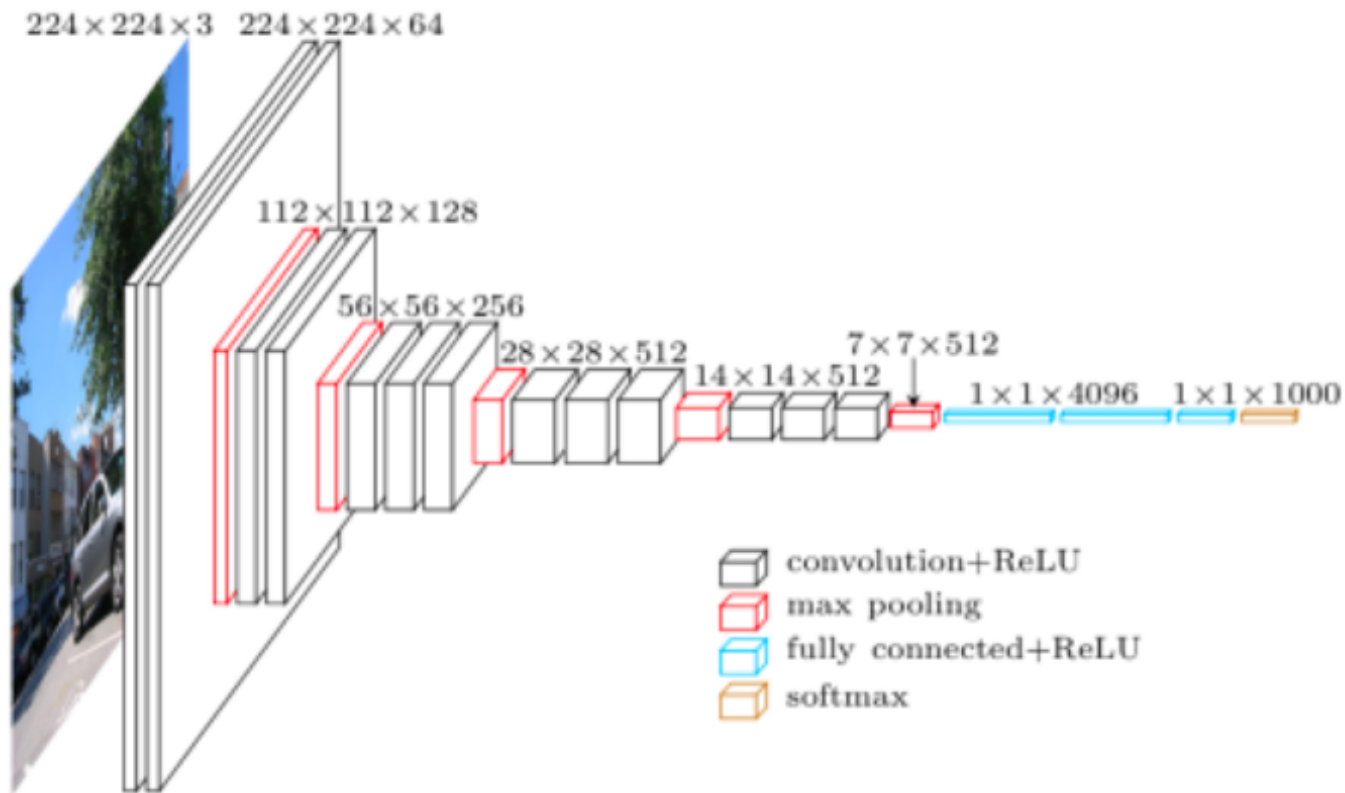# TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2019

# Convolutional Neural Networks (CNNs)

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



Davi Frossard

# Review: Einstein Notation for Linear Threshold Layer
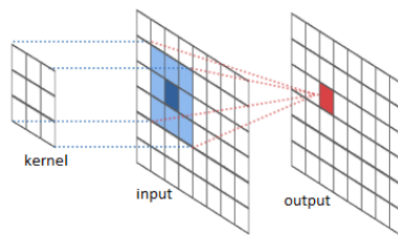
$$y = \sigma\left(W x - B\right)$$

is an abbreviation for

$$y[b,j] = \sigma\left(\left(\sum_i W[j,i]\, x[b,i]\right) - B[j]\right)$$

Think of this as a separate assignment statement for each $(b,j)$.

