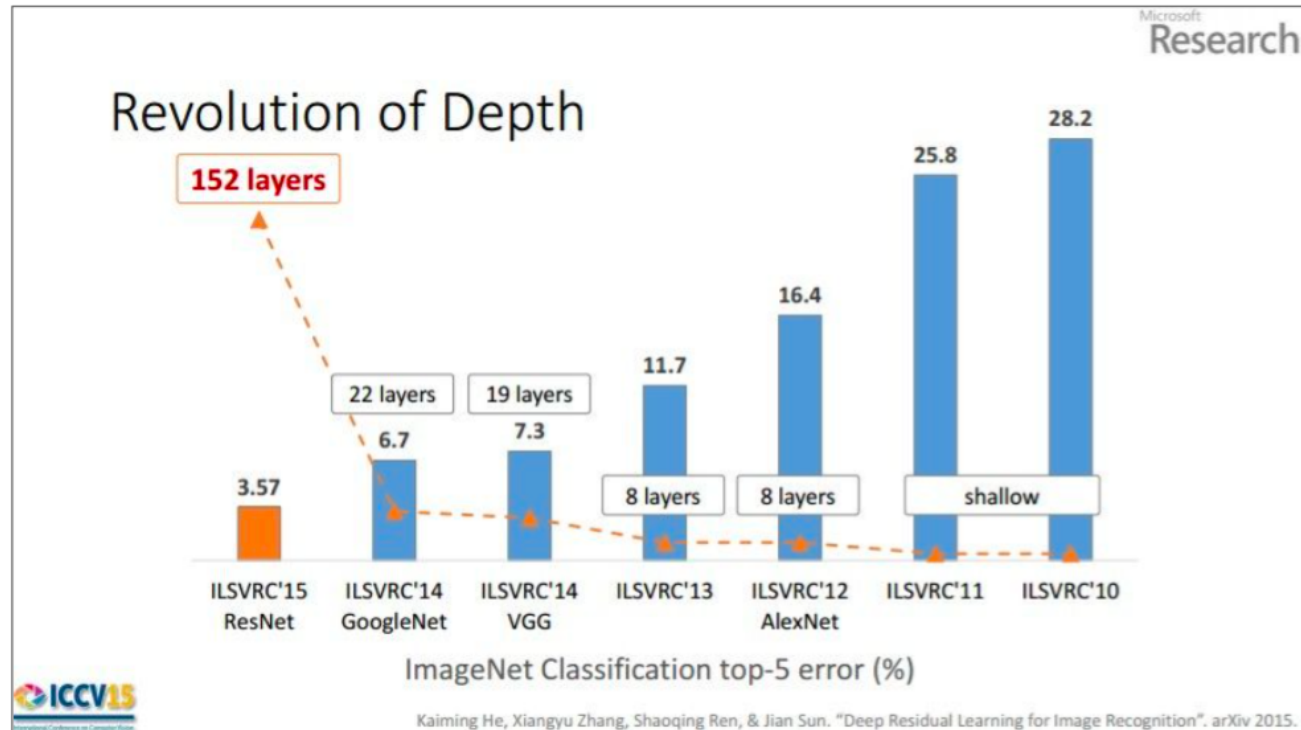# TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2020

