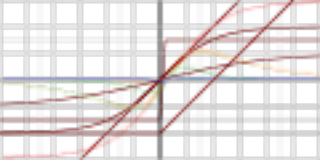


# Neural Network Architectures

Neil Gong

# Artificial Neural Networks

- Input/output
- Weight
- Activation function
- Connection pattern



# Activation function

Name	Plot	Function, $g(x)$	
Identity		$x$	
Binary step		$\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	
Logistic, sigmoid, or soft step		$\sigma(x) \doteq \frac{1}{1 + e^{-x}}$	
Hyperbolic tangent (tanh)		$\tanh(x) \doteq \frac{e^x - e^{-x}}{e^x + e^{-x}}$	
Rectified linear unit (ReLU)		$(x)^+ \doteq \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases} = \max(0, x) = x\mathbf{1}_{x>0}$	
Gaussian Error Linear Unit (GELU)		$\frac{1}{2}x \left( 1 + \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) \right) = x\Phi(x)$	
Leaky rectified linear unit (Leaky ReLU)		$\begin{cases} 0.01x & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$	

Source: Wikipedia

# Connection patterns

- Fully connected
- Softmax
- Convolution
- Residual
- transformer

# Convolution: a 2-D example

input

0	0	0	0	0	0	0	0
0	0	0	0	0	1	1	0
0	1	1	1	1	1	1	0
0	1	1	1	1	1	1	0
0	1	1	1	1	1	1	0
0	0	1	1	1	0	0	0
0	0	1	1	1	0	0	0
0	0	0	0	0	0	0	0

output

filter

1	2	1
0	0	0
-1	-2	-1


Credit: Kaiming He

# Convolution: a 2-D example

input

0	1	0	2	0	1	0	0	0	0
0	0	0	0	0	0	1	1	1	0
0	-1	1	-2	1	-1	1	1	1	0
0	1	1	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	0
0	0	1	1	1	1	0	0	0	0
0	0	1	1	1	1	0	0	0	0
0	0	0	0	0	0	0	0	0	0

output

-3						

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

# Convolution: a 2-D example

input

0	0	1	0	2	0	1	0	0	0	0
0	0	0	0	0	0	0	1	1	1	0
0	1	-1	1	-2	1	-1	1	1	1	0
0	1	1	1	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	1	0
0	0	1	1	1	1	0	0	0	0	0
0	0	1	1	1	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

output

-3	-4						

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

# Convolution: a 2-D example

input

0	0	0	1	0	2	0	1	0	0	0	0
0	0	0	0	0	0	0	0	1	1	0	0
0	1	1	-1	1	-2	1	-1	1	1	0	0
0	1	1	1	1	1	1	1	1	1	0	0
0	1	1	1	1	1	1	1	1	1	0	0
0	0	1	1	1	0	0	0	0	0	0	0
0	0	1	1	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

output

-3	-4	-4				

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

# Convolution: a 2-D example

input

0	0	0	0	1	0	1	0	0	0
0	0	0	0	0	0	1	0	1	0
0	1	1	1	-1	1	-2	1	-1	1
0	1	1	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	0
0	0	1	1	1	0	0	0	0	0
0	0	1	1	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

output

-3	-4	-4	-4			

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

# Convolution: a 2-D example

input

0	0	0	0	0	0	0	0
0	0	0	0	0	1	1	0
0	1	1	1	1	1	1	0
0	1	1	1	1	1	1	0
0	1	1	1	1	1	1	0
0	0	1	1	1	<b>01</b>	<b>02</b>	<b>01</b>
0	0	1	1	1	<b>00</b>	<b>00</b>	<b>00</b>
0	0	0	0	0	<b>-1</b>	<b>-2</b>	<b>-1</b>

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

output

-3	-4	-4	-4	-4	-3
-3	-4	-4	-3	-1	0
0	0	0	0	0	0
2	1	0	1	3	3
2	1	0	1	3	3
1	3	4	3	1	0

# Convolution: a 2-D example

$$y[n, m] = \sum_{i=-r}^r \sum_{j=-r}^r w[i, j] x[n + i, m + j]$$

output map

filter weights

input map

coordinates in a local window

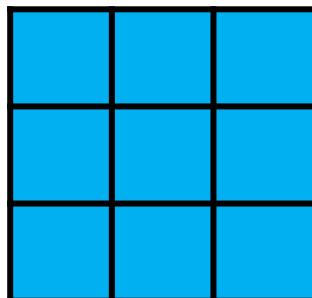
$r$ : kernel radius  
kernel size =  $2r + 1$

\* In ConvNets, convolution is often implemented as **cross-correlation** (no flipping)