Each $y[b,j]$ is the output of a "linear threshold unit".

Einstein notation makes all indeces and summations explicit.

# A Convolution Layer



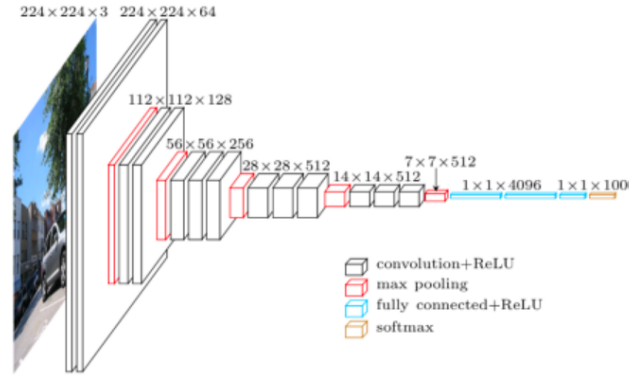$$W[\Delta x, \Delta y, i, j] \qquad L_{\text{IN}}[b, x, y, i] \qquad L_{\text{out}}[b, x, y, j]$$

River Trail Documentation

Procedure CONV($W[\Delta x, \Delta y, i, j], B[j], L_{\text{in}}[b, x, y, i]$)

Return $L_{\text{out}}[b, x, y, j]$

$$= \sigma \left( \left( \sum_{\Delta x, \Delta y, i} W[\Delta x, \Delta y, i, j] \, L_{\text{in}}[b, x + \Delta x, y + \Delta y, i] \right) - B[j] \right)$$
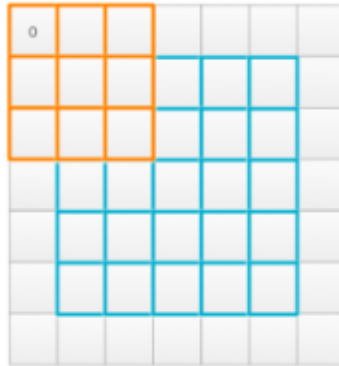
# A Convolution Layer