## Convolutional Neural Networks (CNNs)

# Imagenet Classification

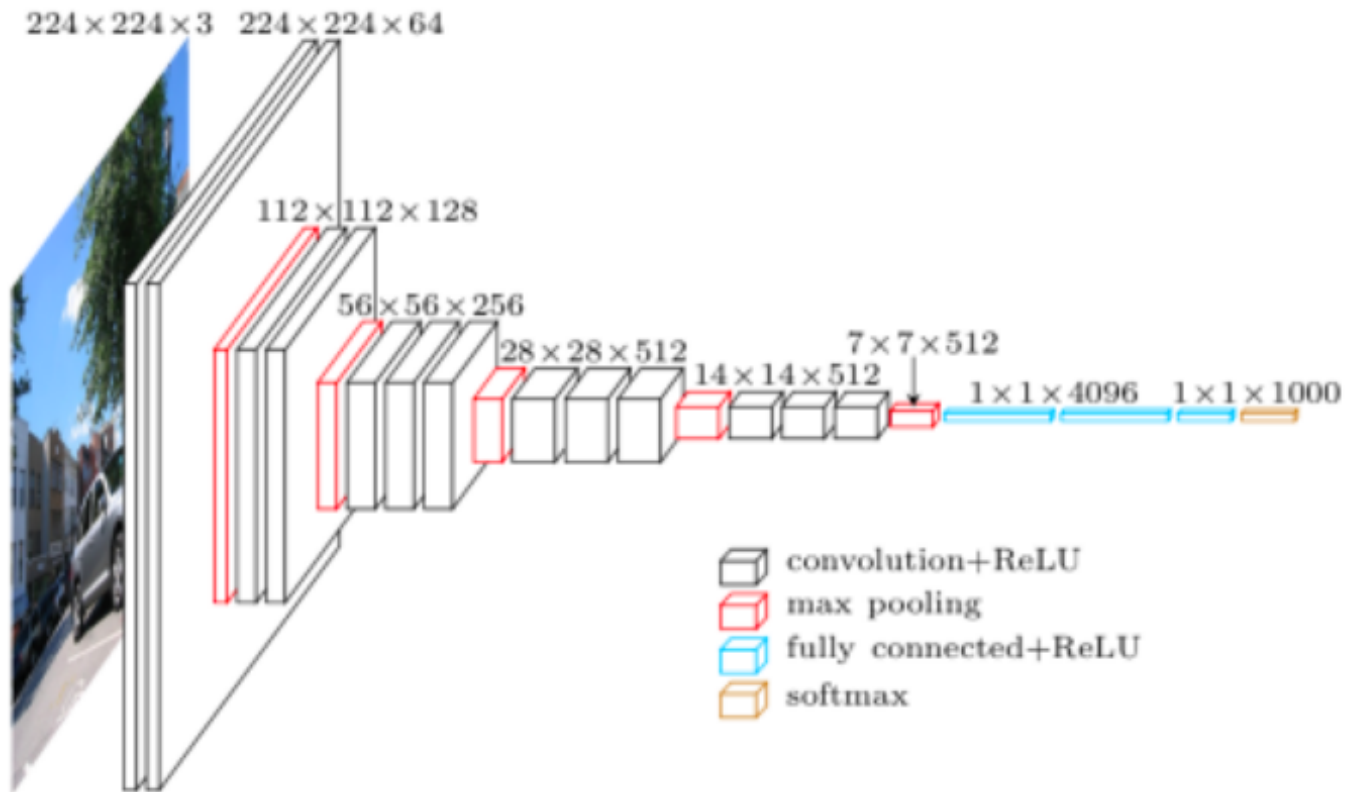1000 kinds of objects.



(slide from Kaiming He's recent presentation)
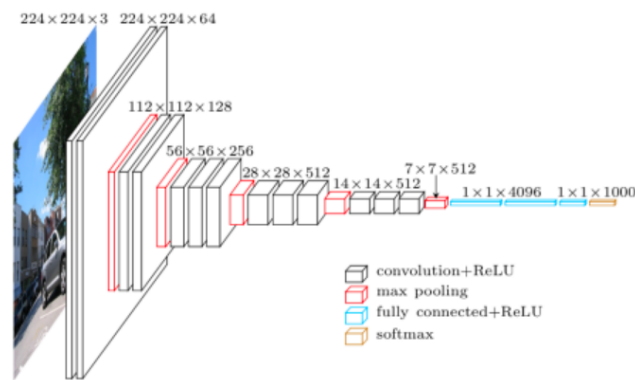
2016 is 3.0%, is 2017 2.25%

SOTA as of January 2020 is 1.3%

# What is a CNN?
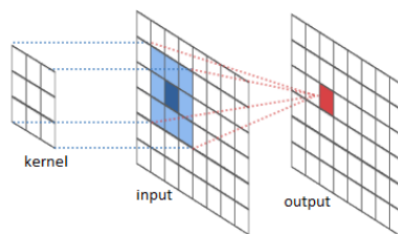## VGG, Zisserman, 2014



Davi Frossard

# A Convolution Layer



Each box is a tensor $L_\ell[b, x, y, i]$

For a convolution layer, each $L_{\ell+1}[b, x, y, j]$ is the output of a single linear threshold unit computed from $L_\ell[b, x, y, i]$.

# A Convolution Layer



kernel

input

output

$$W[\Delta x, \Delta y, i, j] \qquad L_\ell[b, x, y, i] \qquad L_{\ell+1}[b, x, y, j]$$

River Trail Documentation

$$L_{\ell+1}[b, x, y, j]$$

$$= \sigma \left( \left( \sum_{\Delta x, \Delta y, i} W[\Delta x, \Delta y, i, j] \, L_\ell[b, x + \Delta x, y + \Delta y, i] \right) - B[j] \right)$$

5

# Many "Neurons" (Linear Threshold Units)

Each $L_{\ell+1}[b, x, y, j]$ is the output of a single linear threshold unit.

$$L_{\ell+1}[b, x, y, j]$$

$$= \sigma \left( \left( \sum_{\Delta x, \Delta y, i} W[\Delta x, \Delta y, i, j] \, L_{\ell}[b, x + \Delta x, y + \Delta y, i] \right) - B[j] \right)$$

# 2D CNN in PyTorch

conv2d(**input, weight, bias, stride, padding, dilation, groups**)

**input** tensor (minibatch,in-channels,iH,iW)

**weight** filters (out-channels, in-channels/groups, in-channels,kH,kW)

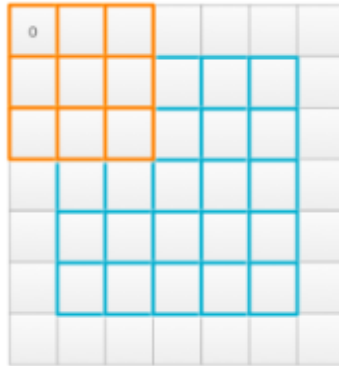**bias** tensor (out-channels) . Default: None

**stride** Single number or (sH, sW). Default: 1

**padding** Single number or (padH, padW). Default: 0

**dilation** Single number or (dH, dW). Default: 1

**groups** split input into groups. Default: 1

# Padding



Jonathan Hui

If we pad the input with zeros then the input and output can have the same spatial dimensions.