# Convolution: padding

input:  $8 \times 8$ , + pad

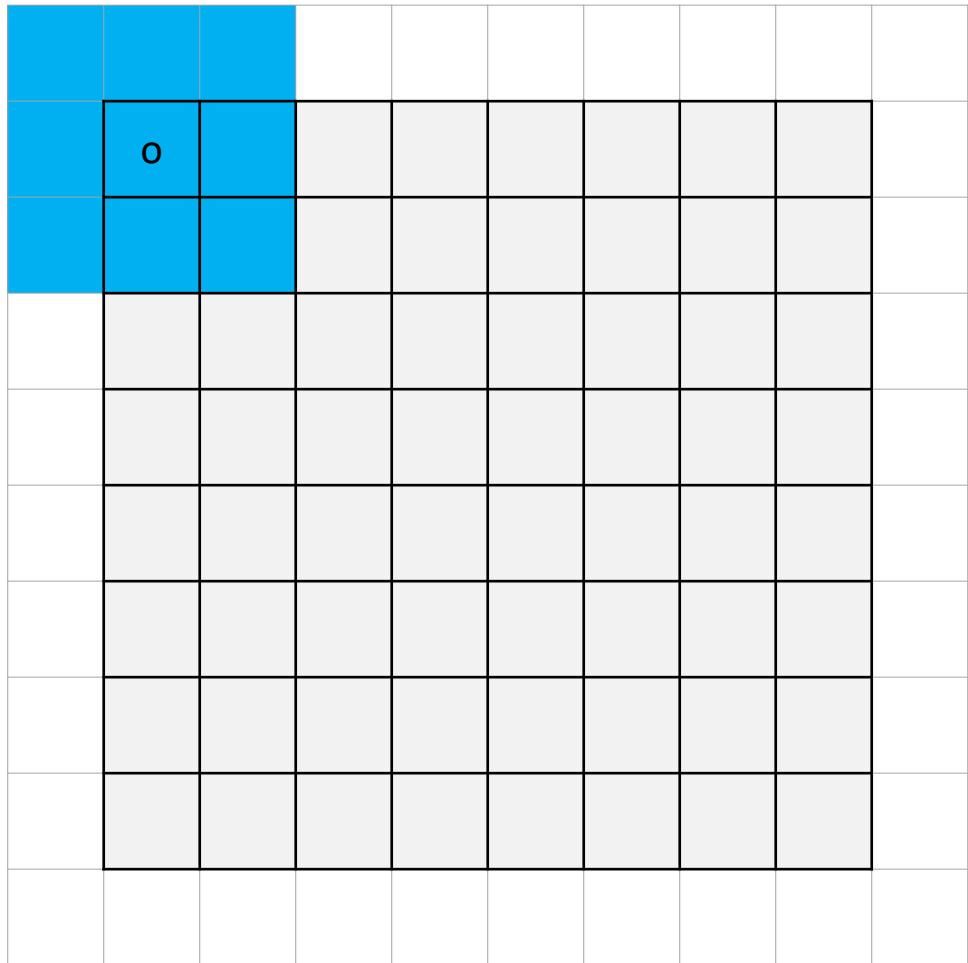
# filter



output:  $H \times W = 8 \times 8$

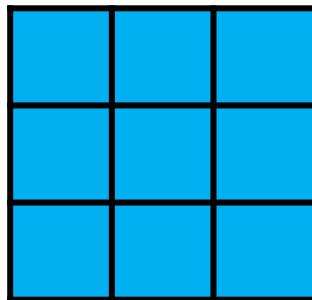
# Convolution: stride

input

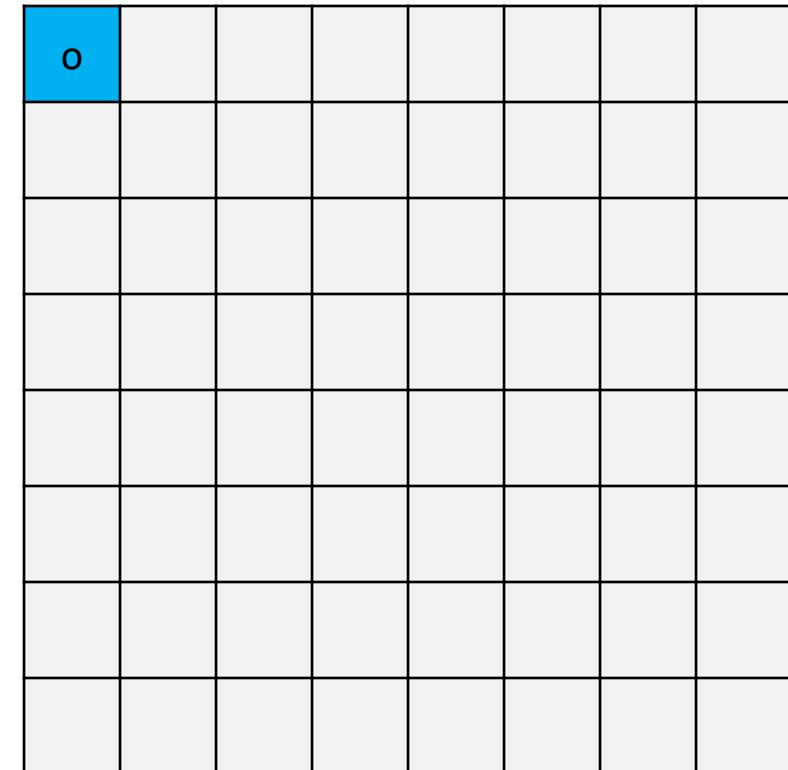


**stride = 2**

# filter

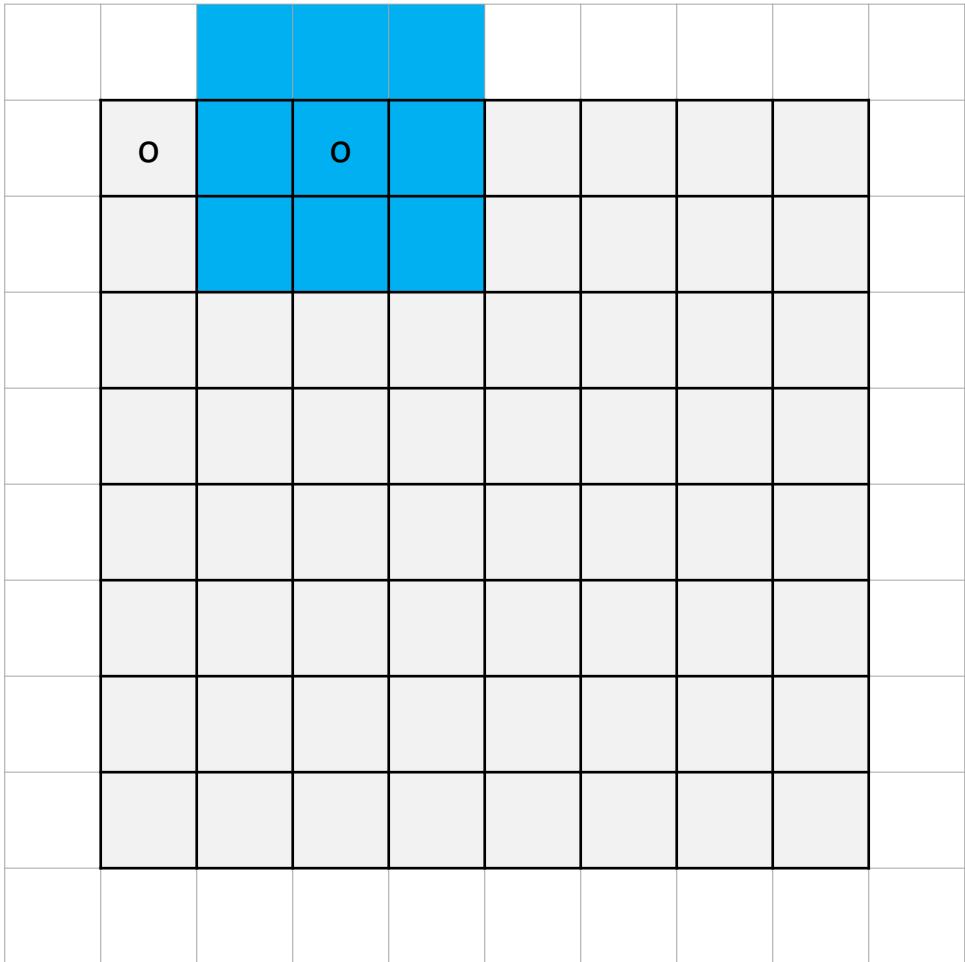


## output



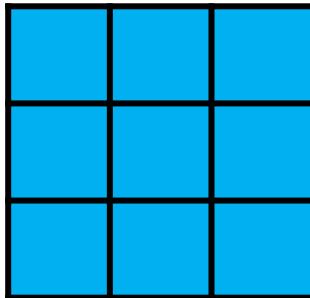
# Convolution: stride

input

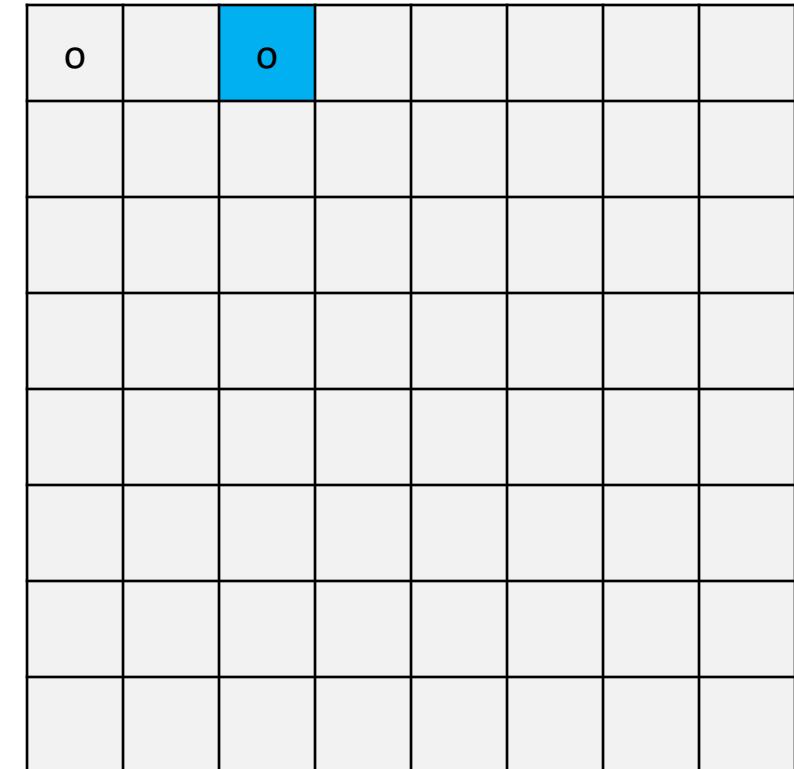


**stride = 2**

# filter



## output



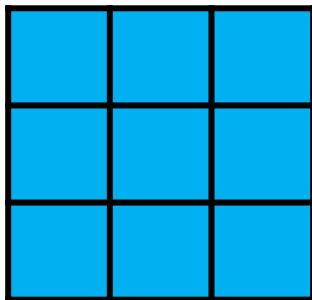
# Convolution: stride

input:  $H \times W = 8 \times 8$

o		o		o		o		
o		o		o		o		
o		o		o		o		
o		o		o		o		

stride = 2

filter



output:  $H \times W = 4 \times 4$

o		o		o		o	
o		o		o		o	
o		o		o		o	
o		o		o		o	

# Convolution: stride

input:  $H \times W = 8 \times 8$

