# Object Classification

- Assumption: This image has a **single** object.
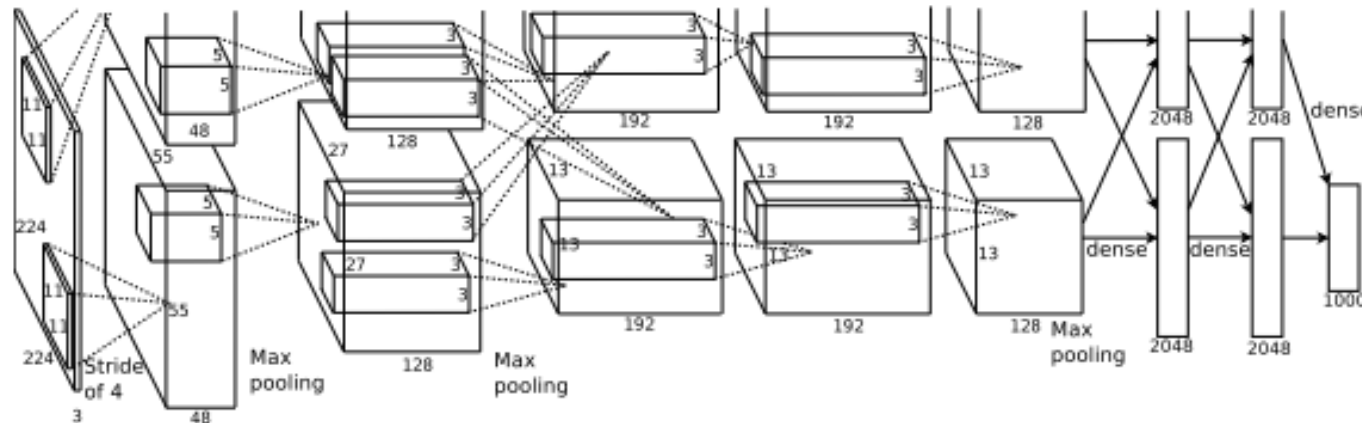
- Question: What is the object?



$$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Dog?
Fish?

Cat?

*SVM (simple) approaches were superior to deep neural network (DNN) approaches until ...*

# AlexNet [NIPS'12]

- Convolutional Neural Network (CNN) works very well on ImageNet!
  - http://www.image-net.org/challenges/LSVRC/
- DNN started receiving significant attention!
- Cited more than 86k times…

*By the way, what is CNN?*

# CNN Summary – Fully Connected Layer

- Fully connected layer (What we have seen so far)
  - Input x: nx1 vector
  - Output y: mx1 vector

x (n=36)    y (m=16)

$$y_1 = x_1 w_{1,1} + x_2 w_{2,1} + x_3 w_{3,1} + \cdots + x_{36} w_{36,1}$$

- $Y = W^T X$
- Weight: **nxm** vector



x (n=36)    x    =    y (m=16)

# CNN Summary – Convolutional Layer

- Convolutional layer (1D)
  - Input x: nx1 vector
  - Output y: mx1 vector

  <br>

  - $Y = X * W$
  - Weight: $\mathbf{f}$x1 filter (kernel)
  - Stride: hopping distance $\mathbf{s}$

Stride $\mathbf{s=2}$

Learnable filter $\mathbf{w}$ (6x1)

$y_1 = x_1\mathbf{w_1} + x_2\mathbf{w_2} + x_3\mathbf{w_3} + x_4\mathbf{w_4} + x_5w_5 + x_6\mathbf{w_6}$

*

=

x (n=36)

y (m=16)

# CNN Summary – Convolutional Layer

- Convolutional layer (1D)
  - Input x: nx1 vector
  - Output y: mx1 vector



Stride **s=2**

$y_2 = x_3 w_1 + x_4 w_2 + x_5 w_3 + x_6 w_4 + x_7 w_5 + x_8 w_6$

Learnable filter **w** (6x1)

x (n=36)

y (m=16)

- $\boldsymbol{Y = X * W}$
- Weight: **f**x1 filter (kernel)
- Stride: hopping distance **s**

# CNN Summary – Convolutional Layer

- Convolutional layer (1D)
  - Input x: nx1 vector
  - Output y: mx1 vector

  - $Y = X * W$
  - Weight: $\mathbf{f}$x1 filter (kernel)
  - Stride: hopping distance $\mathbf{s}$

