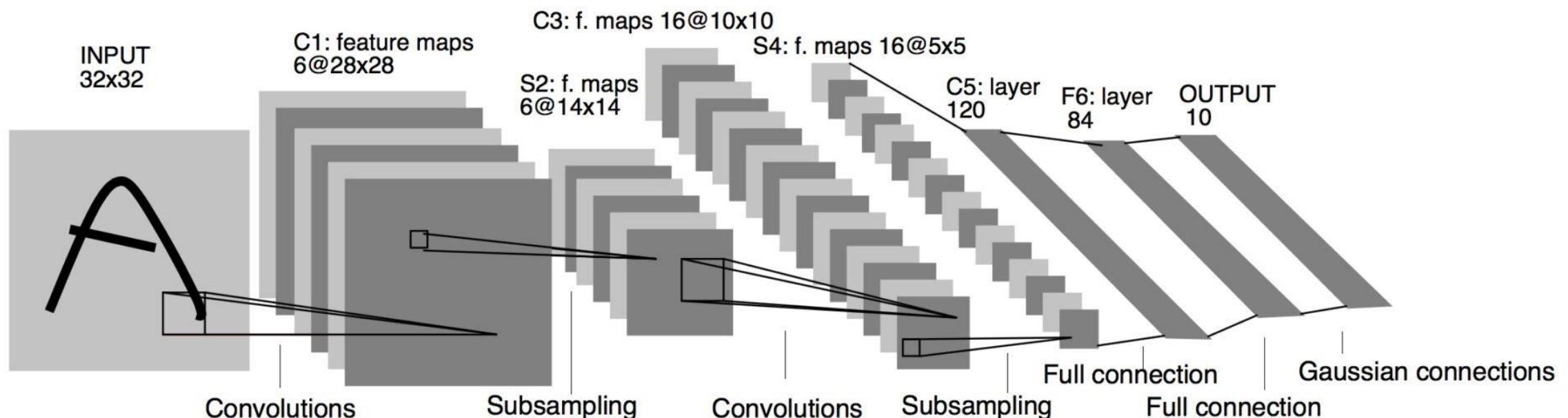
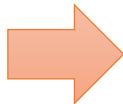


# Deep Learning Architectures: Convolutional Neural Networks



# Neural network with images



[object label]

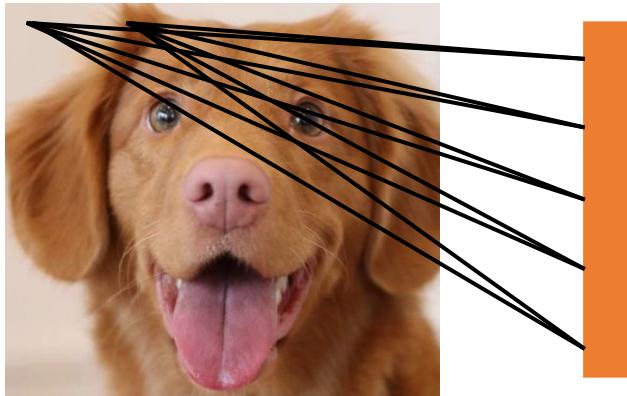
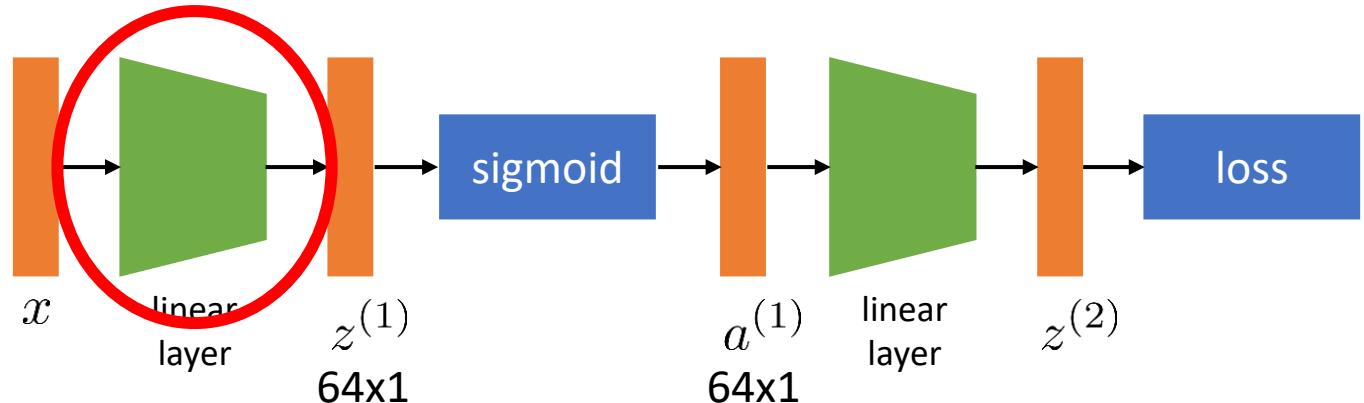


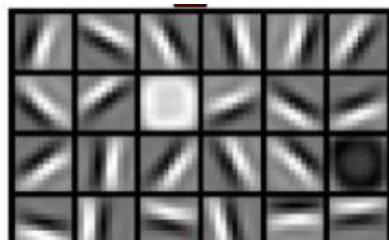
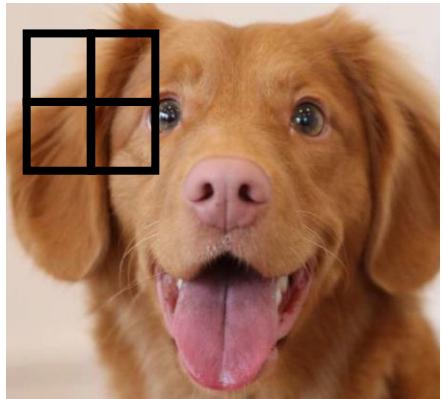
image is  $128 \times 128 \times 3 = 49,152$

$z^{(1)}$  is 64-dim

$64 \times 49,152 \approx 3,000,000$

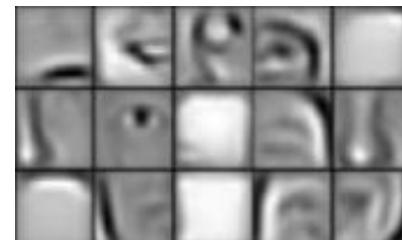
We need a better way!

# An idea...



Layer 1:

edge detectors?



Layer 2:

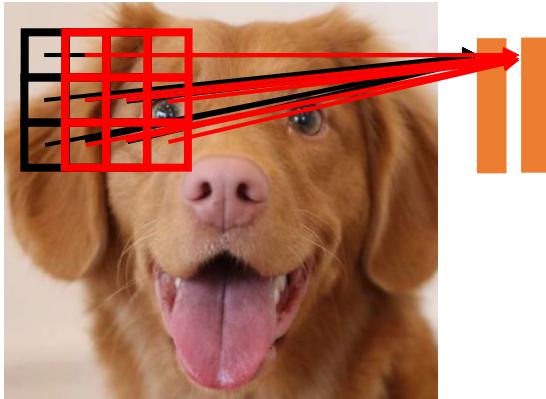
ears? noses?

**Observation:** many useful image features are **local**

to tell if a particular patch of image contains a feature, enough to look at the local patch

# An idea...

**Observation:** many useful image features are **local**

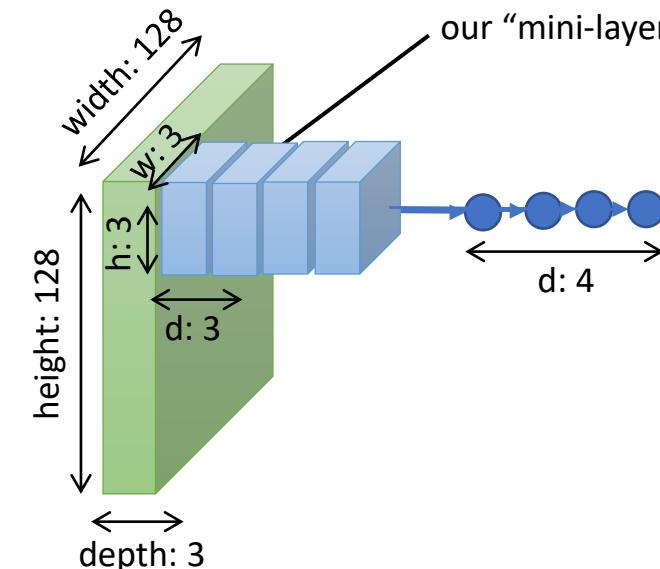


patch is  $3 \times 3 \times 3 = 27$

$z^{(1)}$  is 64-dim

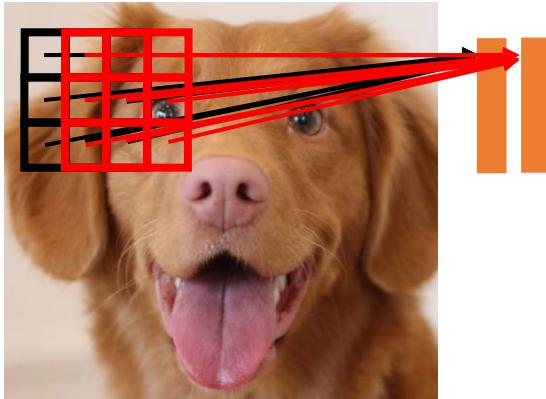
$$64 \times 27 = 1728$$

We get a **different** output at each image location!