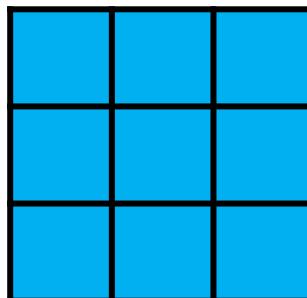
o		o		o		o		
o		o		o		o		
o		o		o		o		
o		o		o		o		

stride = 2

- reduces feature map size
- compress and abstract

output:  $H \times W = 4 \times 4$

filter

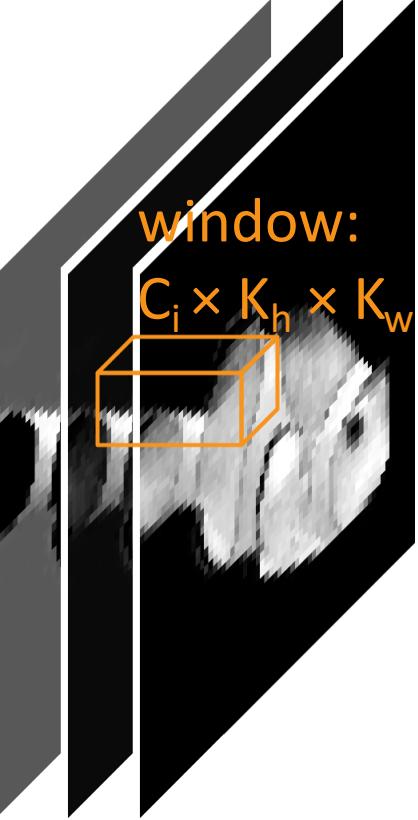


o	o	o	o
o	o	o	o
o	o	o	o
o	o	o	o

$$H_{out} = \lfloor (H_{in} + 2\text{pad}_h - K_h) / \text{str} \rfloor + 1$$

\*rounding operation depends on libraries

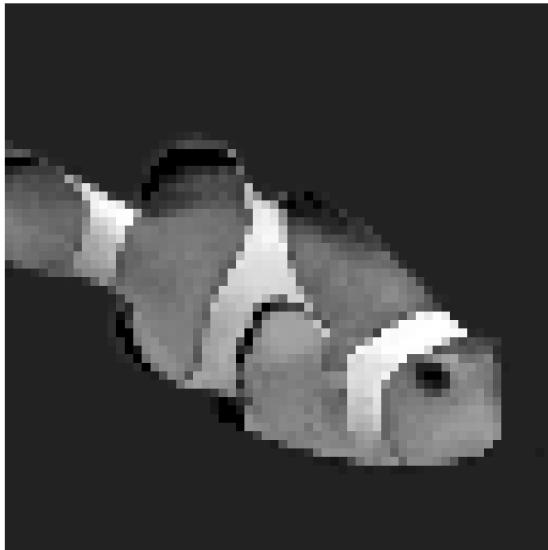
# Convolution: Multi-channel inputs



$$\text{window: } C_i \times K_h \times K_w \quad * \quad \text{filter: } C_i \times K_h \times K_w =$$
A 3D cube representing a filter kernel. It has a front face with a 3x3 grid of smaller squares, representing the spatial dimensions of the filter. The filter is multiplied by the window to produce the output feature map.