Stride $\mathbf{s=2}$

$*$

Learnable filter $\mathbf{w}$ (6x1)

$=$

$y_3 = x_5\boldsymbol{w_1} + x_6\boldsymbol{w_2} + x_7\boldsymbol{w_3} + x_8\boldsymbol{w_4} + x_9\boldsymbol{w_5} + x_{10}\boldsymbol{w_6}$

x (n=36)

y (m=16)

# CNN Summary – Convolutional Layer

- Convolutional layer (1D)
  - Input x: nx1 vector
  - Output y: mx1 vector

  - $Y = X * W$
  - Weight: $f$x1 filter (kernel)
  - Stride: hopping distance **s**

Stride **s=2**

$$y_4 = x_7 w_1 + x_8 w_2 + x_9 w_3 + x_{10} w_4 + x_{11} w_5 + x_{12} w_6$$

*

=

Learnable filter **w** (6x1)

x (n=36)

y (m=16)

# CNN Summary – Convolutional Layer

- Convolutional layer (1D)
  - Input x: nx1 vector
  - Output y: mx1 vector



  - $Y = X * W$
  - Weight: **f**x1 filter (kernel)
  - Stride: hopping distance **s**

Stride **s=2**

*

Learnable filter **w** (6x1)

=

$y_{16} = x_{31}w_1 + x_{32}w_2 + x_{33}w_3 + x_{34}w_4 + x_{35}w_5 + x_{36}w_6$

x (n=36)

y (m=16)

# CNN Summary – Convolutional Layer

- Convolutional layer (2D)
  - Input x: nxn vector
  - Output y: mxm vector

$$y_{11} = x_{11}w_{11} + x_{21}w_{21} + x_{31}w_{31}$$
$$+ x_{12}w_{21} + x_{22}w_{22} + x_{32}w_{32}$$
$$+ x_{13}w_{13} + x_{23}w_{23} + x_{33}w_{33}$$

*

Learnable filter **w** (3x3)

=

x (6x6)

y (4x4)

- $Y = X * W$
- Weight: **f**x**f** filter (kernel)
- Stride: hopping distance **s**

# CNN Summary – Convolutional Layer

- Convolutional layer (2D)
  - Input x: nxn vector
  - Output y: mxm vector

Stride **s=1**

$$y_{12} = x_{12}w_{11} + x_{22}w_{21} + x_{32}w_{31}$$
$$+x_{13}w_{21} + x_{23}w_{22} + x_{33}w_{32}$$
$$+x_{14}w_{13} + x_{24}w_{23} + x_{34}w_{33}$$

Learnable filter **w** (3x3)

x (6x6)

*

=

y (4x4)

- $Y = X * W$
- Weight: **f**x**f** filter (kernel)
- Stride: hopping distance **s**

# CNN Summary – Convolutional Layer

- Convolutional layer (2D)
  - Intuition how it works



Edge detector

Filters can be used to extract specific features

# CNN Summary – Convolutional Layer

- Convolutional layer (2D)
  - Input x: nxn vector
  - Output y: mxm vector



Stride **s=1**

Padding **p=1**

Learnable filter **w** (3x3)

x (6x6)   x (8x8)

y (4x4)

- $Y = X * W$
- Weight: **f**x**f** filter (kernel)
- Stride: hopping distance **s**
- Padding (**p**): adding extra values around the edge
  - Effect: (1) Avoid shrinking image too fast, (2) Avoid throwing away info from the edge
  - Type: (1) Valid Conv (no padding), (2) Same Conv (padding to maintain output size)

# CNN Summary – Convolutional Layer

- Convolutional layer (2D)
  - Input x: nxn vector
  - Output y: mxm vector



Stride **s=1**

Padding **p=0**

Learnable filter **w** (3x3)

x (6x6)

y (4x4)

$$\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \text{x} \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

- $\boldsymbol{Y = X * W}$
- Weight: **f**x**f** filter (kernel)
- Stride: hopping distance **s**
- Padding (**p**): adding extra values around the edge
  - Effect: (1) Avoid shrinking image too fast, (2) Avoid throwing away info from the edge
  - Type: (1) Valid Conv (no padding), (2) Same Conv (padding to maintain output size)

# CNN Summary – Pooling Layer

- Pooling layer
  - Filter (fxf)
    - Max or average pooling
    - Pooling layer has **nothing learnable**
  - Stride (s)