# An idea...

**Observation:** many useful image features are **local**

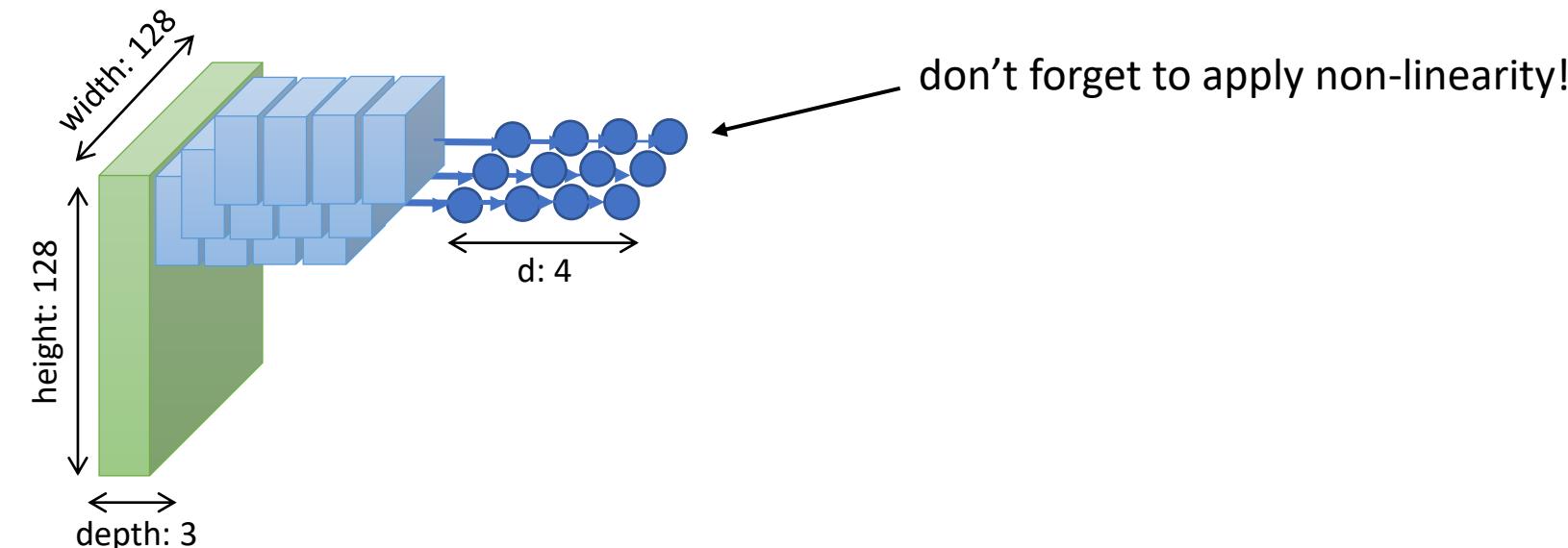


patch is  $3 \times 3 \times 3 = 27$

$z^{(1)}$  is 64-dim

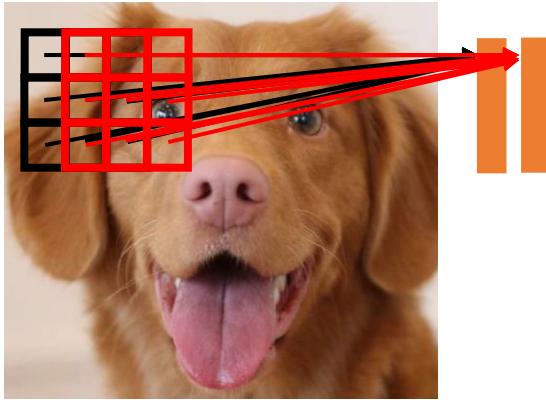
$$64 \times 27 = 1728$$

We get a **different** output at each image location!



# An idea...

**Observation:** many useful image features are **local**

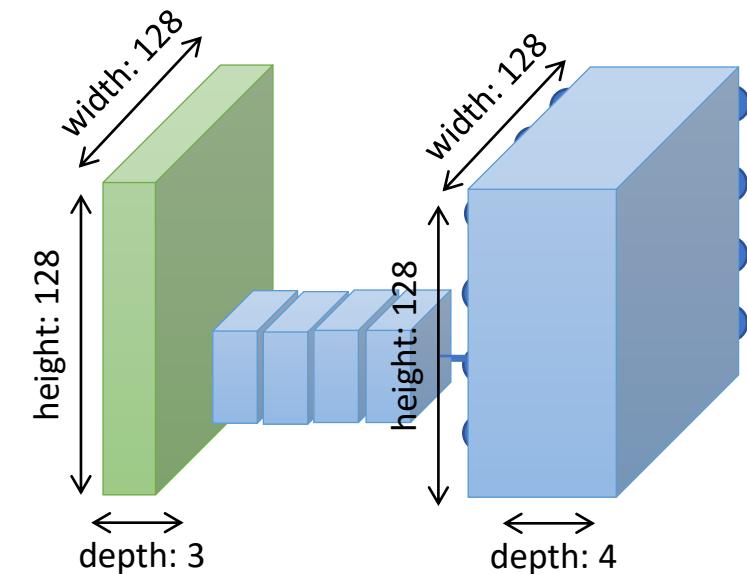


patch is  $3 \times 3 \times 3 = 27$

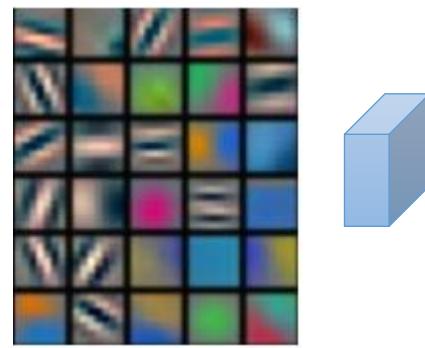
$z^{(1)}$  is 64-dim

$$64 \times 27 = 1728$$

We get a **different** output at each image location!

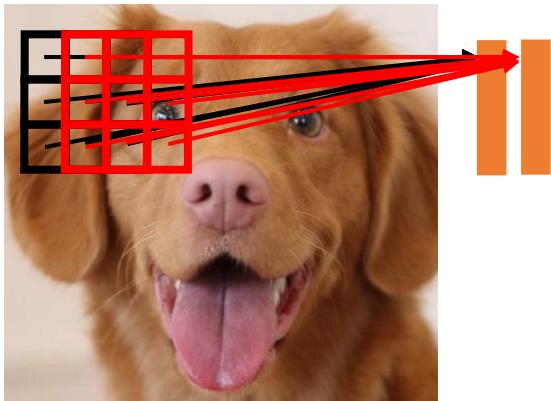


What do they look like?



# An idea...

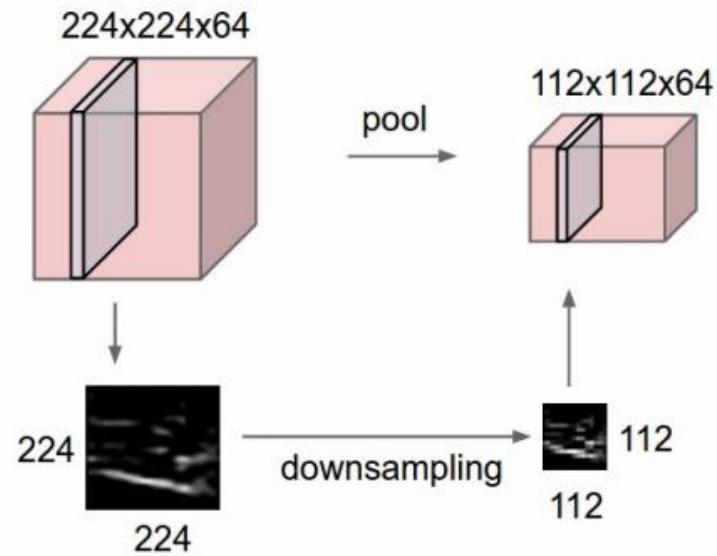
**Observation:** many useful image features are **local**