# Convolution: Multi-channel outputs



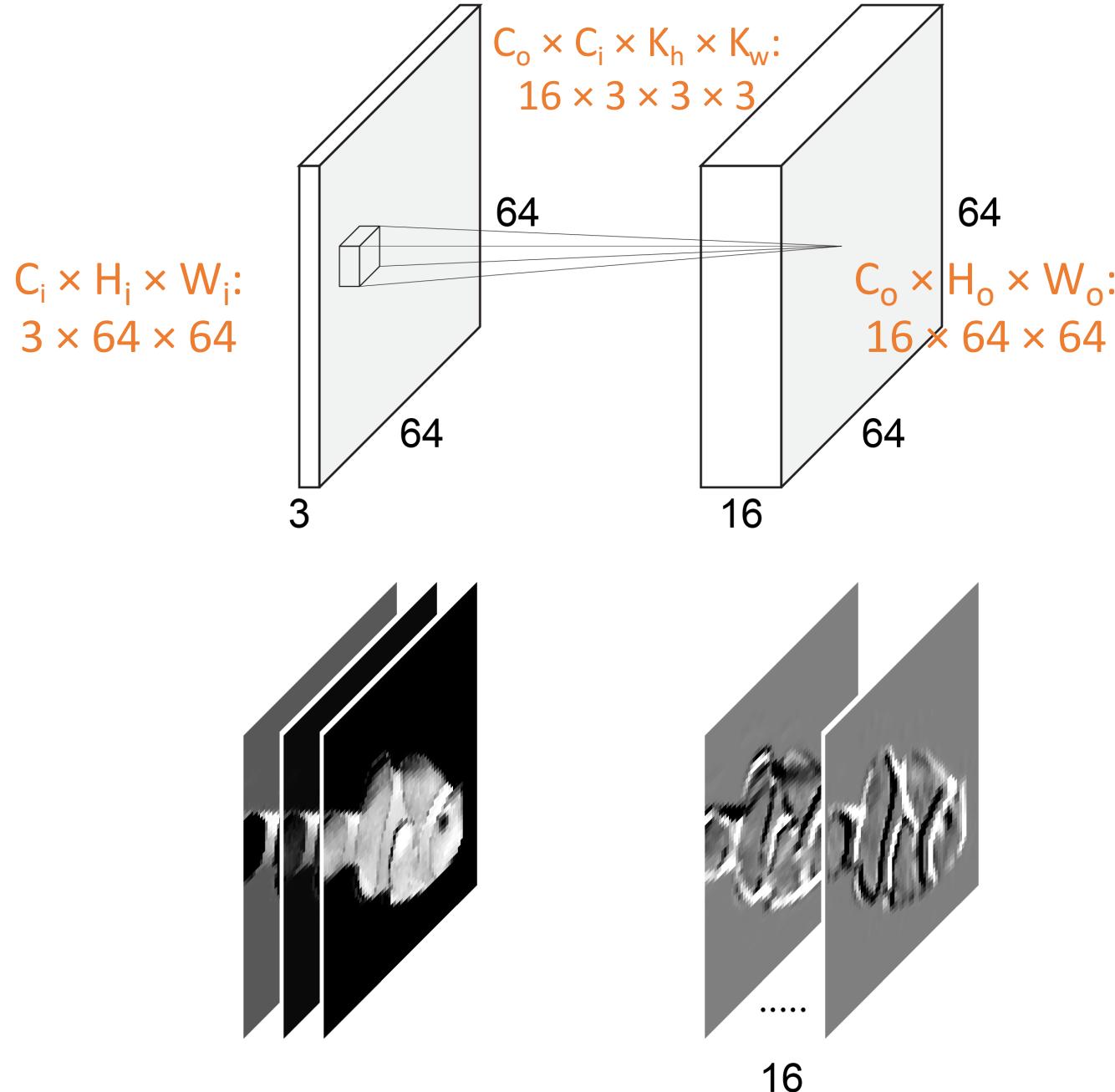
$$\begin{matrix} * & \begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 0 & 0 & 0 \\ \hline -1 & -2 & -1 \\ \hline \end{array} & = \end{matrix}$$

one filter, one feature

$$\begin{matrix} * & \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array} & = \end{matrix}$$