# Zero Padding in NumPy

In NumPy we can add a zero padding of width p to an image as follows:

```
padded = np.zeros(W + 2*p,  H + 2*p)

padded[p:W+p, p:H+p] = x
```
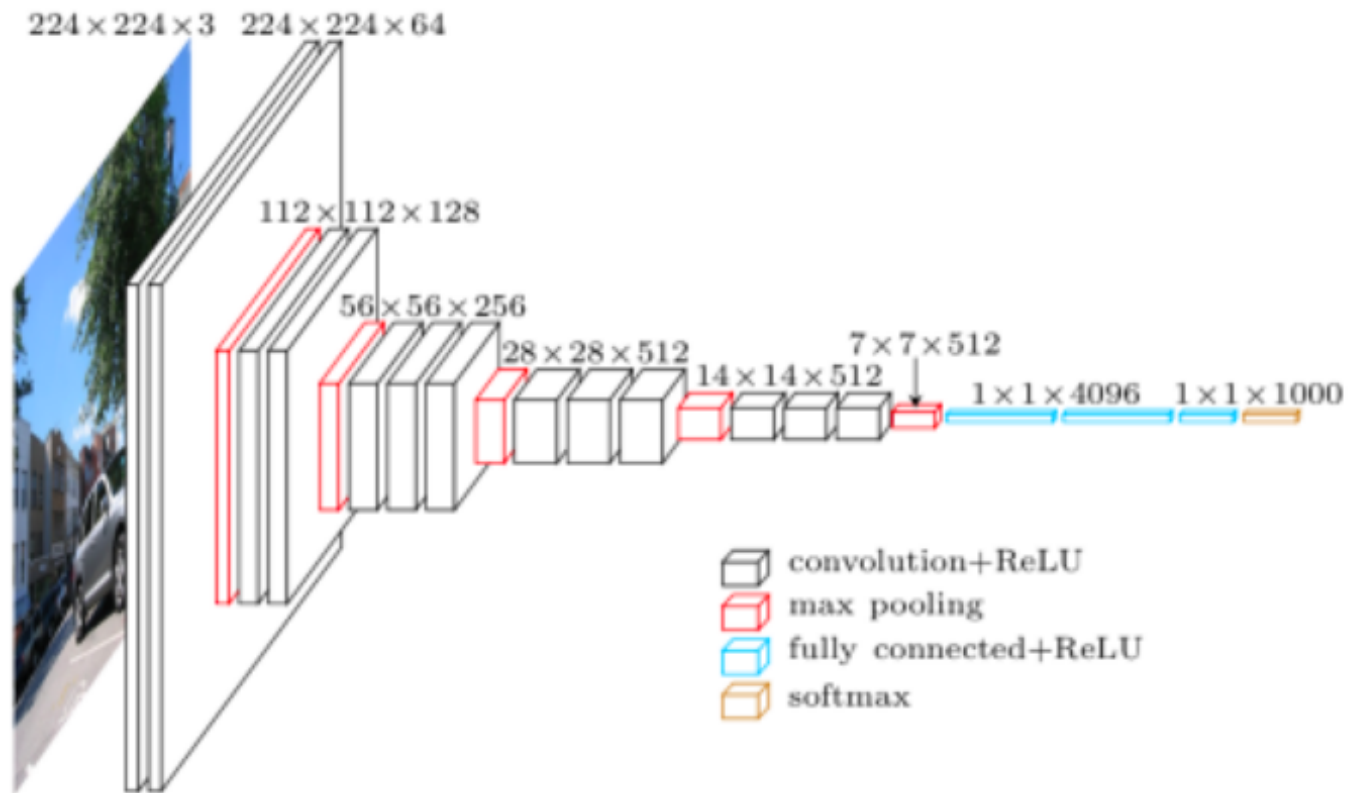
# Padding

$$L'_\ell = \text{Padd}(L_\ell, \ p)$$

$$L_{\ell+1}[b, x, y, j] =$$

$$\sigma\left(\left(\textstyle\sum_{\Delta x, \Delta y, i} W[\Delta x, \Delta y, i, j]\ L'_\ell[b, x + \Delta x, y + \Delta y, i]\right) - B[j]\right)$$

If the input is padded but the output is not padded then $\Delta x$ and $\Delta y$ are non-negative.