patch is  $3 \times 3 \times 3 = 27$

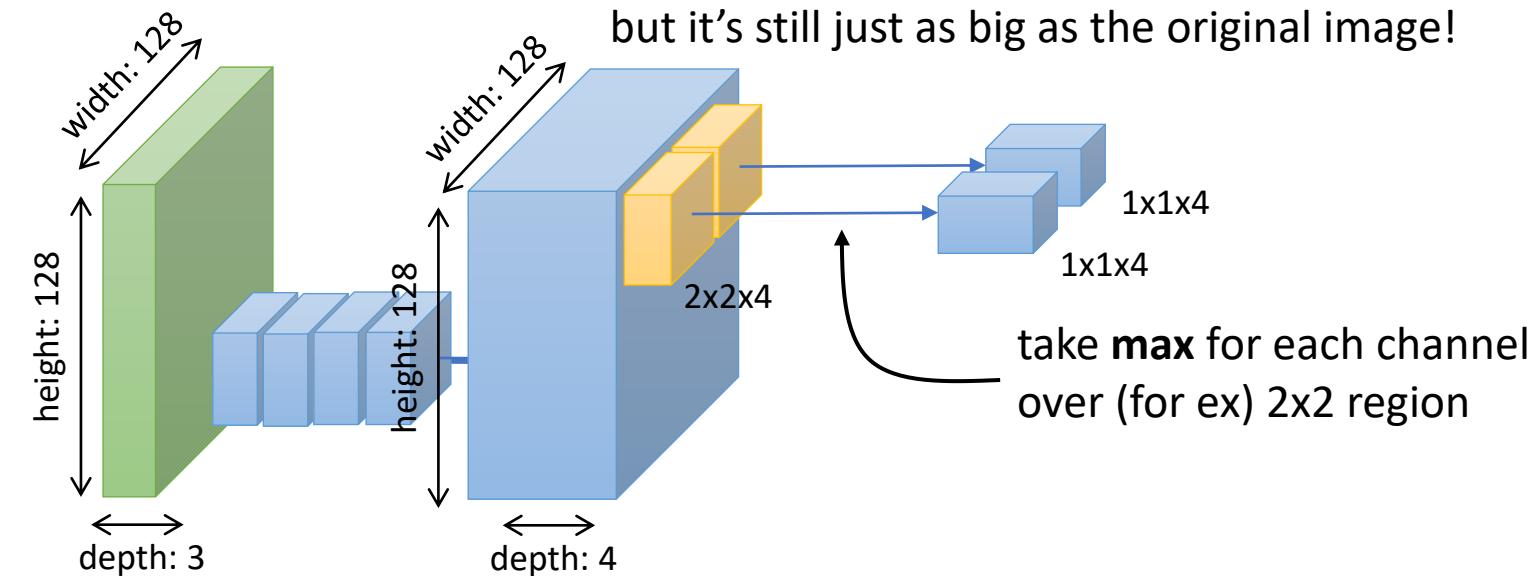
$z^{(1)}$  is 64-dim

$$64 \times 27 = 1728$$



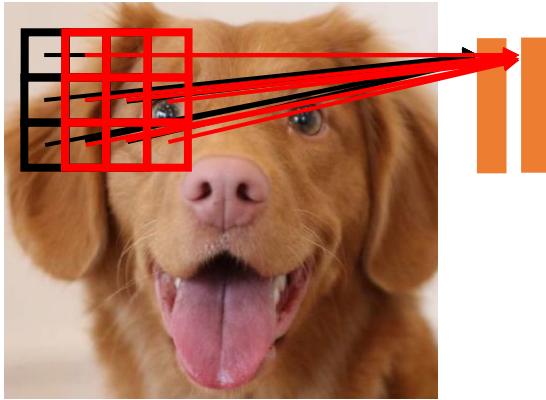
We get a **different** output at each image location!

but it's still just as big as the original image!



# An idea...

**Observation:** many useful image features are **local**

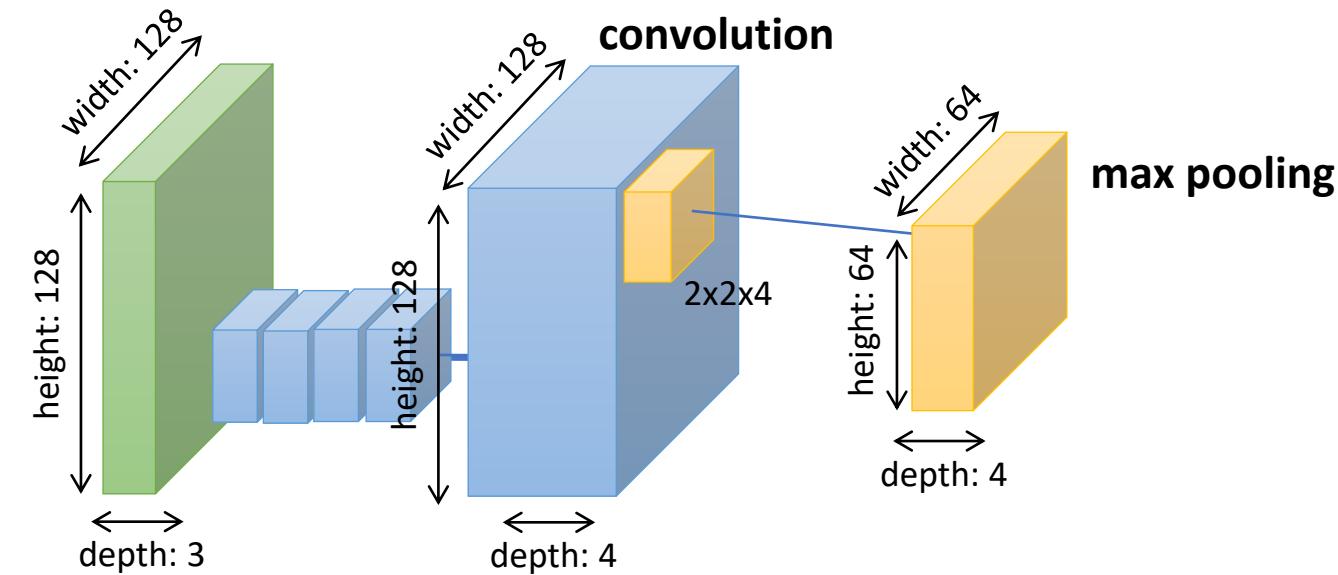


patch is  $3 \times 3 \times 3 = 27$

$z^{(1)}$  is 64-dim

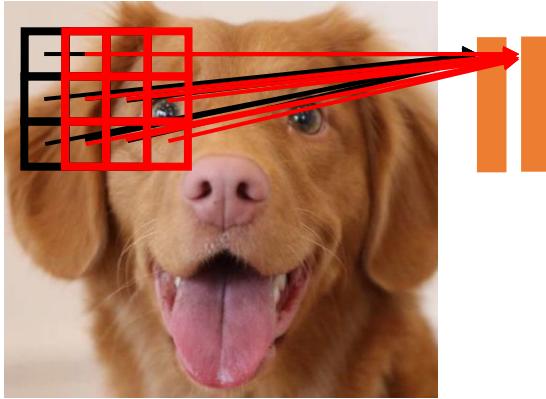
$$64 \times 27 = 1728$$

We get a **different** output at each image location!