# Convolution: tensor views

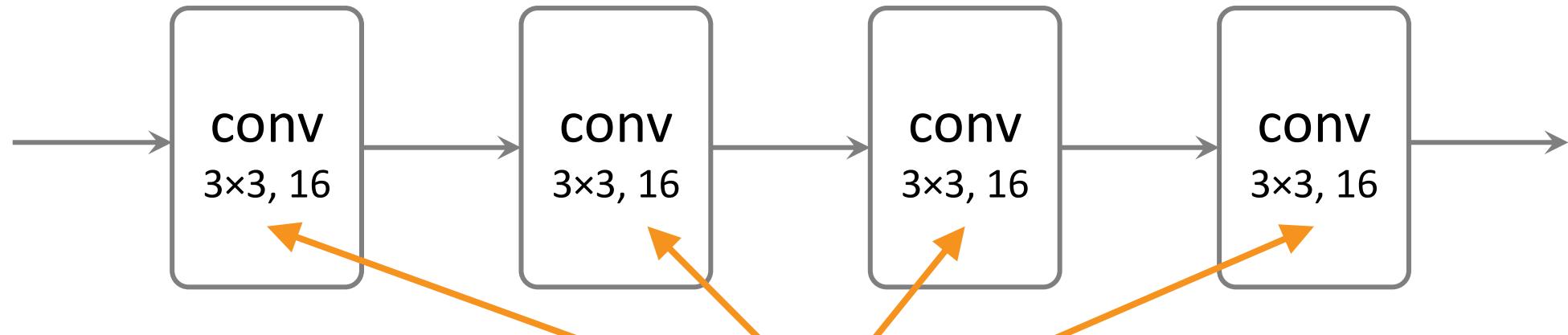
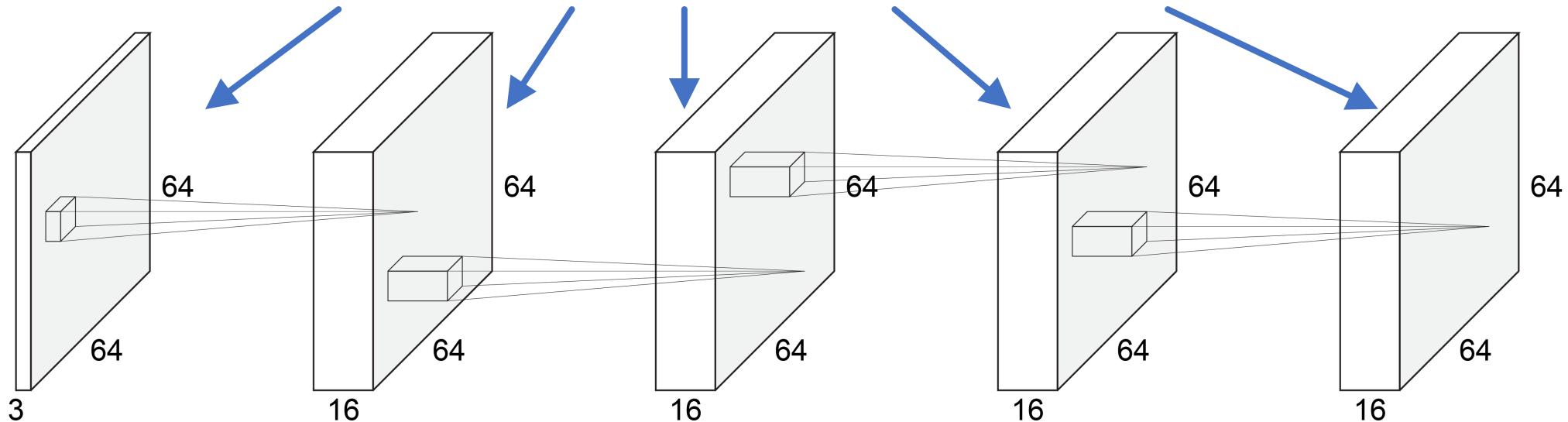


- Tensor: high-dimension array
- feature maps
  - 3-D tensor:  $C \times H \times W$
  - $C$ : channels
  - $H$ : height
  - $W$ : width
- filters
  - 4-D tensor:  $C_o \times C_i \times K_h \times K_w$
  - $C_o$ : output channels
  - $C_i$ : input channels
  - $K_h, K_w$ : filter height, width

two ways of showing  
neural nets

# Composing basic operations

these are activations (features, embeddings, tensors ...)



these are operations (functions, transforms, layers ...)

# Deep Residual Learning

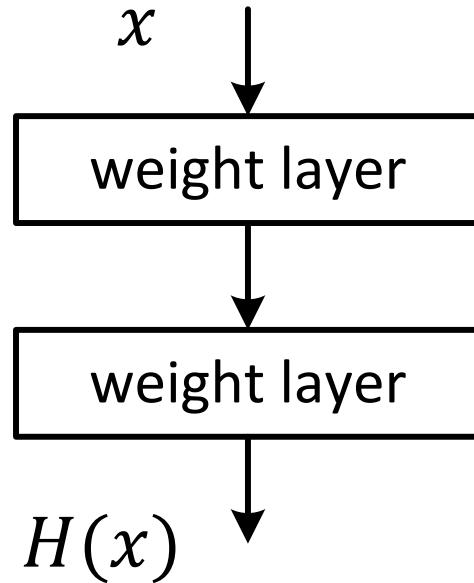
- Deep Learning gets way deeper
- simple component: identity shortcut
- enable networks w/ hundreds of layers

**Compose simple modules into complex functions**



# Deep Residual Learning

a subnet in  
a deep net

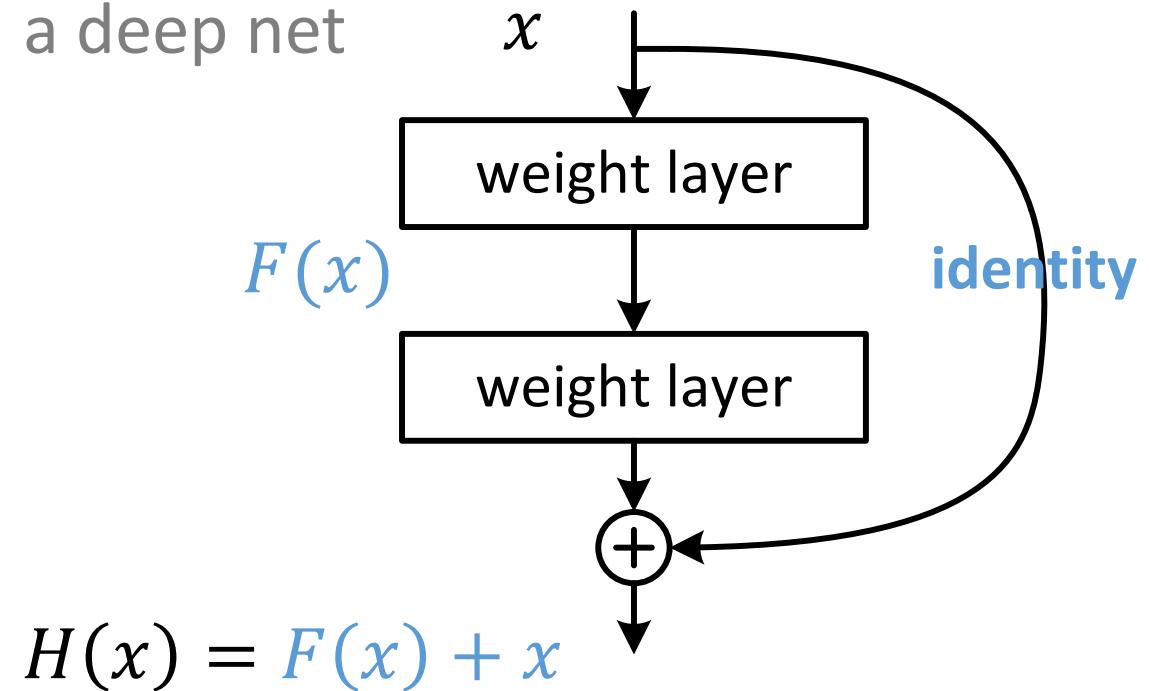


classical network

- $H(x)$ : desired function to be fit by a subnet
- let weight layers fit  $H(x)$