# Reducing Spatial Dimention

# Strides

We can move the filter by a "stride" $s$ for each spatial step.

$$L_{\ell+1}[b, \textcolor{red}{x}, \textcolor{red}{y}, j] =$$

$$\sigma\left(\left(\sum_{\Delta x, \Delta y, i} W[\Delta x, \Delta y, i, j] L_\ell[b, \textcolor{red}{s * x} + \Delta x, \textcolor{red}{s * y} + \Delta y, i]\right) - B[j]\right)$$
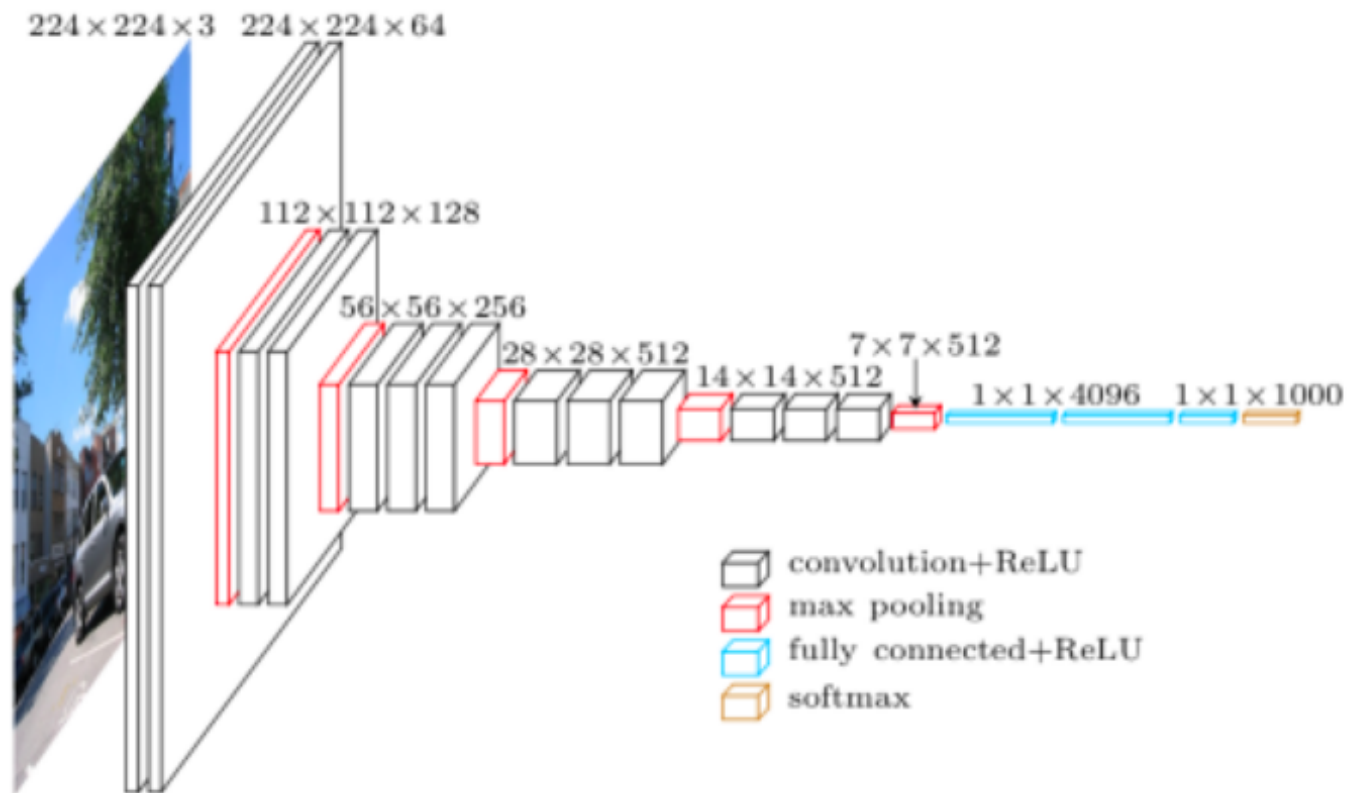
For strides greater than 1 the spatial dimention is reduced.

# Max Pooling

$$L_{\ell+1}[b, \textcolor{red}{x}, \textcolor{red}{y}, i] = \max_{\textcolor{red}{\Delta x, \Delta y}} L_\ell[b, \textcolor{red}{s * x} + \Delta x, \; \textcolor{red}{s * y} + \Delta y, \; i]$$

This is typically done with a stride greater than one so that the image dimension is reduced.