# An idea...

**Observation:** many useful image features are **local**

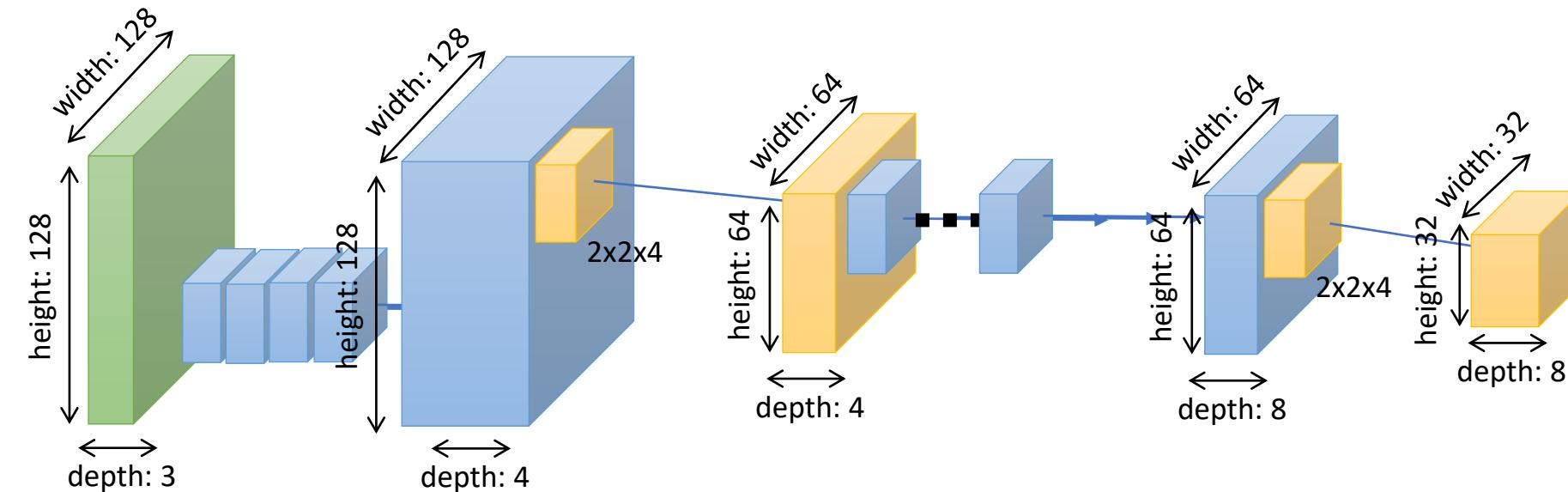


patch is  $3 \times 3 \times 3 = 27$

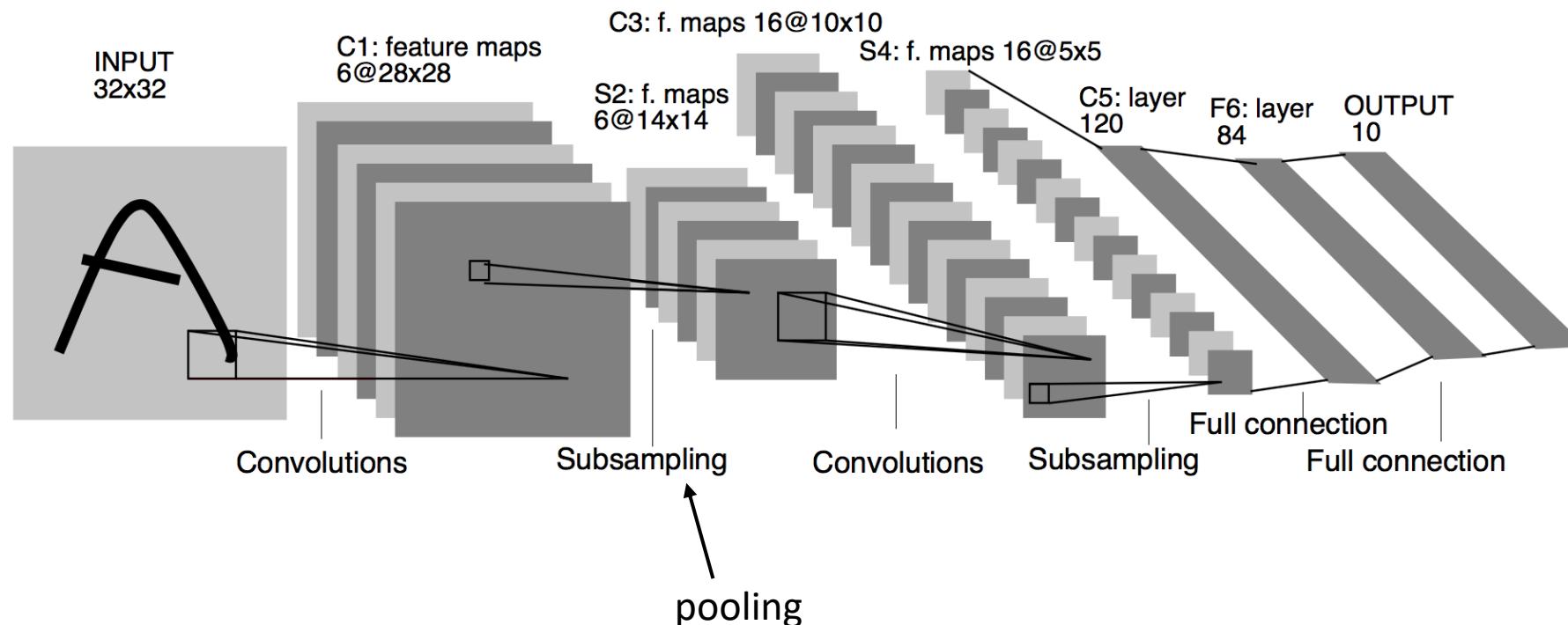
$z^{(1)}$  is 64-dim

$$64 \times 27 = 1728$$

We get a **different** output at each image location!



# What does a real conv net look like?

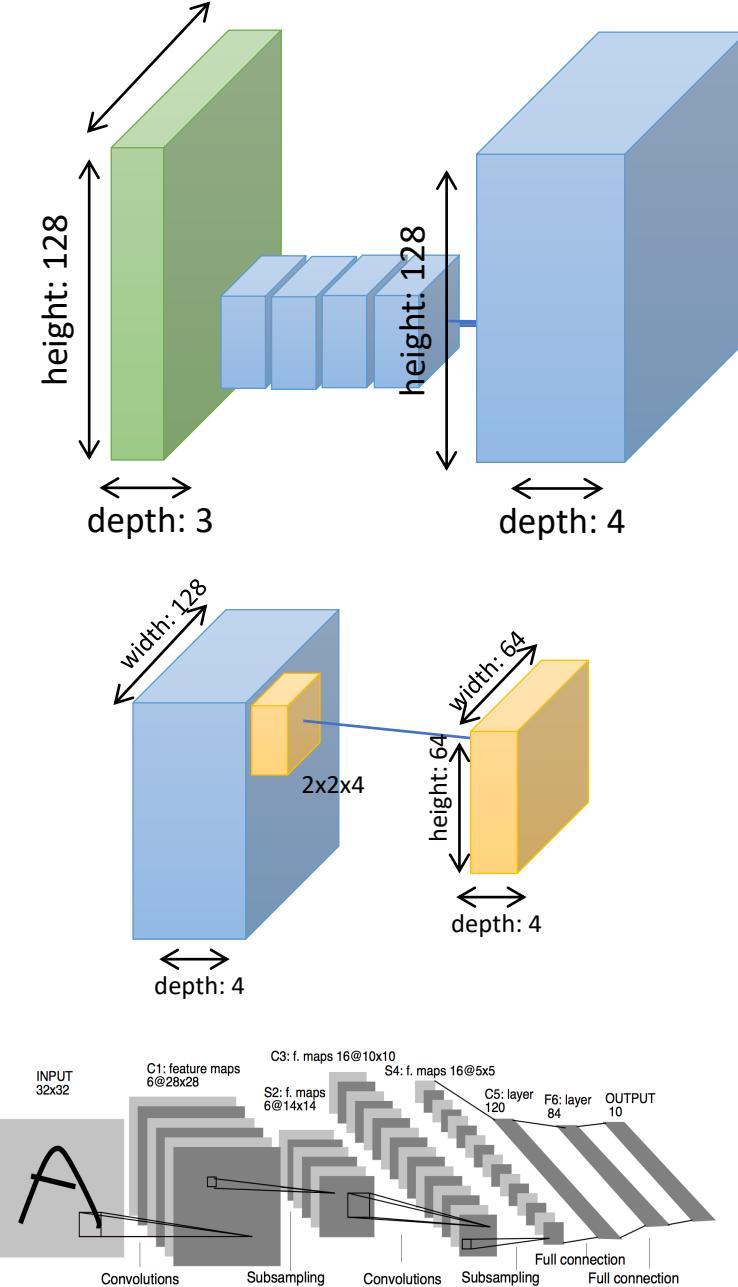


"LeNet" network for handwritten digit recognition

# Implementing convolutional layers

# Summary

- **Convolutional layer**
  - A way to avoid needing millions of parameters with images
  - Each layer is “local”
  - Each layer produces an “image” with (roughly) the same width & height, and number of channels = number of filters
- **Pooling**
  - If we ever want to get down to a single output, we must reduce resolution as we go
  - Max pooling: downsample the “image” at each layer, taking the max in each region
  - This makes it robust to small translation changes
- **Finishing it up**
  - At the end, we get something small enough that we can “flatten” it (turn it into a vector), and feed into a standard fully connected layer



# Convolutional layer in equations

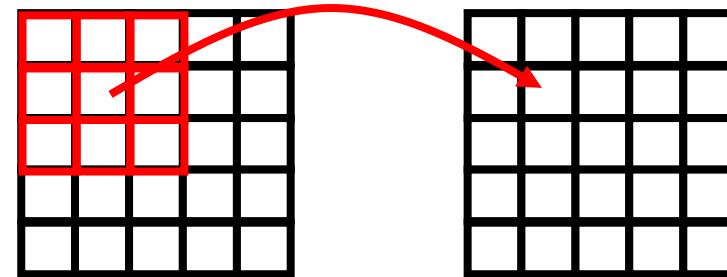
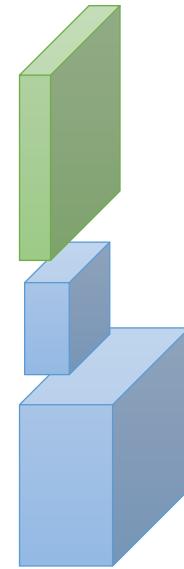
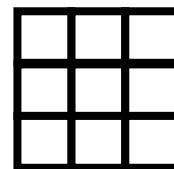
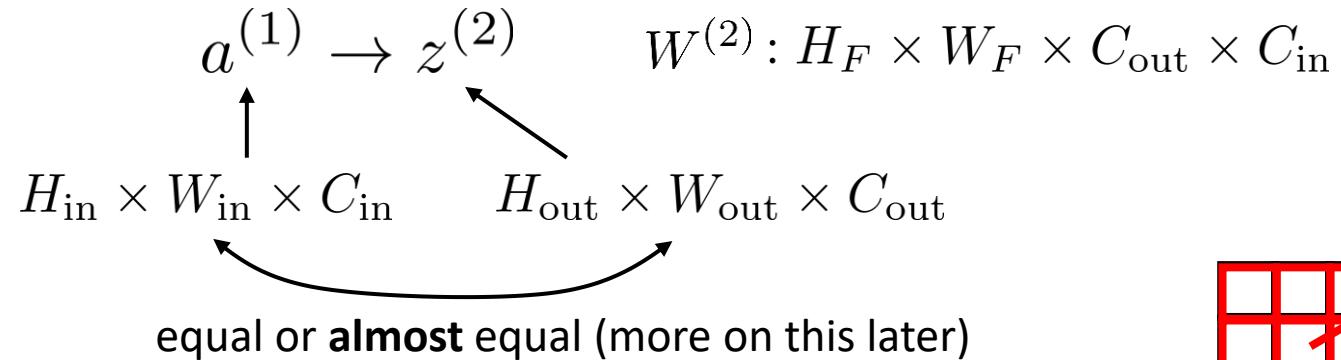
all these operations will involve  $N$ -dimensional arrays

often used synonymously with *tensor*

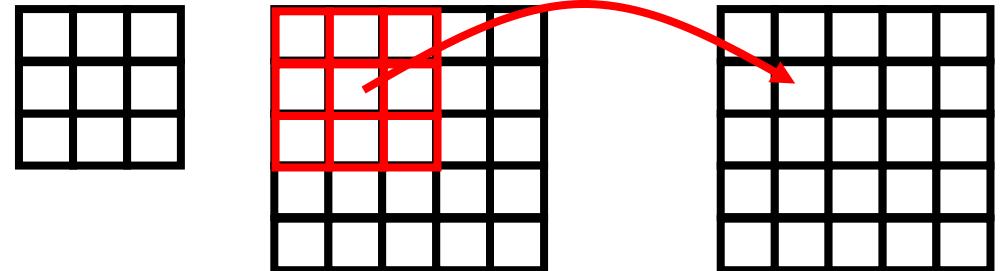
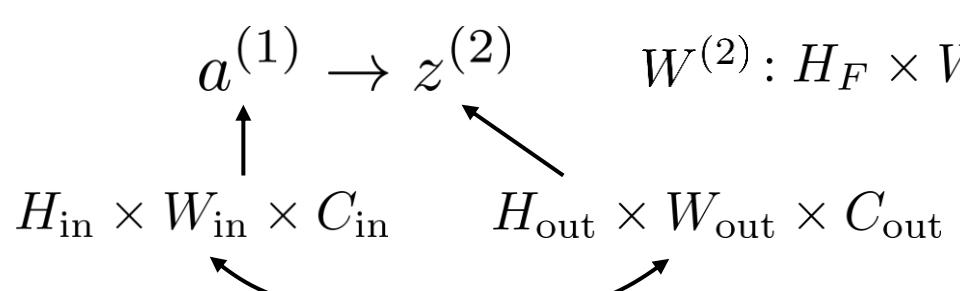
input image: HEIGHT  $\times$  WIDTH  $\times$  CHANNELS

filter: FLT.HEIGHT  $\times$  FLT.WIDTH  $\times$  OUTPUT CHAN  $\times$  INPUT CHAN

activations: HEIGHT  $\times$  WIDTH  $\times$  LAYER.CHANNELS



# Convolutional layer in equations



equal or **almost** equal (more on this later)