# Deep Residual Learning

a subnet in  
a deep net

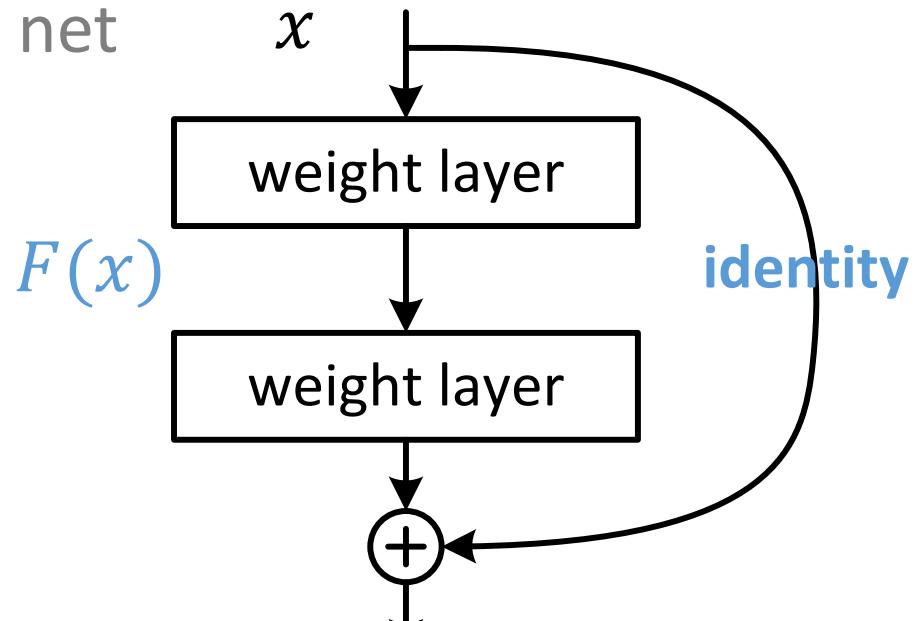


## residual block

- $H(x)$ : desired function to be fit by a subnet
- ~~let weight layers fit  $H(x)$~~
- let weight layers fit  $F(x)$
- set  $H(x) = F(x) + x$

# Deep Residual Learning

a subnet in  
a deep net



$$H(x) = F(x) + x$$

## residual block

- $F(x)$ : residual function
- if  $H(x) = \text{identity}$  is near-optimal
  - push weights to small
  - encourage small changes
- initialization
  - small or zero weights

# Residual Networks (ResNet)

# **Building very deep nets:**

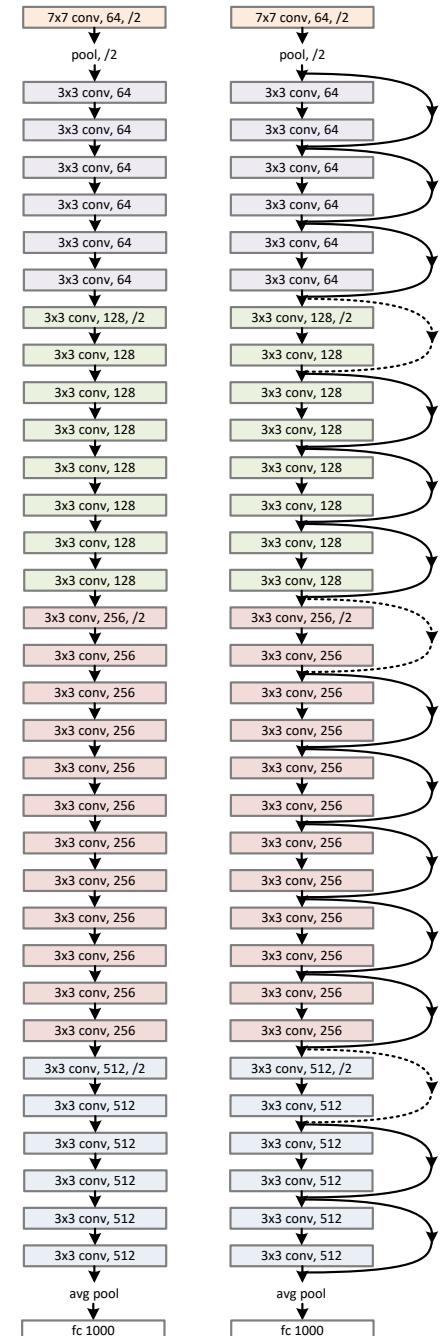
- add **identity connections** to vanilla nets  
(a.k.a. skip/shortcut/residual connections)