Each box is a tensor $L[b, x, y, i]$

$$L_{\mathrm{out}}[b, x, y, j]$$

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

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

# 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

`Procedure CONV`$(\Phi,\ L_{\mathrm{in}}[b,x,y,i],\ \mathrm{padding}\ p)$

$$L'_{\mathrm{in}} = \mathrm{Padd}(L_{\mathrm{in}},\ p)$$

$$L_{\mathrm{out}}[b,x,y,j] =$$

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

Return $L_{\mathrm{out}}[b,x,y,j]$

# Reducing Spatial Dimention

# Strides

We can move the filter by a "stride" $s$ for each spatial step.
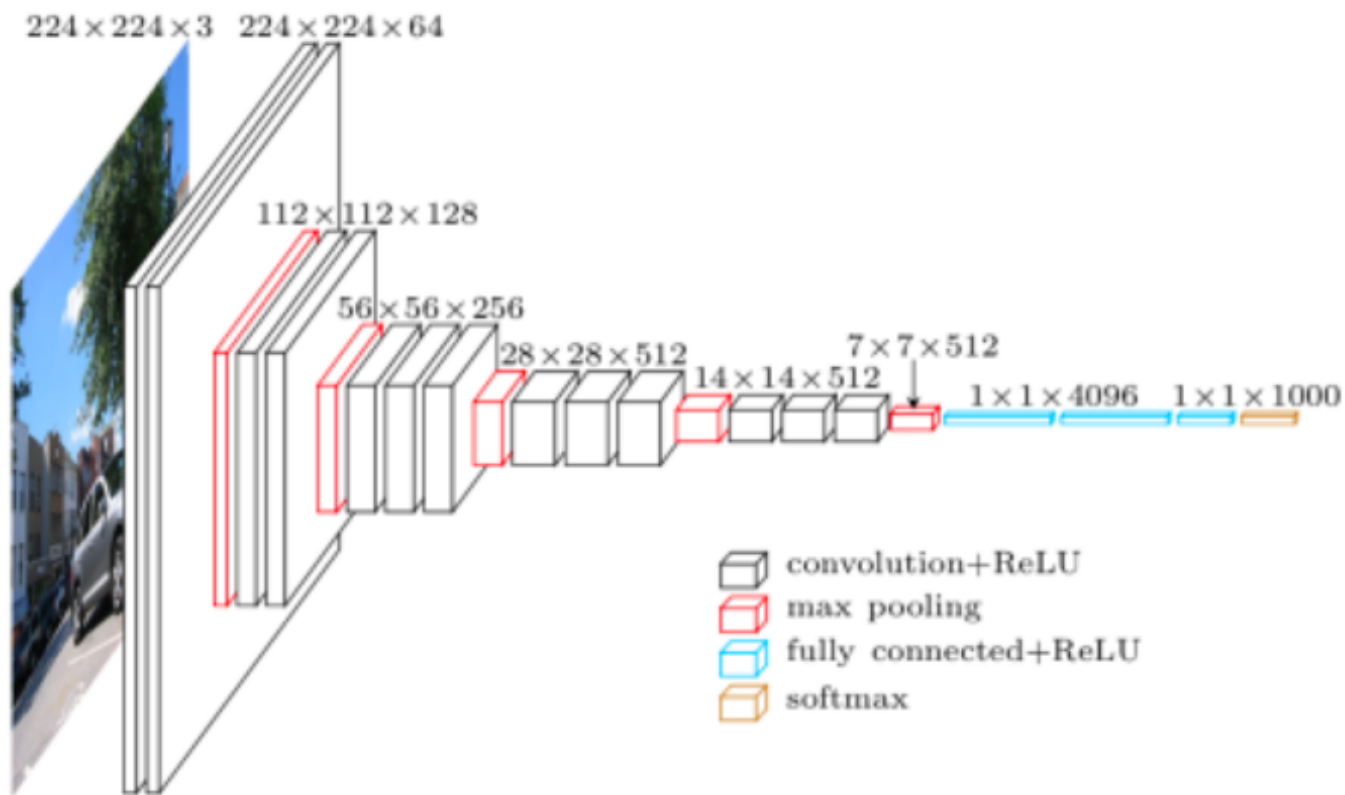
$$L_{\text{out}}[b, x, y, j] =$$

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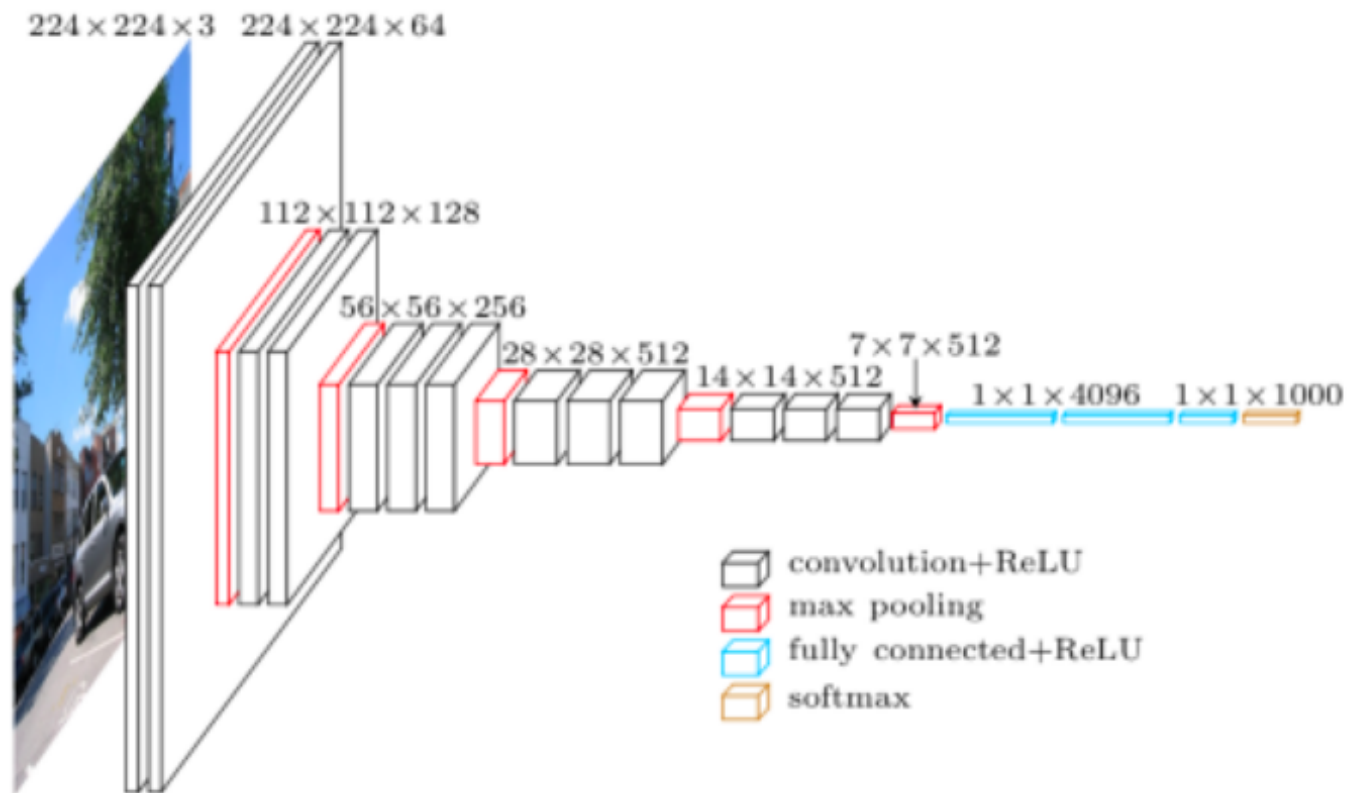
For strides greater than 1 the spatial dimention is reduced.