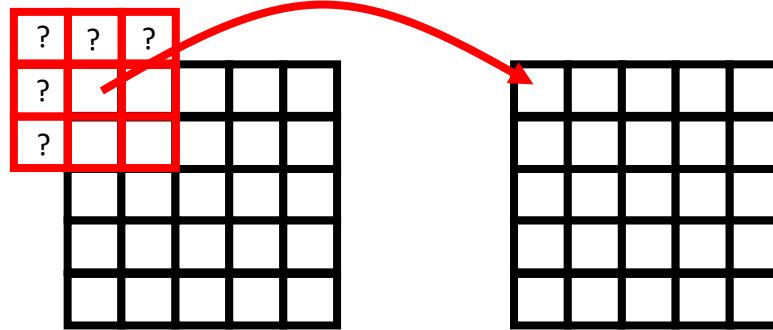
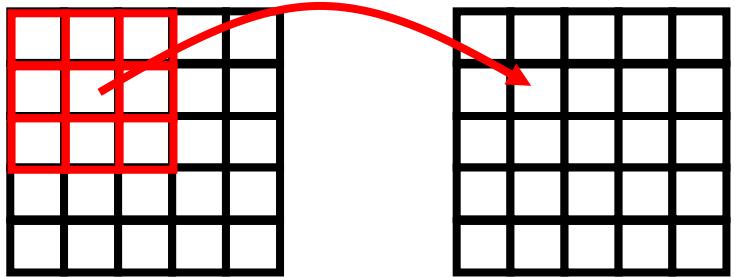
$$z^{(2)}[i, j, k] = \sum_{l=0}^{H_F-1} \sum_{m=0}^{H_W-1} \sum_{n=0}^{C_{\text{in}}-1} W^{(2)}[l, m, k, n] a^{(1)}[i + l - (H_F - 1)/2, j + m - (H_W - 1)/2, n]$$

$$z^{(2)}[i, j] = \sum_{l=0}^{H_F-1} \sum_{m=0}^{H_W-1} W^{(2)}[l, m] a^{(1)}[i + l - (H_F - 1)/2, j + m - (H_W - 1)/2]$$

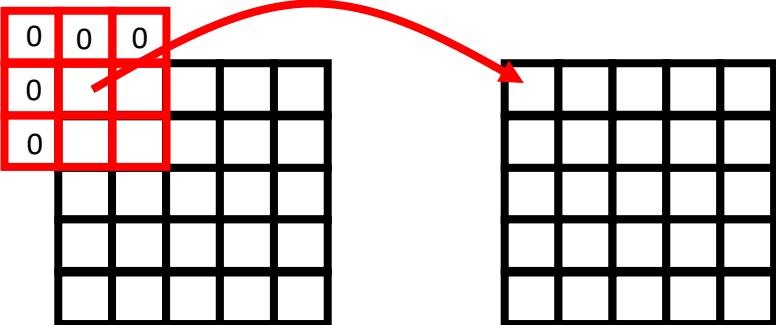
$$a^{(2)}[i, j, k] = \sigma(z^{(2)}[i, j, k]) \quad \text{Activation function applied per element, just like before}$$

Simple principle, but a bit complicated to write

# Padding and edges



**Option 2: zero pad**



**Detail:** remember to subtract the image mean first  
(fancier contrast normalization often used in practice)

**Advantage:** simple, size is preserved

**Disadvantage:** weird effect at boundary  
(this is usually not a problem, hence  
why this method is so popular)

# Strided convolutions

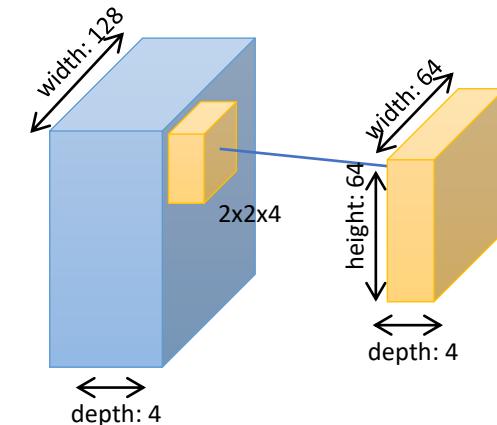
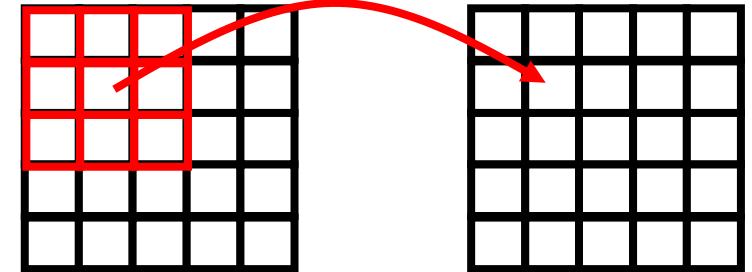
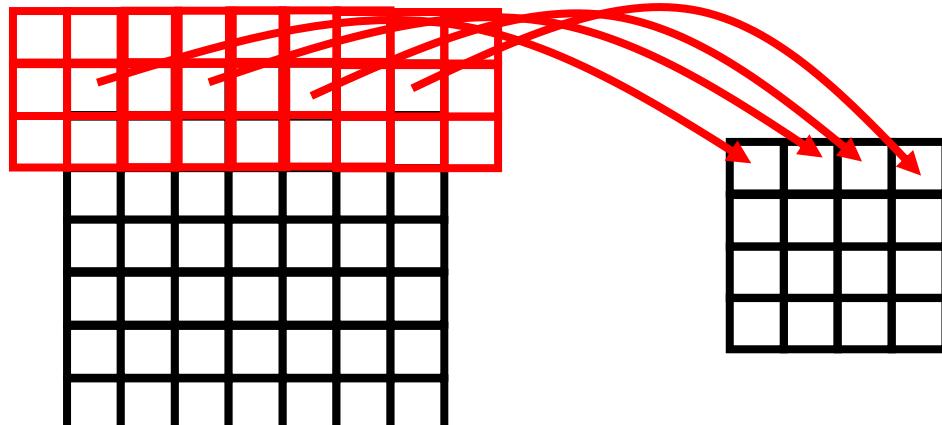
standard conv net structure at each layer:

1. Apply conv,  $H \times W \times C_{\text{in}} \rightarrow H \times W \times C_{\text{out}}$
2. Apply activation func  $\sigma$ ,  $H \times W \times C_{\text{out}} \rightarrow H \times W \times C_{\text{out}}$
3. Apply pooling (width  $N$ ),  $H \times W \times C_{\text{out}} \rightarrow H/N \times W/N \times C_{\text{out}}$

this can be very expensive computationally

$C_{\text{out}} \times C_{\text{in}}$  matrix multiply at each position in  $H \times W$  image!

**Idea:** what if skip over some positions?



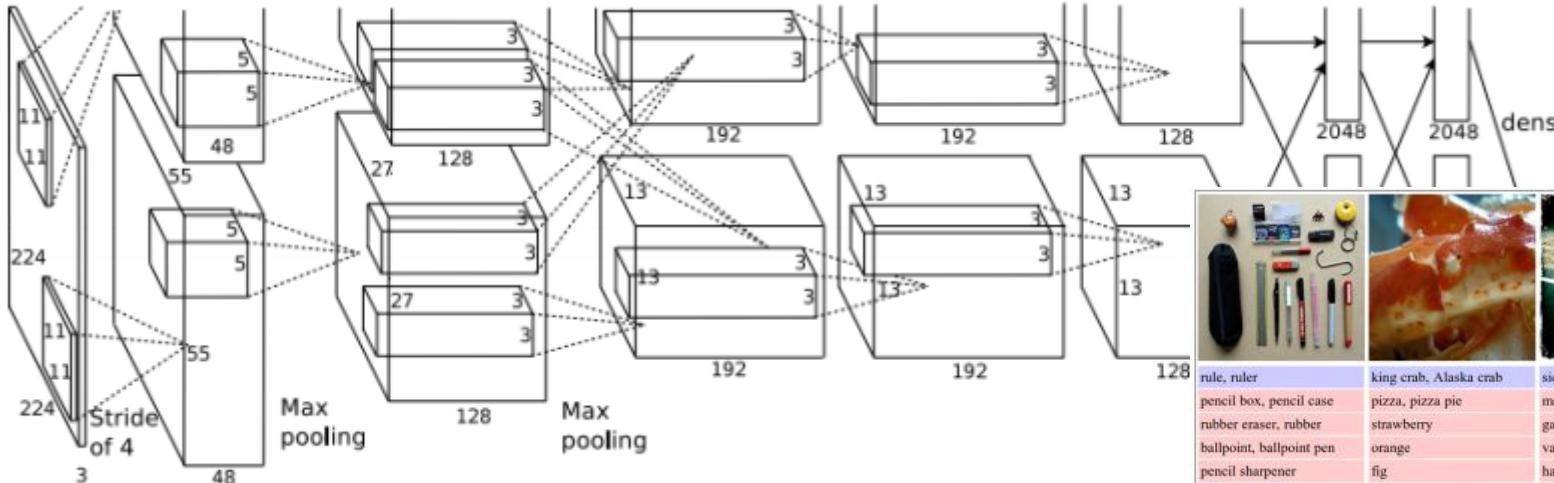
Amount of skipping is called the **stride**