**or:**

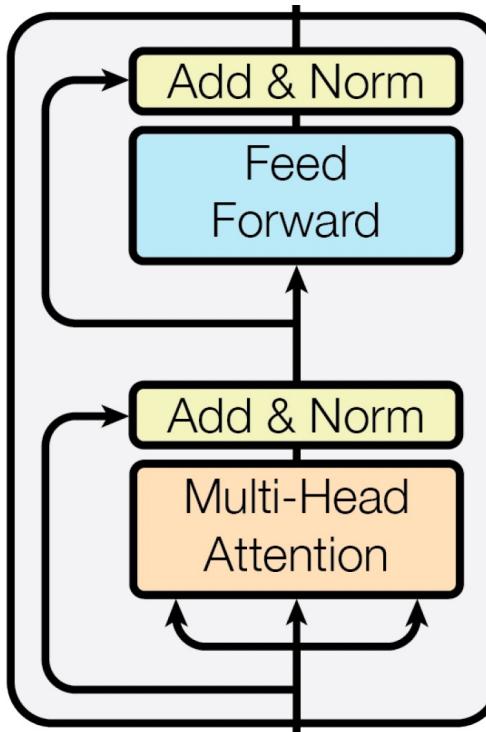
- stack many residual blocks

# Residual Blocks:

- new generic modules for neural nets
  - design blocks and compose them

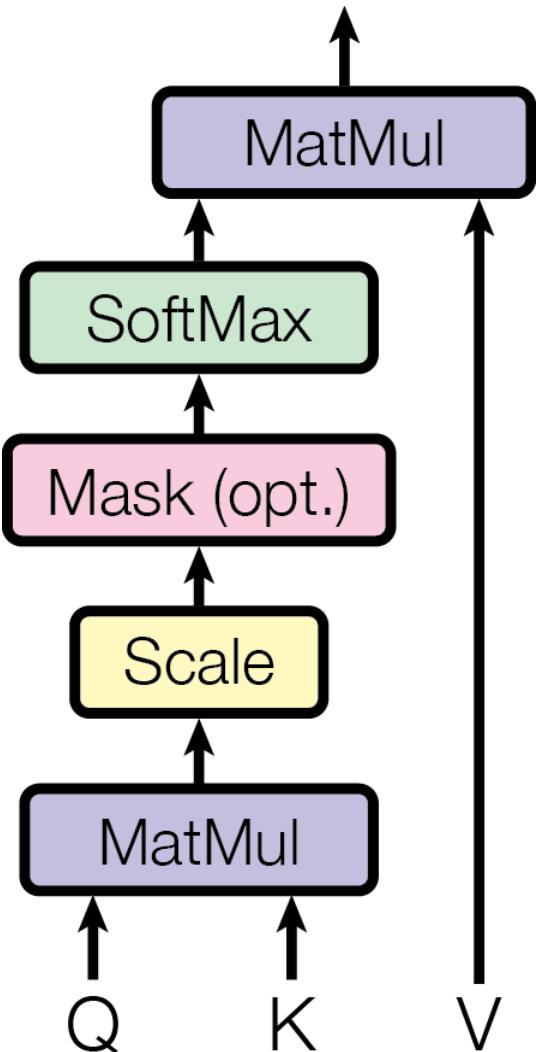


# Residual Block: Transformer



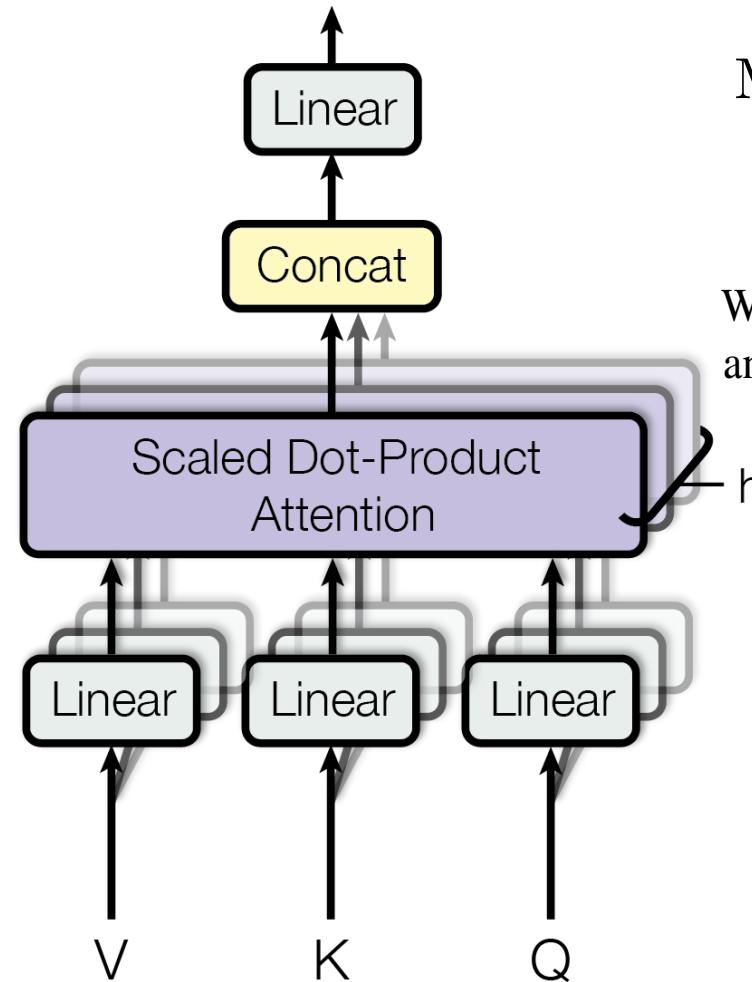
A Transformer Block has two Residual Blocks.

# Scaled Dot-Product Attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

# Multi-Head Attention



$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

Where the projections are parameter matrices  $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$  and  $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$ .

