object Detection

Trigg/Msg/neith.

Image recognition

Speech recognition

Simpler algorithms less hand-engineering

Transfer learning

Two sources of knowledge

→ • Labeled data $(x, y)$

→ • Hand engineered features/network architecture/other components

Andrew Ng

# Tips for doing well on benchmarks/winning competitions

Ensembling        3 - 15   networks                    → $\hat{y}$

- Train several networks independently and average their outputs

Multi-crop at test time

- Run classifier on multiple versions of test images and average results

10-crop



1    +    4    +    1    +    4

Andrew Ng

# Use open source code

- Use architectures of networks published in the literature

- Use open source implementations if possible

- Use pretrained models and fine-tune on your dataset

Andrew Ng