# Max Pooling

$$L_{\text{out}}[b, x, y, i] = \max_{\Delta x, \Delta y} \; L_{\text{in}}[b, s * x + \Delta x, \; s * y + \Delta y, \; i]$$

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 $L_{\mathrm{in}}[b, x, y, i]$ to $L_{\mathrm{in}}[b, i']$ and then

$$L_{\mathrm{out}}[b, j] = \sigma \left( \left( \sum_{i'} W[j, i'] \, L_{\mathrm{in}}[b, i'] \right) - B[j] \right)$$

13

# Alexnet

Given Input$[227, 227, 3]$

$$L_1[55 \times 55 \times 96] = \mathrm{ReLU}(\mathrm{CONV}(\mathrm{Input}, \Phi_1, \mathrm{width}\ 11, \mathrm{pad}\ 0, \mathrm{stride}\ 4))$$

$$L_2[27 \times 27 \times 96] = \mathrm{MaxPool}(L_1, \mathrm{width}\ 3, \mathrm{stride}\ 2))$$

$$L_3[27 \times 27 \times 256] = \mathrm{ReLU}(\mathrm{CONV}(L_2, \Phi_3, \mathrm{width}\ 5, \mathrm{pad}\ 2, \mathrm{stride}\ 1))$$

$$L_4[13 \times 13 \times 256] = \mathrm{MaxPool}(L_3, \mathrm{width}\ 3, \mathrm{stride}\ 2))$$

$$L_5[13 \times 13 \times 384] = \mathrm{ReLU}(\mathrm{CONV}(L_4, \Phi_5, \mathrm{width}\ 3, \mathrm{pad}\ 1, \mathrm{stride}\ 1))$$

$$L_6[13 \times 13 \times 384] = \mathrm{ReLU}(\mathrm{CONV}(L_5, \Phi_6, \mathrm{width}\ 3, \mathrm{pad}\ 1, \mathrm{stride}\ 1))$$

$$L_7[13 \times 13 \times 256] = \mathrm{ReLU}(\mathrm{CONV}(L_6, \Phi_7, \mathrm{width}\ 3, \mathrm{pad}\ 1, \mathrm{stride}\ 1))$$

$$L_8[6 \times 6 \times 256] = \mathrm{MaxPool}(L_7, \mathrm{width}\ 3, \mathrm{stride}\ 2))$$