Stride **s=2**

*Then… why is CNN better than MLP for computer vision?*

# Why CNN?

- CNN



6x6       3x3 =
**9 parameters**       4x4

*Parameter sharing (less memory)*

- Fully connected layer



36x16 = **576 parameters**

36       16

# Why CNN?

- ## CNN

**9 connections**
/output



6x6

3x3 =
**9 parameters**

4x4

- ## Fully connected layer

**36 connections**
/output



36x16 = **576 parameters**

36

16

*Parameter sharing (less memory)*

*Sparsity of connections (less computation)*

*Regularization (less overfitting)*

*Now, let's go back to AlexNet.*

# AlexNet [NIPS'12] - Architecture



227x227x3

[Conv11, s=4]
x96

55x55x96

MaxPool
3x3, s=2

27x27x96

[Conv5,same]
x256

27x27x256

MaxPool
3x3, s=2

13x13x256

x2 layers

[Conv3, same]
x384

13x13x384

[Conv3, same]
x256

13x13x256

MaxPool
3x3, s=2

6x6x256 =

9216

FC
4096

FC
4096

Softmax
1000

# AlexNet [NIPS'12] – Pros and Cons

- Similar to LeNet (or LeNet-5) but has **1000x more** parameters (~60M)
  - More parameters can be trained well due to more data and computation power
- ReLu activation instead of sigmoid

- <span style="color:red">Complex architecture</span>
  - Various filter sizes: 3x3, 5x5, 11x11
  - Various strides: 1, 2, 4

# VGG-16 [ICLR'15]

- Simple and uniform architecture!
  - One type convolution layer: 3x3, s=1
  - One type pooling layer: max pooling, 2x2, s=2
  - Doubling channels

# VGG-16 [ICLR'15] - Architecture



224x224x3

x2 layers
[Conv3, same] x64 → 224x224x64 → MaxPool 2x2, s=2 → 112x112x64

x2 layers
[Conv3, same] x128 → 112x112x128 → MaxPool 2x2, s=2 → 56x56x128

x3 layers
[Conv3, same] x256 → 56x56x256 → MaxPool 2x2, s=2 → 28x28x256

x3 layers
[Conv3, same] x512 → 28x28x512 → MaxPool 2x2, s=2 → 14x14x512

# VGG-16 [ICLR'15] - Architecture

x3 layers

[Conv3, same] x512 → 14x14x512 → MaxPool, 2x2, s=2 → 7x7x512

FC 4096 → FC 4096 → FC 1000

SoftMax 1000

# VGG-16 [ICLR'15]

- Simple and uniform architecture!
  - One type convolution layer: 3x3, s=1
  - One type pooling layer: max pooling, 2x2, s=2
  - Doubling channels


- So many parameters (16 layers, ~138M parameters)

# Inception [CVPR'15]

- Problems
  - Computation increases with # of layers
  - CNN design is complex
    - You should determine layer type (conv or pooling), filter size, and stride…

- Don't worry, we will do these all at each layer, with **efficient computation!**
  - Then, we can go deeper conveniently

# 1x1 Convolution

- Cross-channel fully connected layer



28x28x192         [1x1x192 filter] x 192         28x28x192

Filter 1

ReLu

192 inputs     192 parameters (filter 1)     192 outputs

# 1x1 Convolution

- Cross-channel fully connected layer

Filter 1

Filter 2

**ReLu**

\*

=

28x28x192

[1x1x192 filter] x 192

28x28x192

192 inputs

192 parameters
(filter 2)

192 outputs

# 1x1 Convolution

- What if we change the number of filters?



28x28x192        [1x1x192 filter] x **64**        28x28x**64**

*Changing # of channels while maintaining width and height*

# 1x1 Convolution

- Finally, 1x1 conv is actually **faster** than normal conv even when the number of operations is same, because it does not require **memory reordering**



IF YOU DON'T UNDERSTAND, DON'T WORRY ABOUT IT

makeameme.org

# Inception [CVPR'15]

- Inception module – naïve (impractical) version
  - Mix various filters in a single layer and let them trained automatically
  - Maintain width and height by using padding

# Inception [CVPR'15]

- Inception module – Reducing computation cost
  - Passing a conv filter is quite expensive when the input (or output) has **many channels**