Some people think that strided convolutions are just as good as conv + pooling

# Examples of convolutional neural networks

# AlexNet

[Krizhevsky et al. 2012]



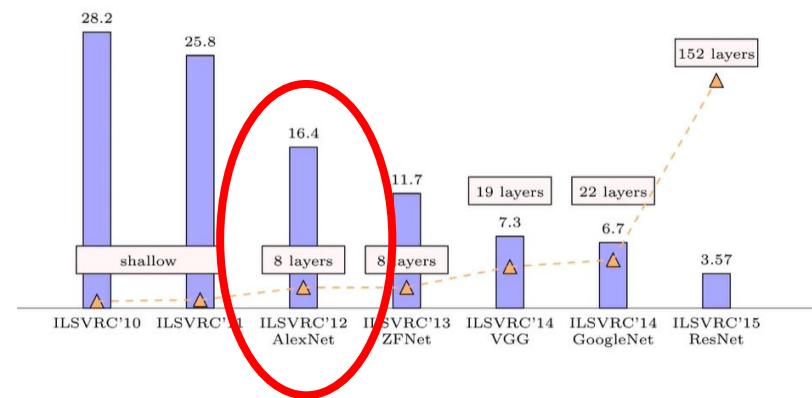
ILSVRC (ImageNet), 2009: 1.5 million images  
1000 categories



pencil box, pencil case	pizza, pizza pie	maze, labyrinth	pill bottle	stethoscope	vase	<span style="background-color: #90EE90;">schipperke</span>
rubber eraser, rubber	strawberry	gar, garfish	water bottle	whistle	pitcher, ewer	<span style="background-color: #FFB6C1;">groenendael</span>
ballpoint, ballpoint pen	orange	valley, vale	lotion	ice lolly, lolly	coffeepot	<span style="background-color: #F0E6F6;">doormat, welcome mat</span>
pencil sharpener	fig	hammerhead	hair spray	hair spray	mask	<span style="background-color: #F0E6F6;">teddy, teddy bear</span>
carpenter's kit, tool kit	ice cream, icecream	sea snake	beer bottle	maypole	cup	<span style="background-color: #F0E6F6;">jigsaw puzzle</span>

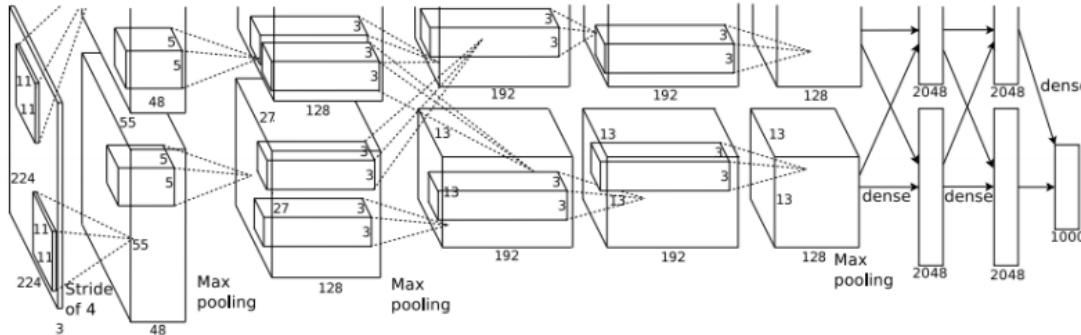
## Why is this model important?

- “Classic” medium-depth convolutional network design (a bit like a modernized version of LeNet)
- Widely known for being the first neural network to attain state-of-the-art results on the ImageNet large-scale visual recognition challenge (ILSVRC)



# AlexNet

[Krizhevsky et al. 2012]



**CONV1:** 11x11x96, Stride 4, maps 224x224x3 -> 55x55x96 [without zero padding]

**POOL1:** 3x3x96, Stride 2, maps 55x55x96 -> 27x27x96

**NORM1:** Local normalization layer

**CONV2:** 5x5x256, Stride 1, maps 27x27x96 -> 27x27x256 [**with** zero padding]

**POOL2:** 3x3x256, Stride 2, maps 27x27x256 -> 13x13x256

**NORM2:** Local normalization layer

**CONV3:** 3x3x384, Stride 1, maps 13x13x256 -> 13x13x384 [**with** zero padding]

**CONV4:** 3x3x384, Stride 1, maps 13x13x384 -> 13x13x384 [**with** zero padding]

**CONV5:** 3x3x256, Stride 1, maps 13x13x256 -> 13x13x256 [**with** zero padding]

**POOL3:** 3x3x256, Stride 2, maps 13x13x256 -> 6x6x256

**FC6:** 6x6x256 -> 9,216 -> 4,096 [matrix is 4,096 x 9,216]

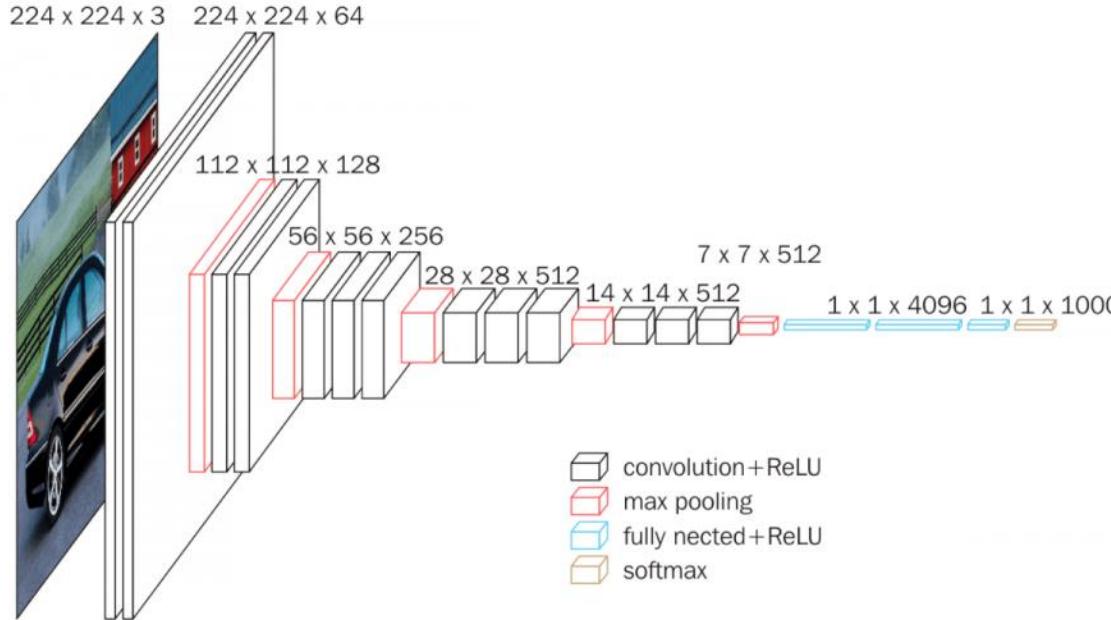
**FC7:** 4,096 -> 4,096

**FC8:** 4,096 -> 1,000

**SOFTMAX**

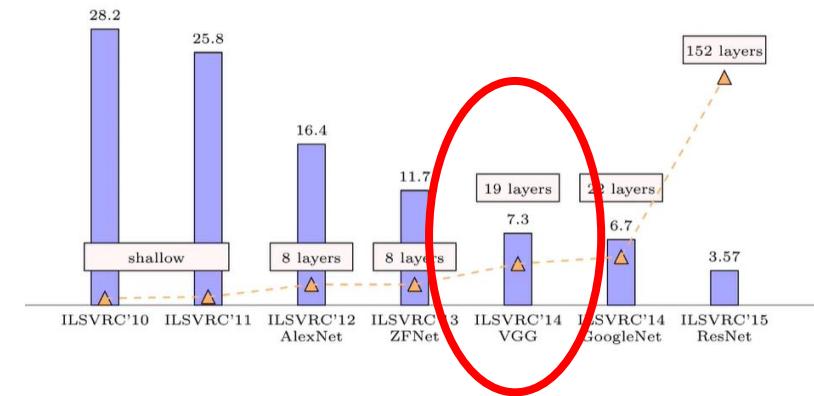
- Don't forget: ReLU nonlinearities after every CONV or FC layer (except the last one!)
- Trained with regularization (we'll learn about these later):
  - Data augmentation
  - Dropout
- Local normalization (not used much anymore, but there are other types of normalization we do use)

# VGG



## Why is this model important?

- Still often used today
- Big increase in **depth** over previous best model
- Start seeing “homogenous” stacks of multiple convolutions interspersed with resolution reduction