$$L_9[4096] = \mathrm{ReLU}(\mathrm{FC}(L_8, \Phi_9))$$

$$L_{10}[4096] = \mathrm{ReLU}(\mathrm{FC}(L_9, \Phi_{10}))$$

$$s[1000] = \mathrm{ReLU}(\mathrm{FC}(L_{10}, \Phi_s))\ \ \text{class scores}$$

# 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

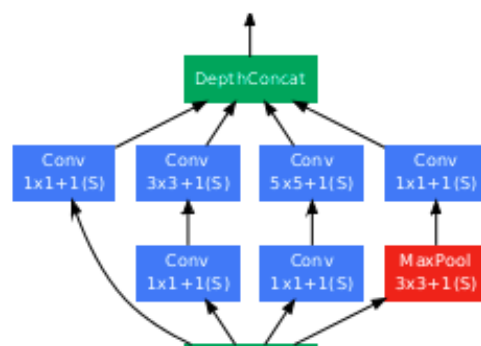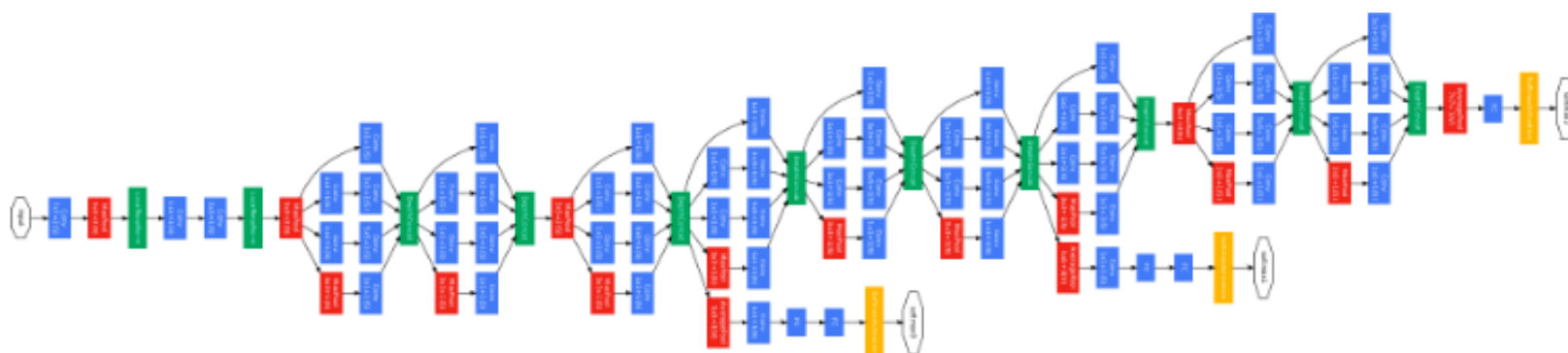softmax