28x28x192 → [Conv5x192, same] x32 → 28x28x32

\# of multiplications
(5x5x**192**)x(28x28)x32 = **120.4M**

# Inception [CVPR'15]

- Inception module – Reducing computation cost
  - Passing a conv filter is quite expensive when the input (or output) has **many channels**
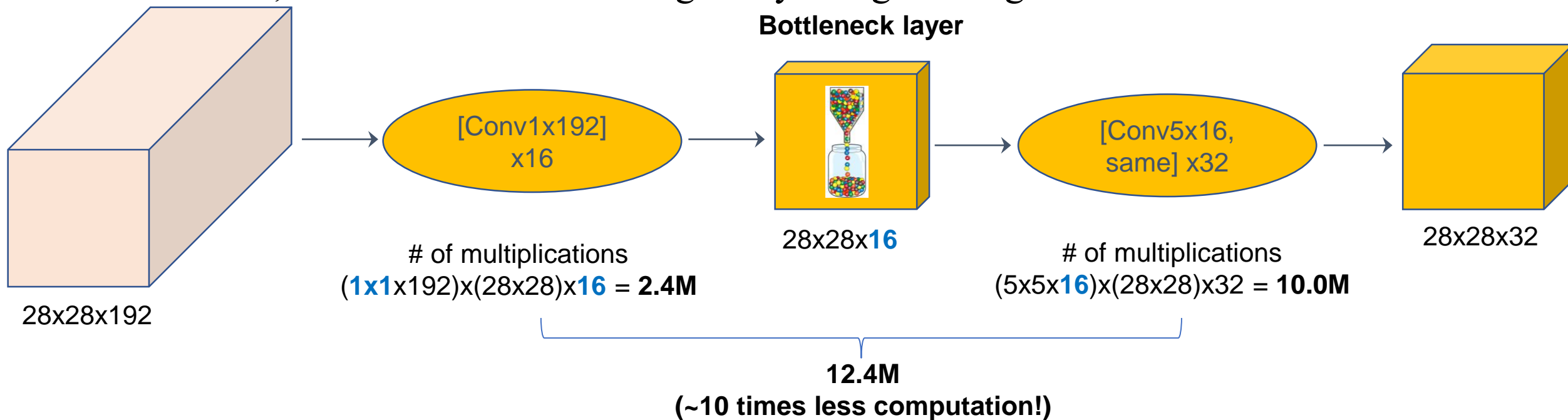  - Reduce # of channels by **inserting a 1x1 Conv layer**, making an intermediate **bottleneck** layer
  - Then, increase # of channels again by using the original filter size