# VGG

**CONV:** 3x3x64, maps 224x224x3 -> 224x224x64

**CONV:** 3x3x64, maps 224x224x64 -> 224x224x64

**POOL:** 2x2, maps 224x224x64 -> 112x112x64

**CONV:** 3x3x128, maps 112x112x64 -> 112x112x128

**CONV:** 3x3x128, maps 112x112x128 -> 112x112x128

**POOL:** 2x2, maps 112x112x128 -> 56x56x128

**CONV:** 3x3x256, maps 56x56x128 -> 56x56x256

**CONV:** 3x3x256, maps 56x56x256 -> 56x56x256

**CONV:** 3x3x256, maps 56x56x256 -> 56x56x256

**POOL:** 2x2, maps 56x56x256 -> 28x28x256

**CONV:** 3x3x512, maps 28x28x256 -> 28x28x512

**CONV:** 3x3x512, maps 28x28x512 -> 28x28x512

**CONV:** 3x3x512, maps 28x28x512 -> 28x28x512

**POOL:** 2x2, maps 28x28x512 -> 14x14x512

**CONV:** 3x3x512, maps 14x14x512 -> 14x14x512

**CONV:** 3x3x512, maps 14x14x512 -> 14x14x512

**CONV:** 3x3x512, maps 14x14x512 -> 14x14x512

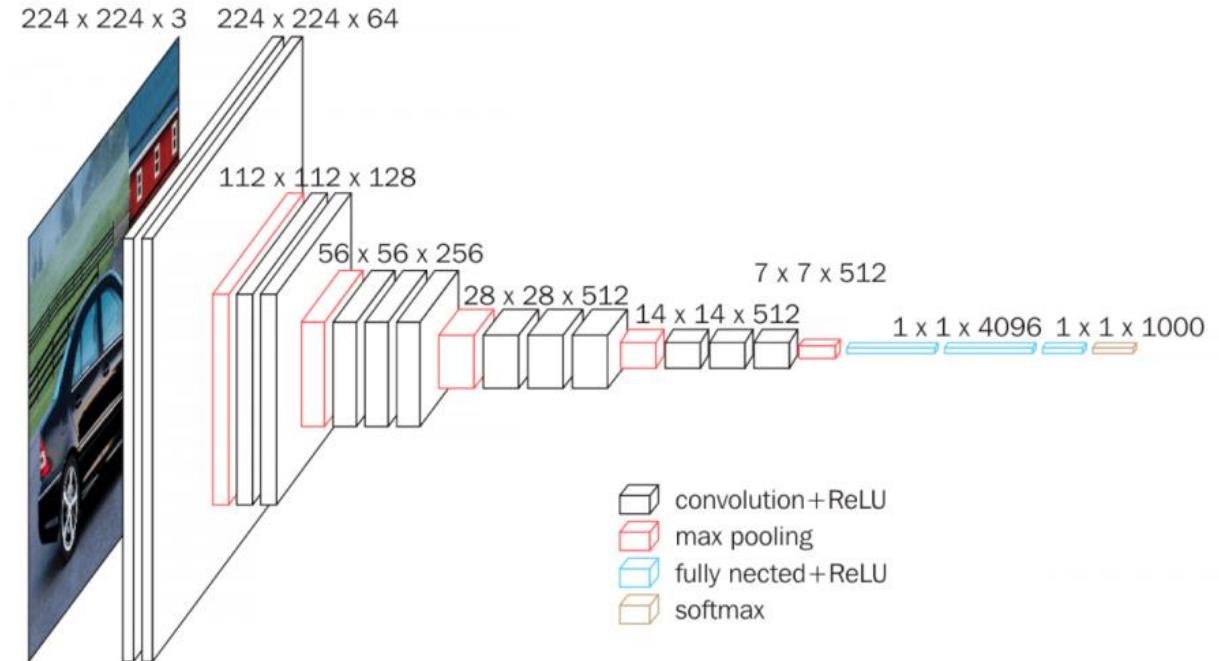
**POOL:** 2x2, maps 14x14x512 -> 7x7x512

**FC:** 7x7x512 -> 25,088 -> 4,096 ← almost all parameters are here

**FC:** 4,096 -> 4,096

**FC:** 4,096 -> 1,000

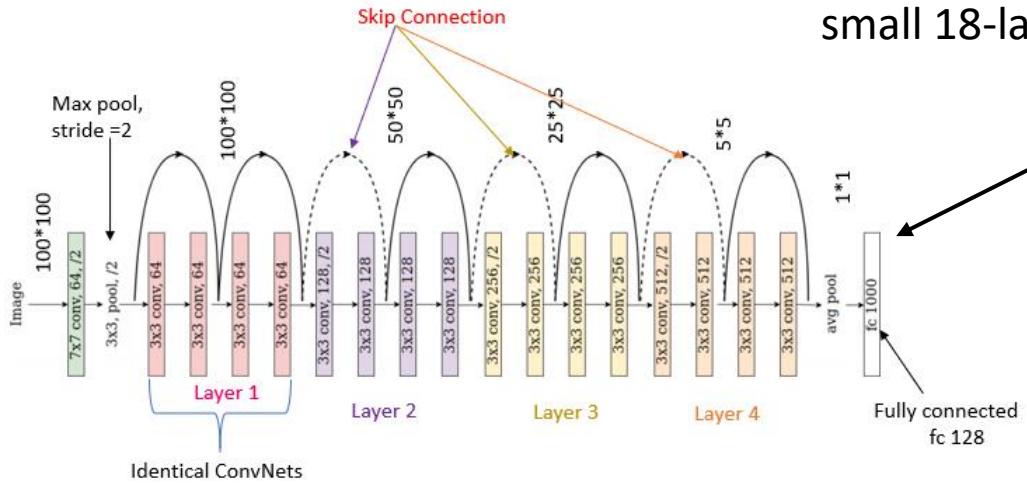
**SOFTMAX**



- More layers = more processing, which is why we see repeated blocks
- Which parts use the most memory?
- Which parts have the most parameters?

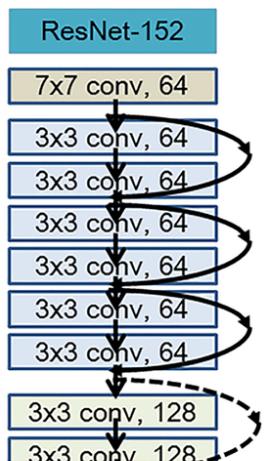
# ResNet

## 152 layers!

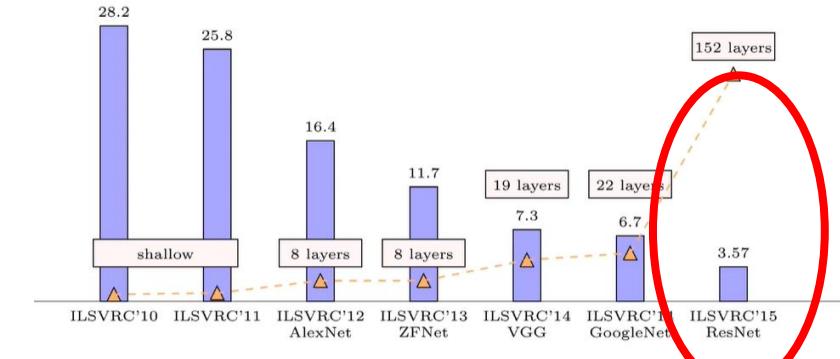
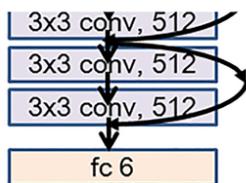


small 18-layer prototype (ResNet-18)

don't bother with huge FC layer at the end, just average pool over all positions and have one linear layer