# Position-wise feed-forward network

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

# One last detail: layer normalization

**Main idea:** batch normalization is very helpful, but hard to use with sequence models

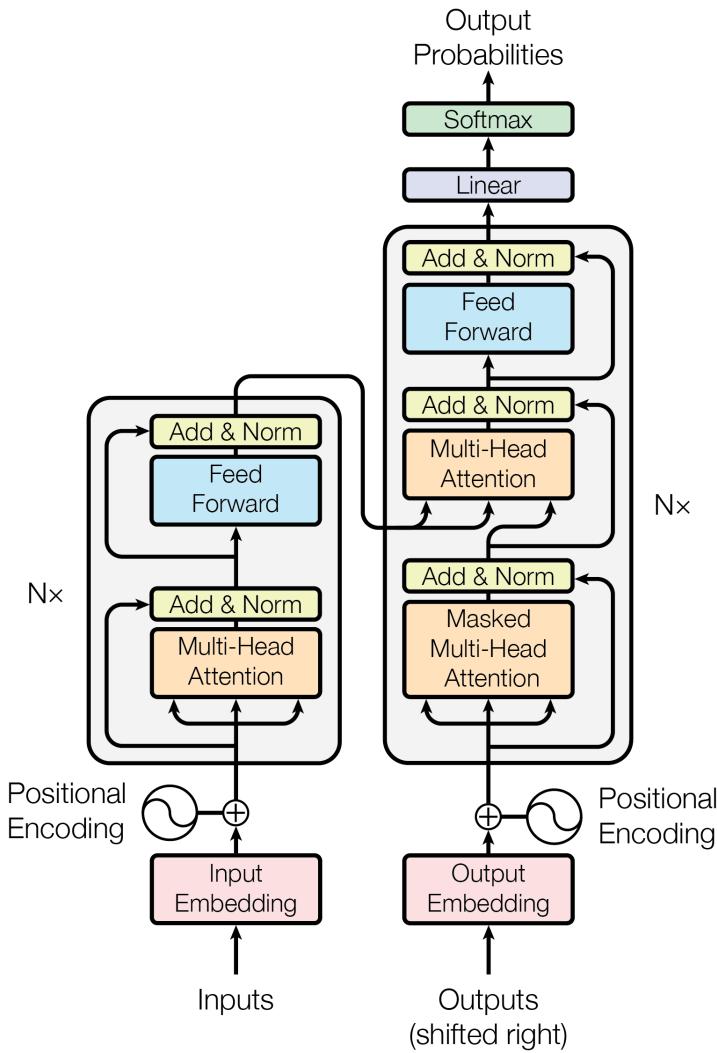
Sequences are different lengths, makes normalizing across the batch hard

Sequences can be very long, so we sometimes have small batches

**Simple solution:** “layer normalization” – like batch norm, but not across the batch

<p><b>Batch norm</b></p> <p><math>d\text{-dim}</math></p> $\mu = \frac{1}{B} \sum_{i=1}^B a_i$ $\sigma = \sqrt{\frac{1}{B} \sum_{i=1}^B (a_i - \mu)^2}$ $\bar{a}_i = \frac{a_i - \mu}{\sigma} \gamma + \beta$	<p><math>d\text{-dimensional vectors}</math></p> <p>for each sample in batch</p> <p><b>Layer norm</b></p> <p><math>a</math></p> <p><math>1\text{-dim}</math></p> $\mu = \frac{1}{d} \sum_{j=1}^d a_j$ $\sigma = \sqrt{\frac{1}{d} \sum_{j=1}^d (a_j - \mu)^2}$ $\bar{a} = \frac{a - \mu}{\sigma} \gamma + \beta$
---	--

# Transformer architecture



# Positional encoding: sin/cos

**Naïve positional encoding:** just append  $t$  to the input

$$\bar{x}_t = \begin{bmatrix} x_t \\ t \end{bmatrix}$$

This is not a great idea, because **absolute** position is less important than **relative** position

I walk my dog every day



every single day I walk my dog



The fact that “my dog” is right after “I walk” is the important part, not its absolute position

we want to represent **position** in a way that tokens with similar **relative** position have similar **positional encoding**

**Idea:** what if we use **frequency-based** representations?

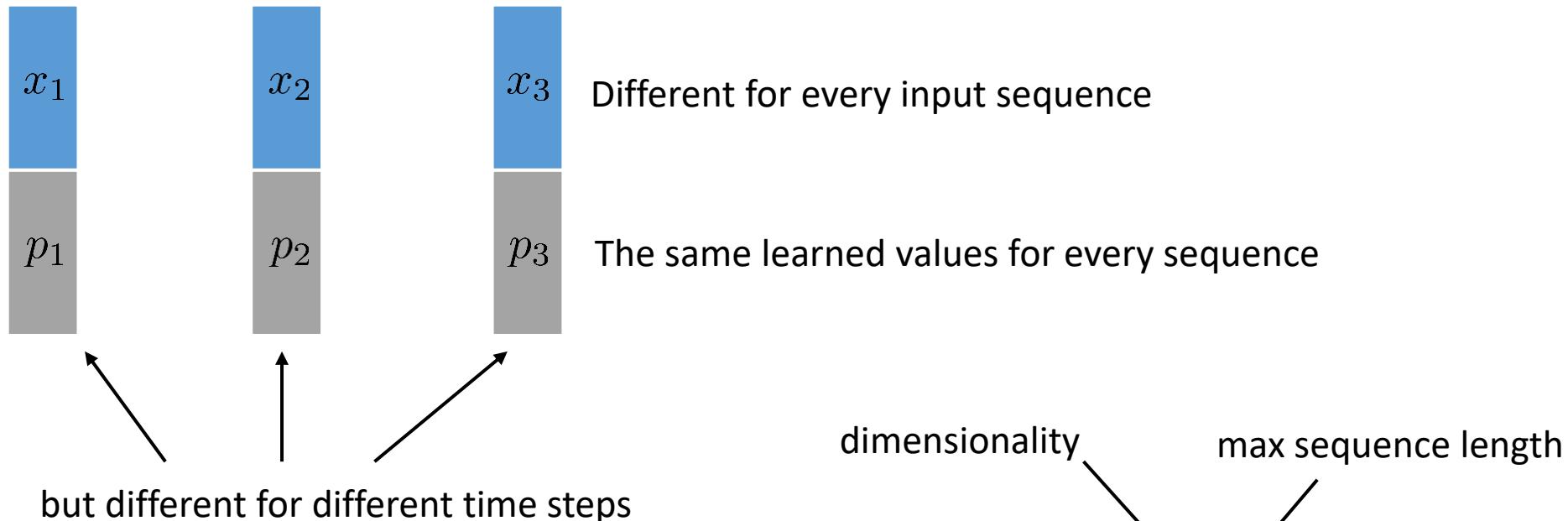
$$p_t = \begin{bmatrix} \sin(t/10000^{2*1/d}) \\ \cos(t/10000^{2*1/d}) \\ \sin(t/10000^{2*2/d}) \\ \cos(t/10000^{2*2/d}) \\ \dots \\ \sin(t/10000^{2*\frac{d}{2}/d}) \\ \cos(t/10000^{2*\frac{d}{2}/d}) \end{bmatrix}$$

dimensionality of positional encoding



# Positional encoding: learned

Another idea: just learn a positional encoding



How many values do we need to learn?

$$P = [p_1, p_2, \dots, p_T] \in R^{d \times T}$$

+ more flexible (and perhaps more optimal) than sin/cos encoding

+ a bit more complex, need to pick a max sequence length (and can't generalize beyond it)

# Vision Transformer (ViT)