**Bottleneck layer**

[Conv1x192] x16

28x28x**16**

[Conv5x16, same] x32

28x28x192

28x28x32

# of multiplications
(**1x1**x192)x(28x28)x**16** = **2.4M**

# of multiplications
(5x5x**16**)x(28x28)x32 = **10.0M**

**12.4M**
**(~10 times less computation!)**

# Inception [CVPR'15]

- Inception module – with less computational overhead

# Inception [CVPR'15]

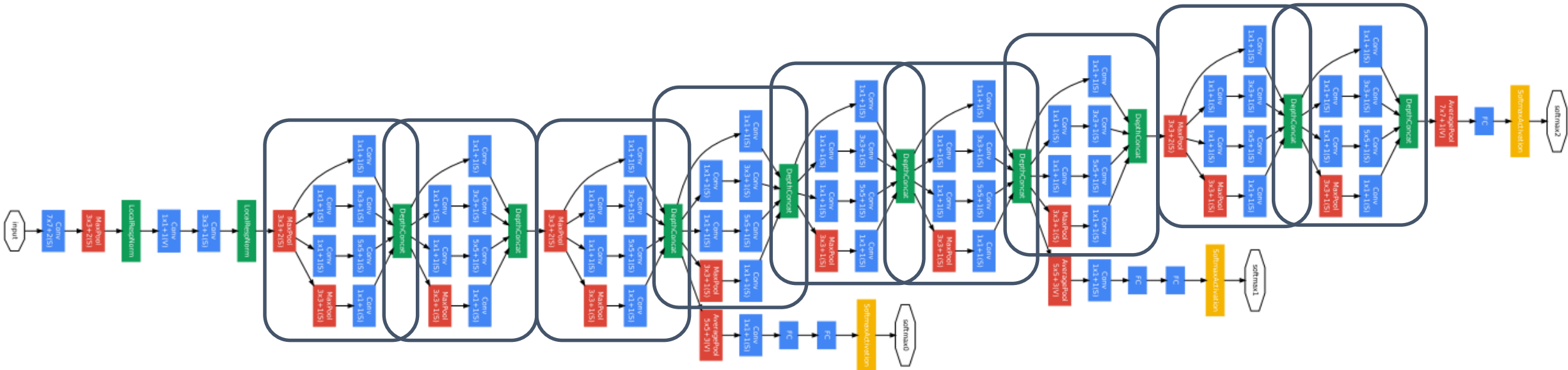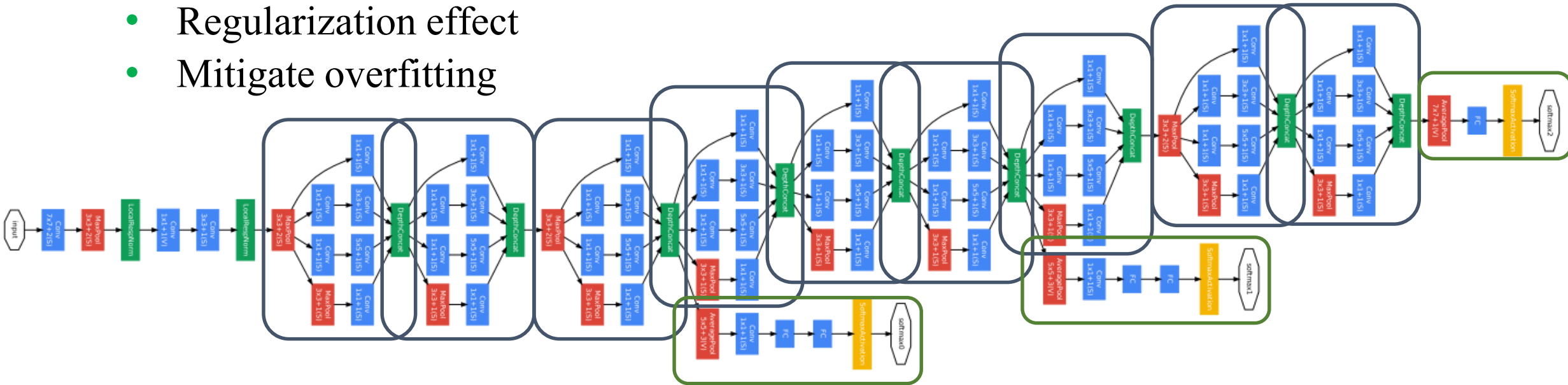- GoogLeNet – a **22-layer** DNN using inception modules

# Inception [CVPR'15]

- GoogLeNet – a **22-layer** DNN using inception modules
- Extract results in the middle of network and include them in the loss function
  - Regularization effect
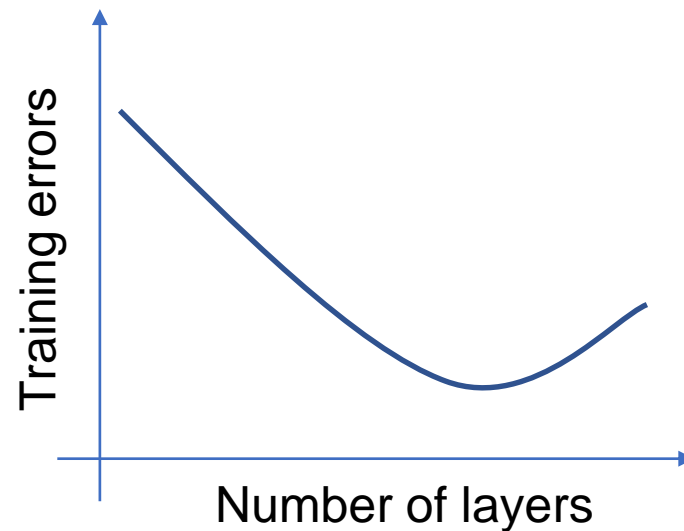  - Mitigate overfitting

# Inception [CVPR'15]

- Inception shows that, even though inserting a bottleneck layer could lose lots of information, doing it is actually **NOT** that harmful


- Implication
  - CNN is more efficient than MLP
    - Parameter sharing and sparsity of connections
  - However, this architecture is also designed by **human intuition**
    - Yes, it works, but there has been no mathematical/empirical verification showing if it is "THE" optimal architecture
  - A significant amount of CNN computation might still be **unnecessary**
  - We might be able to build a lighter DNN without sacrificing accuracy