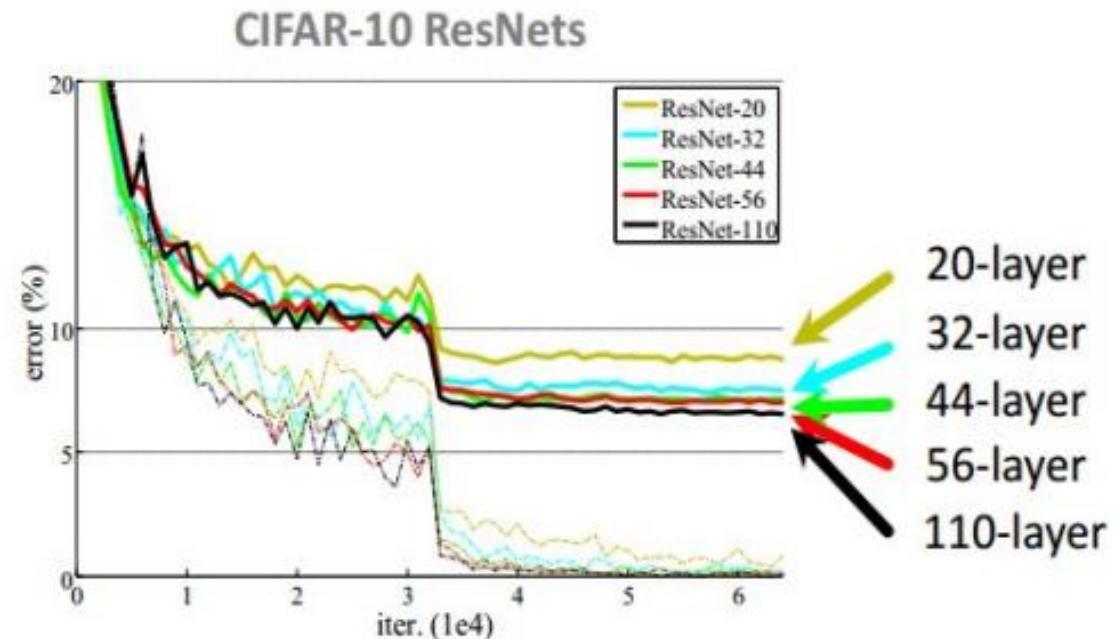
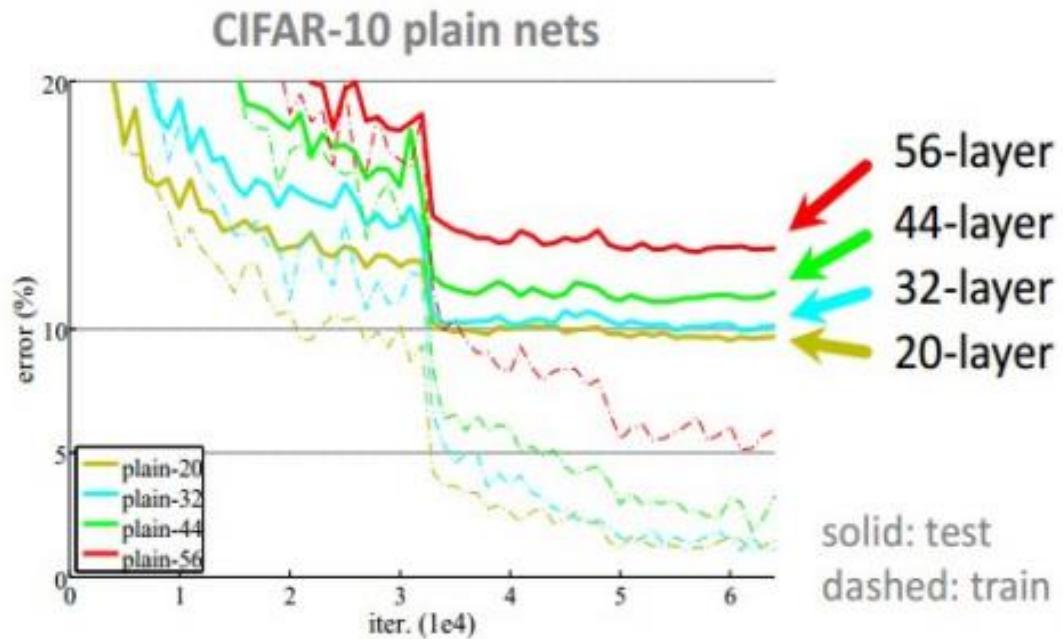


152 layers

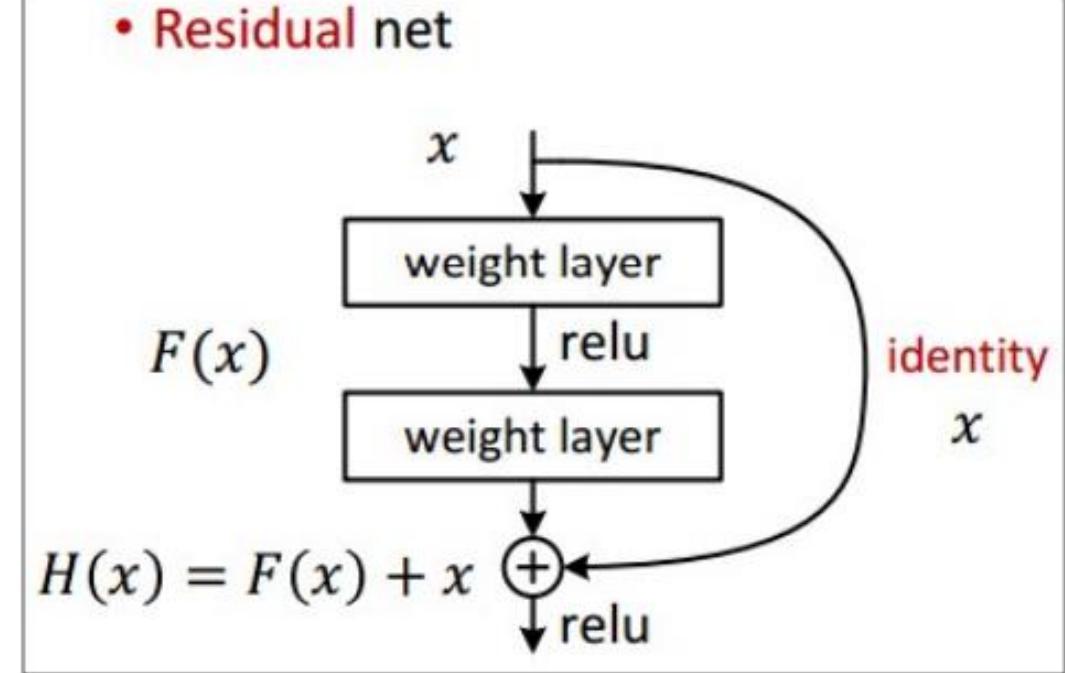
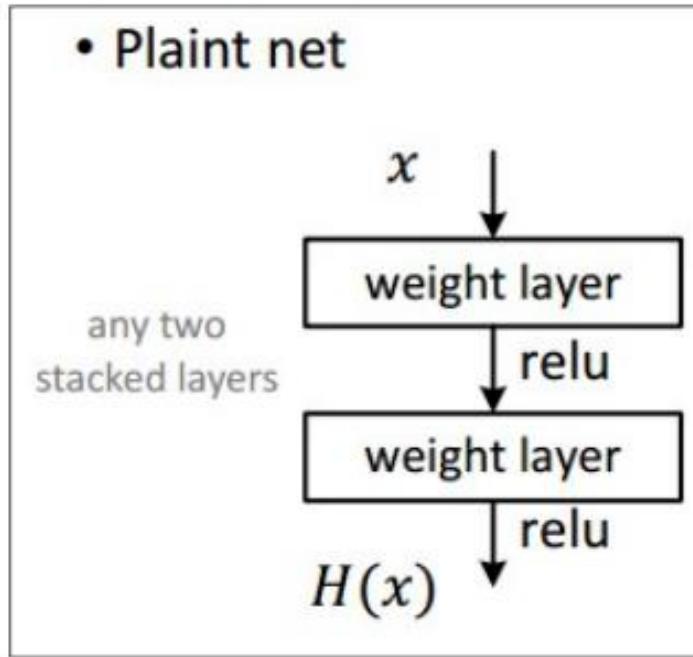


# ResNet

## CIFAR-10 experiments

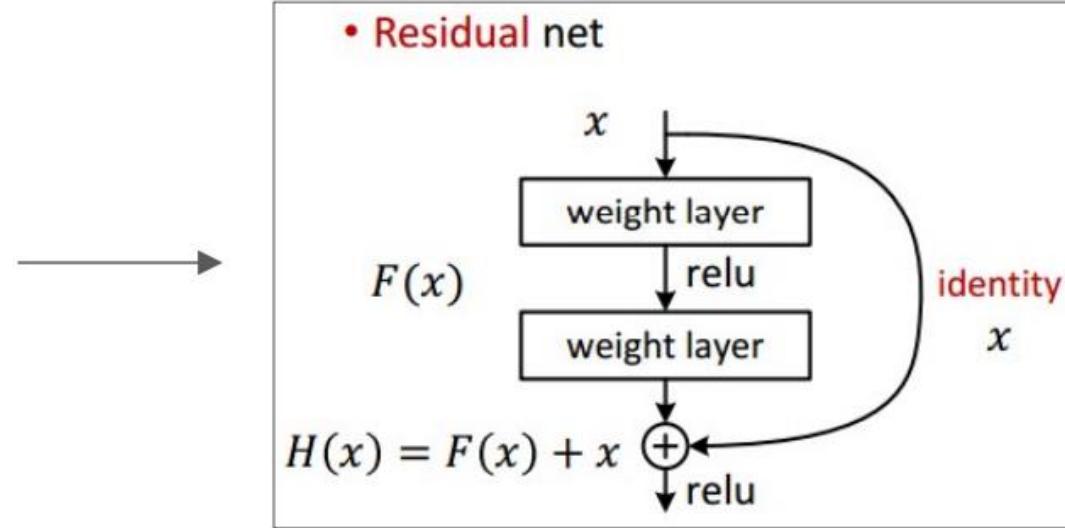
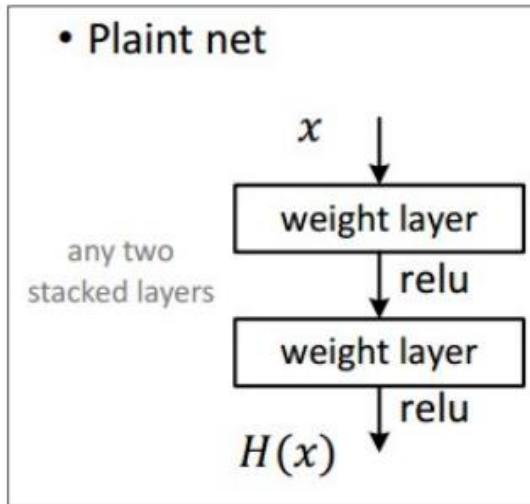


# What's the main idea?



Why is this a good idea?

# So why is this a good idea?



$$\frac{dH}{dx}$$

could be **big** or **small**  
not close to  $\mathbf{I}$

$$\frac{dH}{dx} = \frac{dF}{dx} + \mathbf{I}$$

If weights are not too big,  
this will be small(ish)