# Fully Connected (FC) Layers

# Fully Connected (FC) Layers

We reshape $L_\ell[b, x, y, i]$ to $L_\ell[b, i']$ and then

$$L_{\ell+1}[b, j] = \sigma \left( \left( \sum_{i'} W[j, i'] \, L_\ell[b, i'] \right) - B[j] \right)$$

# Image to Column (Im2C)

Reduce convolution to matrix multiplication

more space but faster.

$$\tilde{L}_{\ell+1}[b, x, y, j]$$

$$= \left( \sum_{\Delta x, \Delta y, i} W[\Delta x, \Delta y, i, j] * L_{\ell}[b, x + \Delta x, \ y + \Delta y, \ i] \right) + B[j]$$

We make a bigger tensor $\tilde{L}$ with two additional indeces.

$$\tilde{L}_{\ell}[b, x, y, \Delta x, \Delta y, i] = L_{\ell}[b, x + \Delta x, y + \Delta y, i]$$

# Image to Column (Im2C)

$$\tilde{L}_{\ell+1}[b, x, y, j]$$

$$= \left( \sum_{\Delta x, \Delta y, i} W[\Delta x, \Delta y, i, j] * L_\ell[b, x + \Delta x, \ y + \Delta y, \ i] \right) + B[j]$$
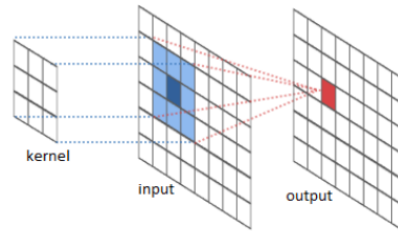
$$= \left( \sum_{\Delta x, \Delta y, i} \tilde{L}_\ell[b, x, y, \Delta x, \Delta y, i] * W[\Delta x, \Delta y, i, j] \right) + B[j]$$

$$= \left( \sum_{(\Delta x, \Delta y, i)} \tilde{L}_\ell[(b, x, y), (\Delta x, \Delta y, i)] * W[(\Delta x, \Delta y, i), j] \right) + B[j]$$
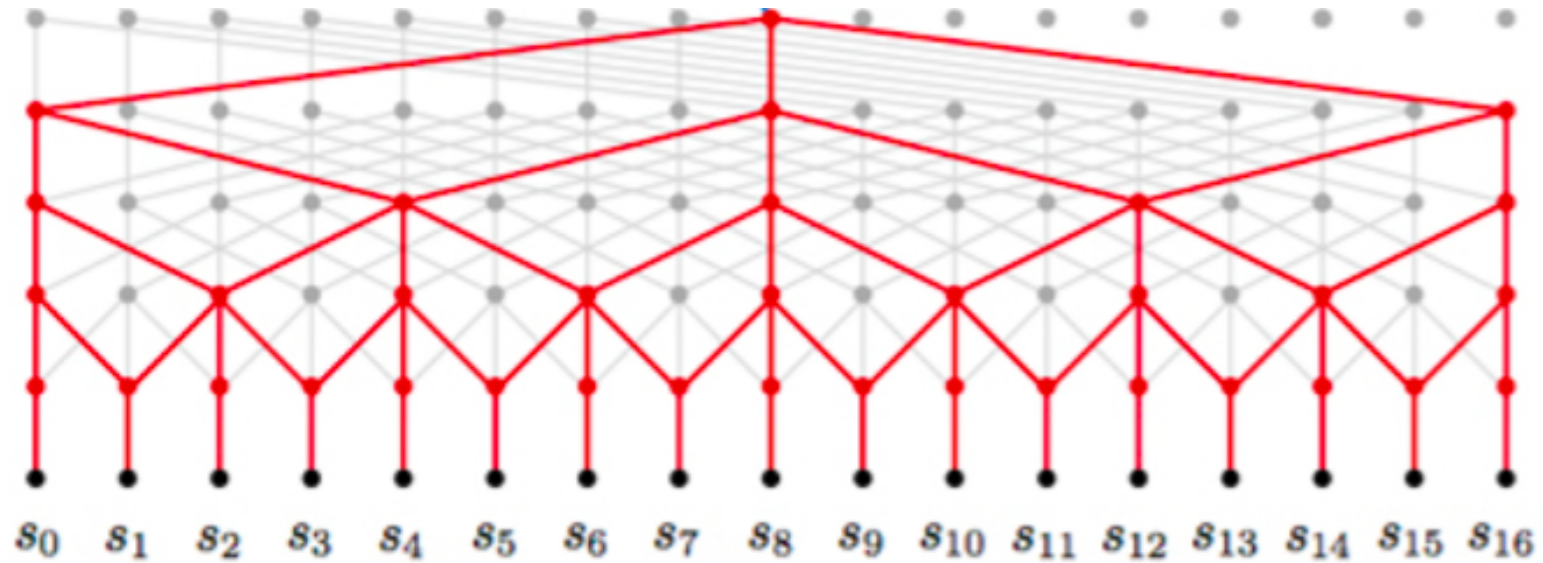
# Dilation

A CNN for image classification typically reduces an $N \times N$ image to a single feature vector.

Dilation is a trick for treating the whole CNN as a "filter" that can be passed over an $M \times M$ image with $M > N$.



An output tensor with full spatial dimension can be useful in, for example, image segmentation.