# ResNet [CVPR'16] - Motivation

- You know what? I want to train an even ***deeper and heavier*** neural network

- Hmm… but when a model is too deep, performance is degraded



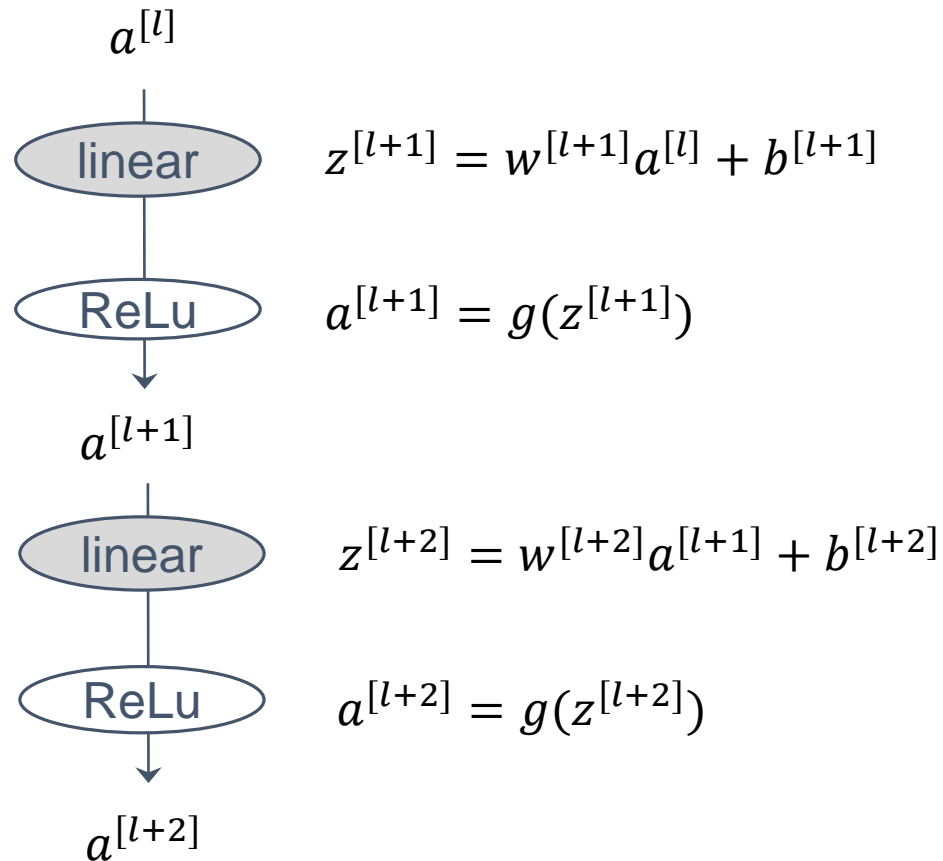*How can we continuously decrease training errors as the number of layers increases?*