Davi Frossard

# VGG



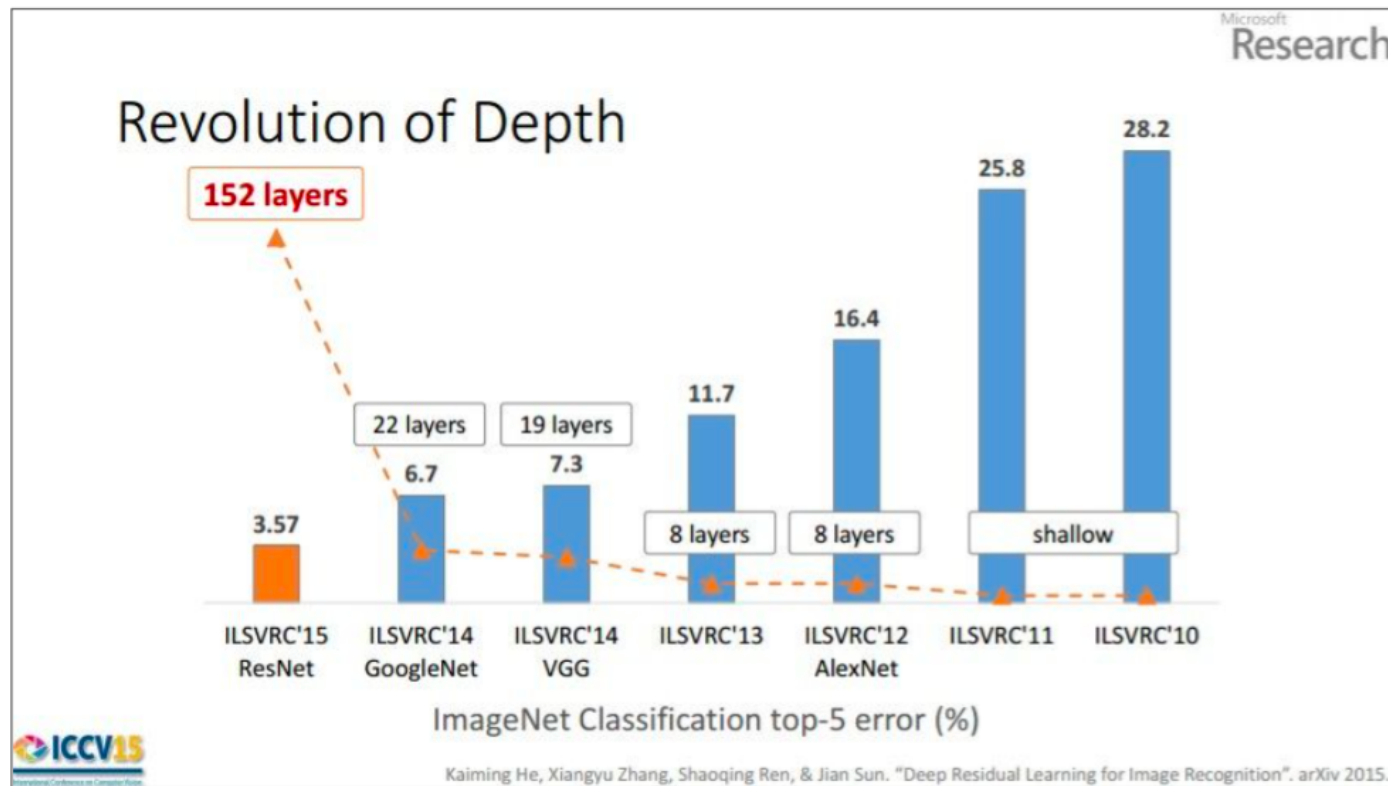AlexNet | VGG16 | VGG19

Stanford CS231

# Inception, Google, 2014

# Imagenet Classification

1000 kinds of objects.



(slide from Kaiming He's recent presentation)

2016 error rate is 3.0%          2017 error rate is 2.25%

# Review of The Swap Rule

$$\tilde{y}[b,j] = \sum_i W[j,i]\, x[b,i]$$

$$y[b,j] = \sigma(\tilde{y}[b,j] - B[j])$$

$$x.\text{grad}[b,i] \mathrel{+}= \sum_j \tilde{y}.\text{grad}[b,j]W[j,i]$$

$$W.\text{grad}[j,i] \mathrel{+}= \sum_b \tilde{y}.\text{grad}[b,j]x[b,i]$$

# Alternative Swap Rule

Intialize all computed tensors to zero and write the program using only `+=`.

$$\text{for } b, i, j \qquad \tilde{y}[b, j] \mathrel{+}= W[j, i]\, x[b, i]$$

$$\text{for } b, i, j \quad x.\text{grad}[b, i] \mathrel{+}= \tilde{y}.\text{grad}[b, j] W[j, i]$$

$$\text{for } b, i, j \; W.\text{grad}[j, i] \mathrel{+}= \tilde{y}.\text{grad}[b, j] x[b, i]$$

one swaps the output with one of the inputs inside the body of the loop.