# Dilation



$s_0 \quad s_1 \quad s_2 \quad s_3 \quad s_4 \quad s_5 \quad s_6 \quad s_7 \quad s_8 \quad s_9 \quad s_{10} \quad s_{11} \quad s_{12} \quad s_{13} \quad s_{14} \quad s_{15} \quad s_{16}$

This is called a "fully convolutional" CNN.

# Dilation

To implement a fully convolutional CNN we can "dilate" the filters by a dilation parameter $d$.

$$\tilde{L}_{\ell+1}[b, x, y, j] = W[\Delta x, \Delta y, i, j] L_{\ell}[b, x + \textcolor{red}{d * \Delta x}, y + \textcolor{red}{d * \Delta y}, i] + B[j]$$

# Hypercolumns

An alternative to dilation and fully convolutional networks is hypercolumns.

$$L_{\ell+1}[b, x, y] = L_1[b, x, y]; \cdots ; L_n[b, x/W_\ell, y/W_\ell]; \cdots ; L_N[b]$$

where x;y denotes the concatenation of vectors $x$ and $y$.

# Grouping

$$L_{\ell+1}[b, x, y] = L^1_{\ell+1}[b, x, y]; \cdots ; L^G_{\ell+1}[b, x, y]$$

$$\tilde{L}^g_{\ell+1}[b, x, y, j] = \left( \sum_{\Delta x, \Delta y, i} W^g[\Delta x, \Delta y, i, j] * L_\ell[b, x + \Delta x, \ y + \Delta y, \ i] \right)$$

For a fixed number of features $j$ in the total output $L_{\ell+1}[b, x, y, j]$ using $G$ groups reduces the number of weight parameters by a factor of $G$.

# 2D CNN in PyTorch

conv2d(**input, weight, bias, stride, padding, dilation, groups**)

**input**  tensor (minibatch,in-channels,iH,iW)

**weight** filters (out-channels, in-channels/groups, in-channels,kH,kW)

**bias**  tensor (out-channels) . Default: None

**stride**  Single number or (sH, sW). Default: 1

**padding**  Single number or (padH, padW). Default: 0

**dilation**  Single number or (dH, dW). Default: 1

**groups**  split input into groups. Default: 1

# Modern Trends

Modern Convolutions use 3X3 filters. This is faster and has fewer parameters. Expressive power is preserved by increasing depth with many stride 1 layers.

Max pooling and dilation seem to have disappeared.

Resnet and resnet-like architectures are now dominant (next lecture).