# ResNet [CVPR'16] – Core Element

- Residual connection

$$a^{[l]}$$

$$\text{linear} \qquad z^{[l+1]} = w^{[l+1]}a^{[l]} + b^{[l+1]}$$

$$\text{ReLu} \qquad a^{[l+1]} = g(z^{[l+1]})$$

$$a^{[l+1]}$$

$$\text{linear} \qquad z^{[l+2]} = w^{[l+2]}a^{[l+1]} + b^{[l+2]}$$

$$\text{ReLu} \qquad a^{[l+2]} = g(z^{[l+2]})$$

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

# ResNet [CVPR'16] – Core Element

- Residual connection



$a^{[l]}$

linear $\qquad z^{[l+1]} = w^{[l+1]}a^{[l]} + b^{[l+1]}$

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

**Short cut!**

$a^{[l+1]}$

linear $\qquad z^{[l+2]} = w^{[l+2]}a^{[l+1]} + b^{[l+2]}$

$\oplus$

ReLu $\qquad a^{[l+2]} = g(z^{[l+2]} + \boldsymbol{a^{[l]}})$

$a^{[l+2]}$

# ResNet [CVPR'16]

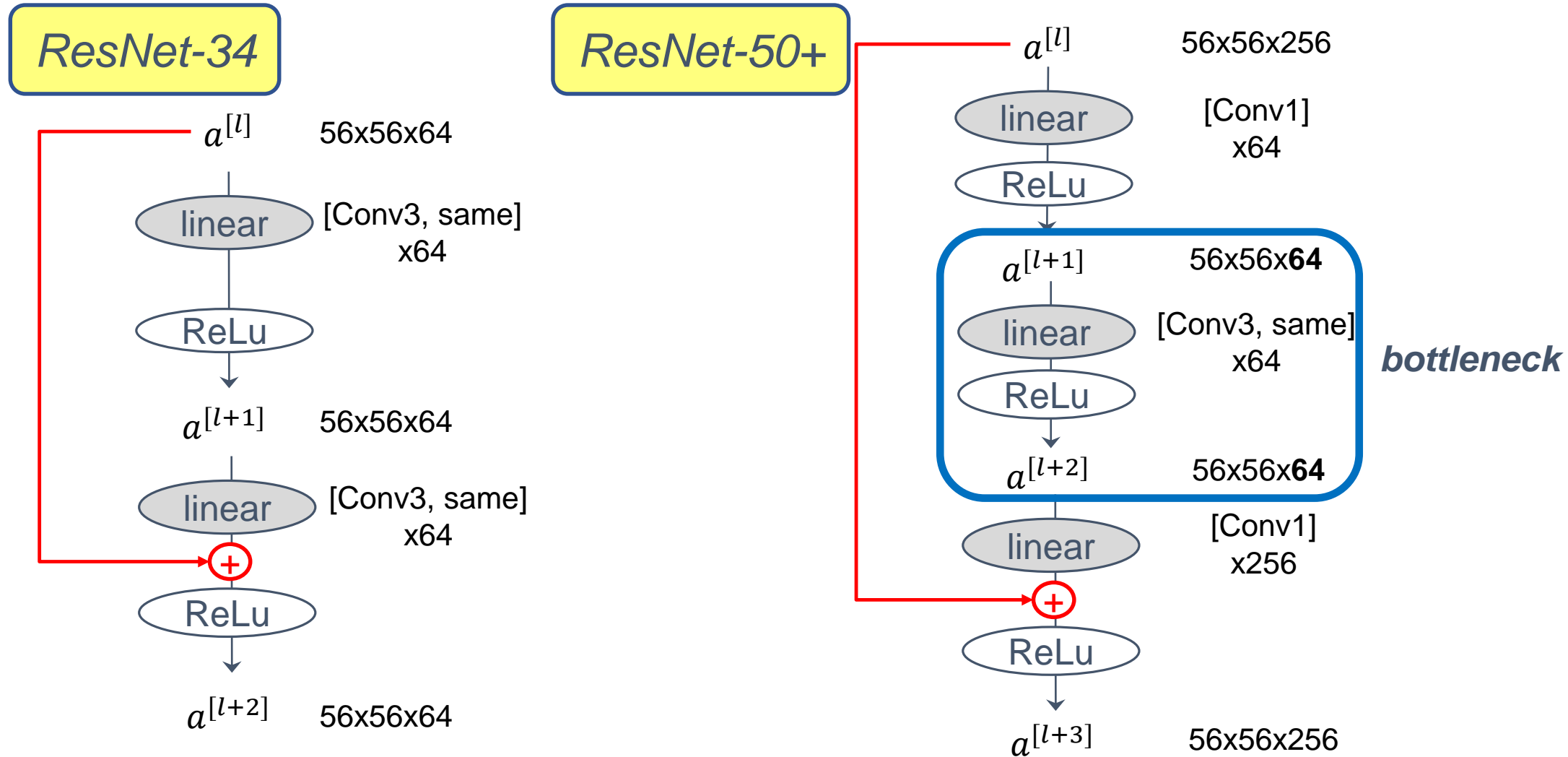- ResNet: DNN architecture with residual connection

# ResNet [CVPR'16] – Why does it work?

- Why is the residual connection good for training a deep network?
  - In a negative case, weights of a new layer can be trained as the <span style="color:red">identity function</span>, which does not degrade performance at least
    - $a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$
      $= g(w^{[l+2]}a^{[l+1]} + b^{[l+2]} + a^{[l]})$
      $= g(a^{[l]})$    when $w^{[l+2]} = 0, b^{[l+2]} = 0$
      $= a^{[l]}$

  - In a fortunate case, weights of a new layer can be trained in a positive way, improving performance
  - Less risk, but more opportunity as a DNN goes deeper **(>100 layers)**

*But how about training overhead for a super deep DNN?*

*We have the **bottleneck** approach!*

# ResNet [CVPR'16] – Bottleneck Architecture

Thanks!