# Alternative Swap Rule for Convolution

for $b, x, y, i, j, \Delta x, \Delta y$

$\quad \tilde{L}_{\text{out}}[b, x, y, j]$ += $W[\Delta x, \Delta y, i, j] \, L_{\text{in}}[b, x + \Delta x, y + \Delta y, i]$

for $b, x, y, i, j, \Delta x, \Delta y$

$\quad W.\text{grad}[\Delta x, \Delta y, i, j]$ += $\tilde{L}_{\text{out}}.\text{grad}[b, x, y, j] \, L_{\text{in}}[b, x + \Delta x, y + \Delta y, i]$

for $b, x, y, i, j, \Delta x, \Delta y$

$L_{\text{in}}.\text{grad}[b, x + \Delta x, y + \Delta y, i, j]$ += $\tilde{L}_{\text{out}}.\text{grad}[b, x, y, j] \, W[\Delta x, \Delta y, i, j]$
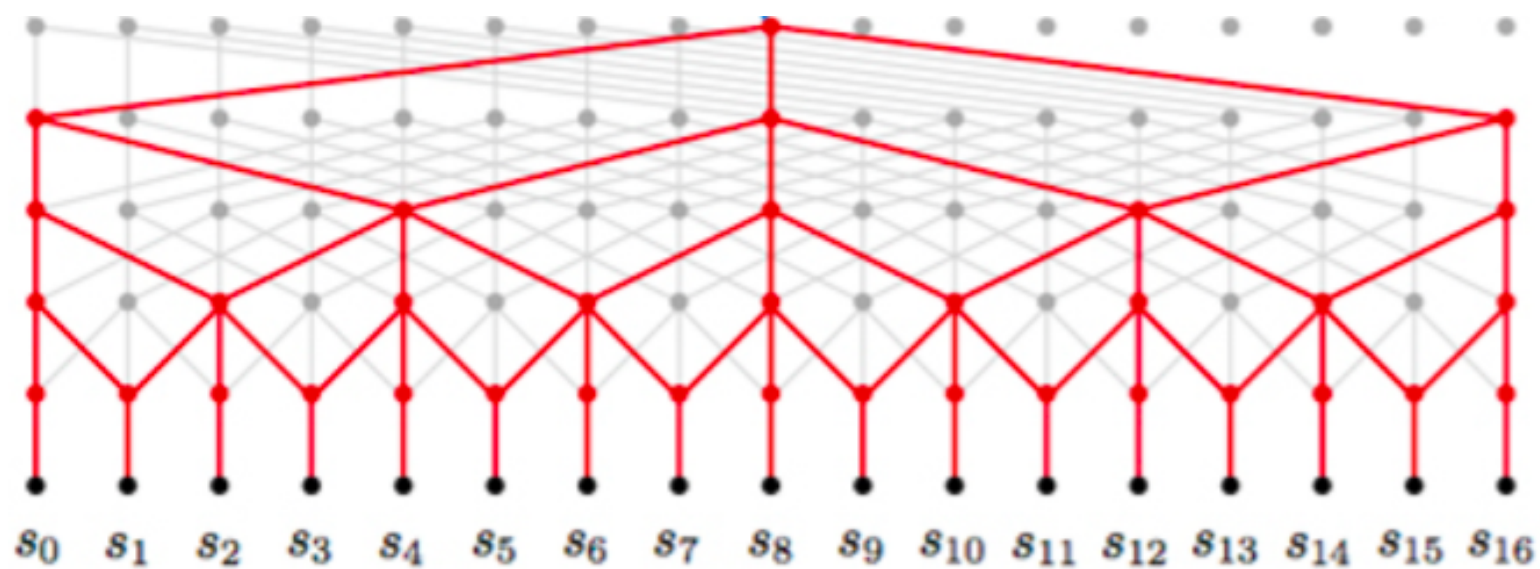
# Image to Column (Im2C)

Reduce convolution to matrix multiplication — more space but faster.

$$L_{\text{in}}[b, x, y, \Delta x, \Delta y, i] = L_{\text{in}}[b, x + \Delta x, y + \Delta y, i]$$

$$\tilde{L}_{\text{out}}[b, x, y, j]$$

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

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

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

# Fully Convolutional Networks

# Dilation

We can "dilate" the filter by introducing an image step size $d$ for each step in the filter coordinates.

$$\tilde{L}_{\text{out}}[b, x, y, j] = W[\Delta x, \Delta y, i, j]L_{\text{in}}[b, x + d * \Delta x, y + d * \Delta y, i] + B[j]$$

This is used for "fully convolutional" CNNs.

# Summary

- Convolution

- Padding

- Stides

- Max Pooling

- Fully Connected Layers

- Dilation

# 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).

END