# Alexnet

Given Input$[227, 227, 3]$

$$L_1[55 \times 55 \times 96] = \text{ReLU}(\text{CONV}(\text{Input}, \Phi_1, \text{width } 11, \text{pad } 0, \text{stride } 4))$$
$$L_2[27 \times 27 \times 96] = \text{MaxPool}(L_1, \text{width } 3, \text{stride } 2))$$
$$L_3[27 \times 27 \times 256] = \text{ReLU}(\text{CONV}(L_2, \Phi_3, \text{width } 5, \text{pad } 2, \text{stride } 1))$$
$$L_4[13 \times 13 \times 256] = \text{MaxPool}(L_3, \text{width } 3, \text{stride } 2))$$
$$L_5[13 \times 13 \times 384] = \text{ReLU}(\text{CONV}(L_4, \Phi_5, \text{width } 3, \text{pad } 1, \text{stride } 1))$$
$$L_6[13 \times 13 \times 384] = \text{ReLU}(\text{CONV}(L_5, \Phi_6, \text{width } 3, \text{pad } 1, \text{stride } 1))$$
$$L_7[13 \times 13 \times 256] = \text{ReLU}(\text{CONV}(L_6, \Phi_7, \text{width } 3, \text{pad } 1, \text{stride } 1))$$
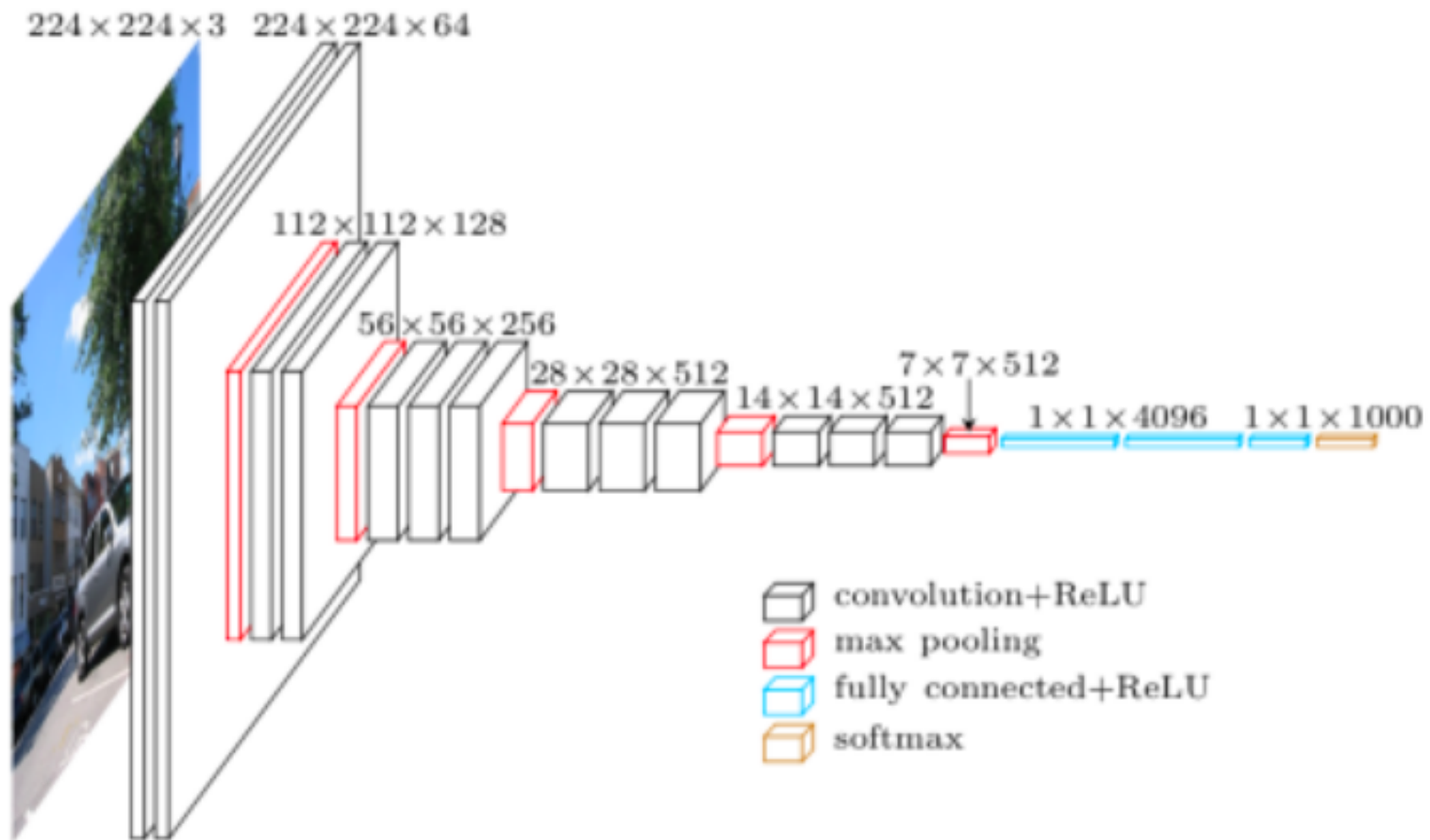$$L_8[6 \times 6 \times 256] = \text{MaxPool}(L_7, \text{width } 3, \text{stride } 2))$$
$$L_9[4096] = \text{ReLU}(\text{FC}(L_8, \Phi_9))$$
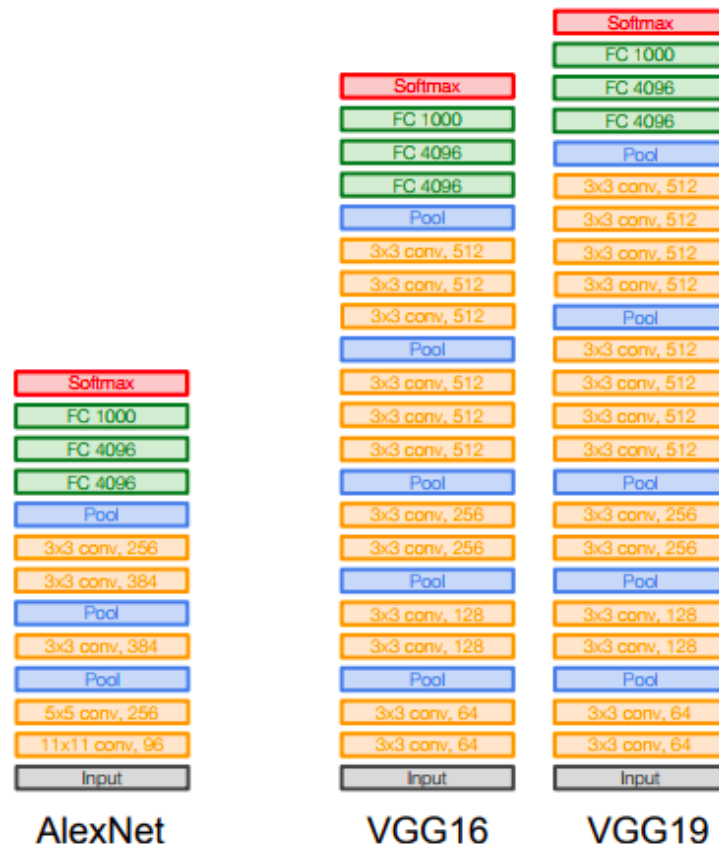$$L_{10}[4096] = \text{ReLU}(\text{FC}(L_9, \Phi_{10}))$$
$$s[1000] = \text{ReLU}(\text{FC}(L_{10}, \Phi_s)) \quad \text{class scores}$$
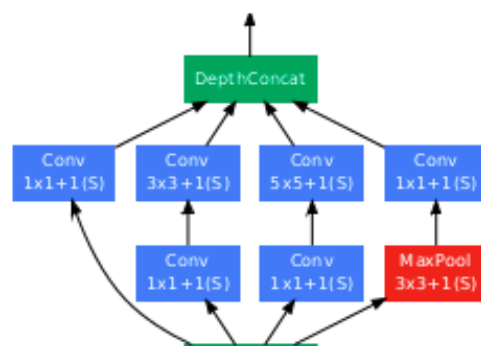
# VGG, Zisserman, 2014



$224 \times 224 \times 3$   $224 \times 224 \times 64$

$112 \times 112 \times 128$

$56 \times 56 \times 256$

$28 \times 28 \times 512$

$14 \times 14 \times 512$

$7 \times 7 \times 512$

$1 \times 1 \times 4096$   $1 \times 1 \times 1000$

convolution+ReLU
max pooling
fully connected+ReLU
softmax

Davi Frossard

26

# VGG



AlexNet | VGG16 | VGG19
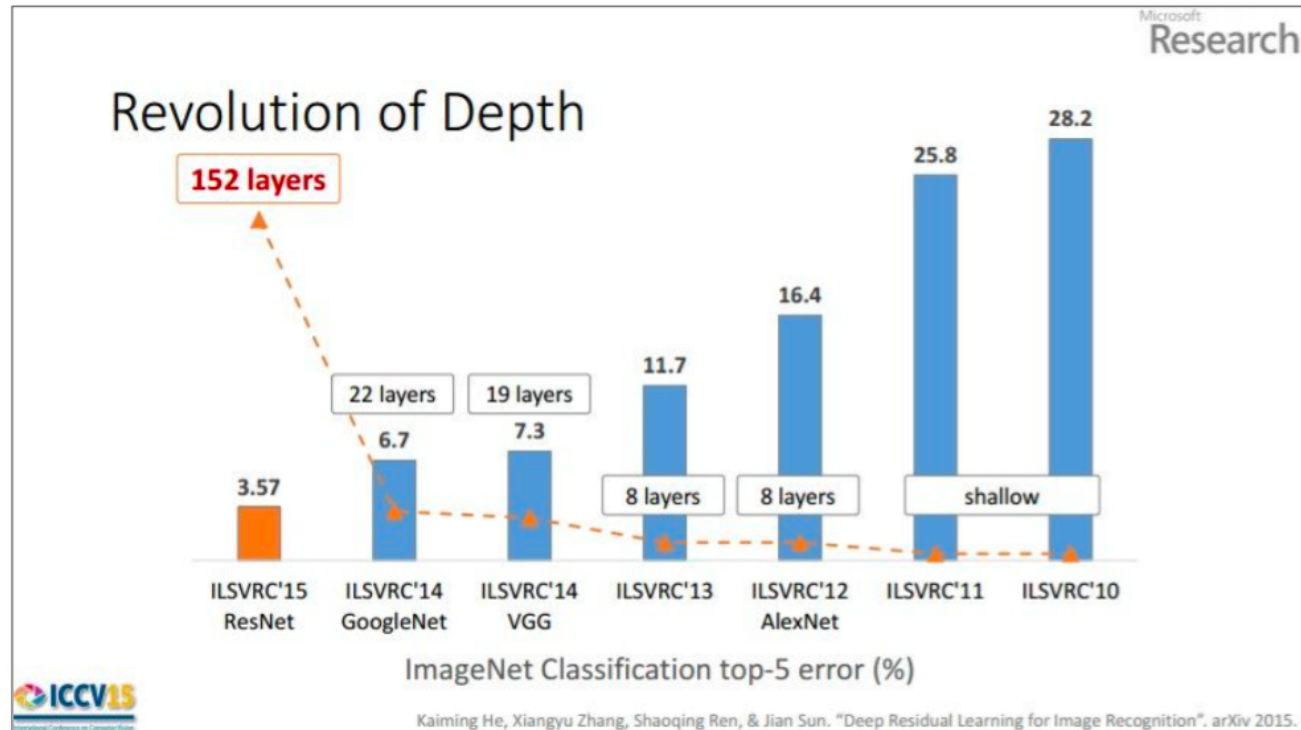
# Inception, Google, 2014

# Models for Image Classification in PyTorch

- AlexNet

- VGG

- ResNet

- SqueezeNet

- DenseNet

- Inception v3

- GoogLeNet

- ShuffleNet v2

- MobileNet v2

- ResNeXt

- Wide ResNet

- MNASNet

# Imagenet Classification

1000 kinds of objects.



(slide from Kaiming He's recent presentation)

2016 is 3.0%, is 2017 2.25%

SOTA as of January 2020 is 1.3%

END