

위성 격자 자료와 합성곱 신경망(CNN)을 통한 공간 정보 학습

차세대수치예보모델개발사업단

전현주

2026.2.23

실습자료





Hyeon-Ju Jeon

Research Scientist at KIAPS

Contact Information

- hjeon@kiaps.org, hgd963@gmail.com
- <https://hgd963.github.io/>

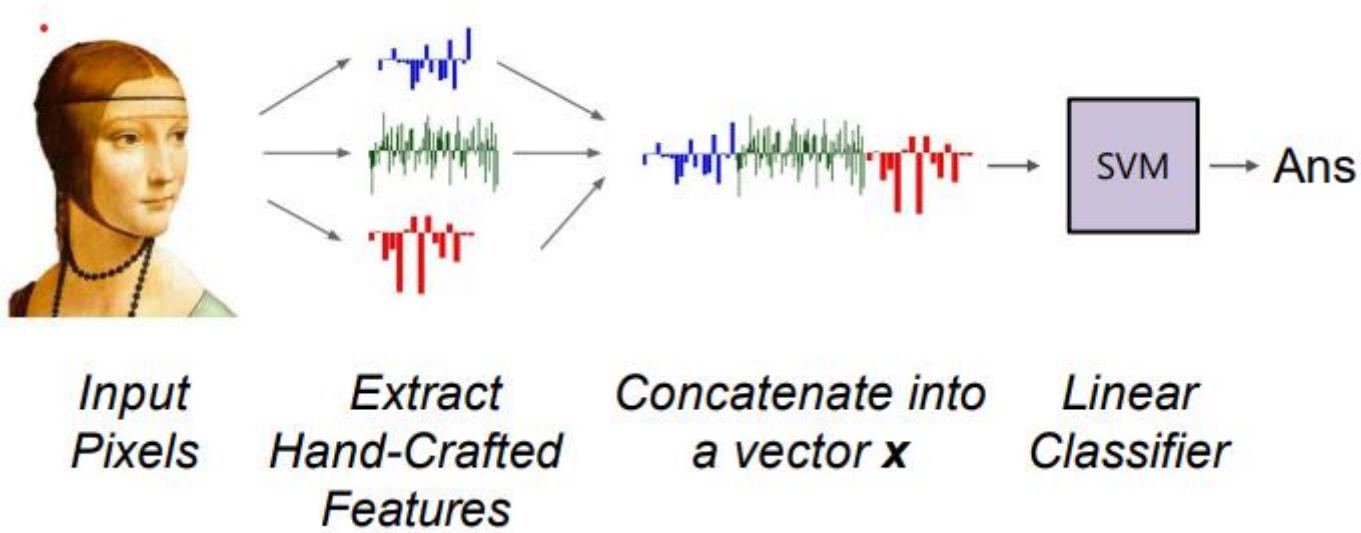
Research Interest

- Multimodal and Irregular Real-world Data Mining
- Applied Machine Learning
- Spatiotemporal Graph Neural Networks
- Application domains: Weather prediction, Disease prediction, Biomedical, Bibliographic network analysis, Anomaly detection, etc.

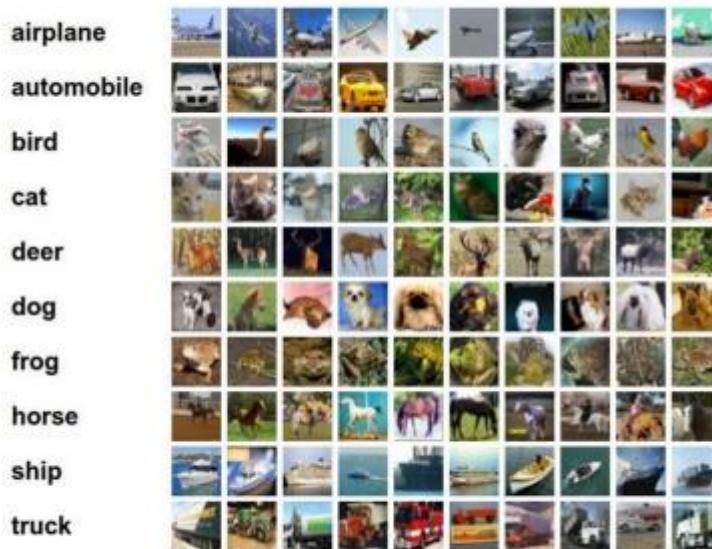
Professional Experience

- Ph.D., Chung-Ang University (Leave of absence)
- Research Scientist, KIAPS (2021.10 – Present)
- Lecturer, LG CNS (2021 – 2022)

Life Before Deep Learning



Why use features? Why not pixels?

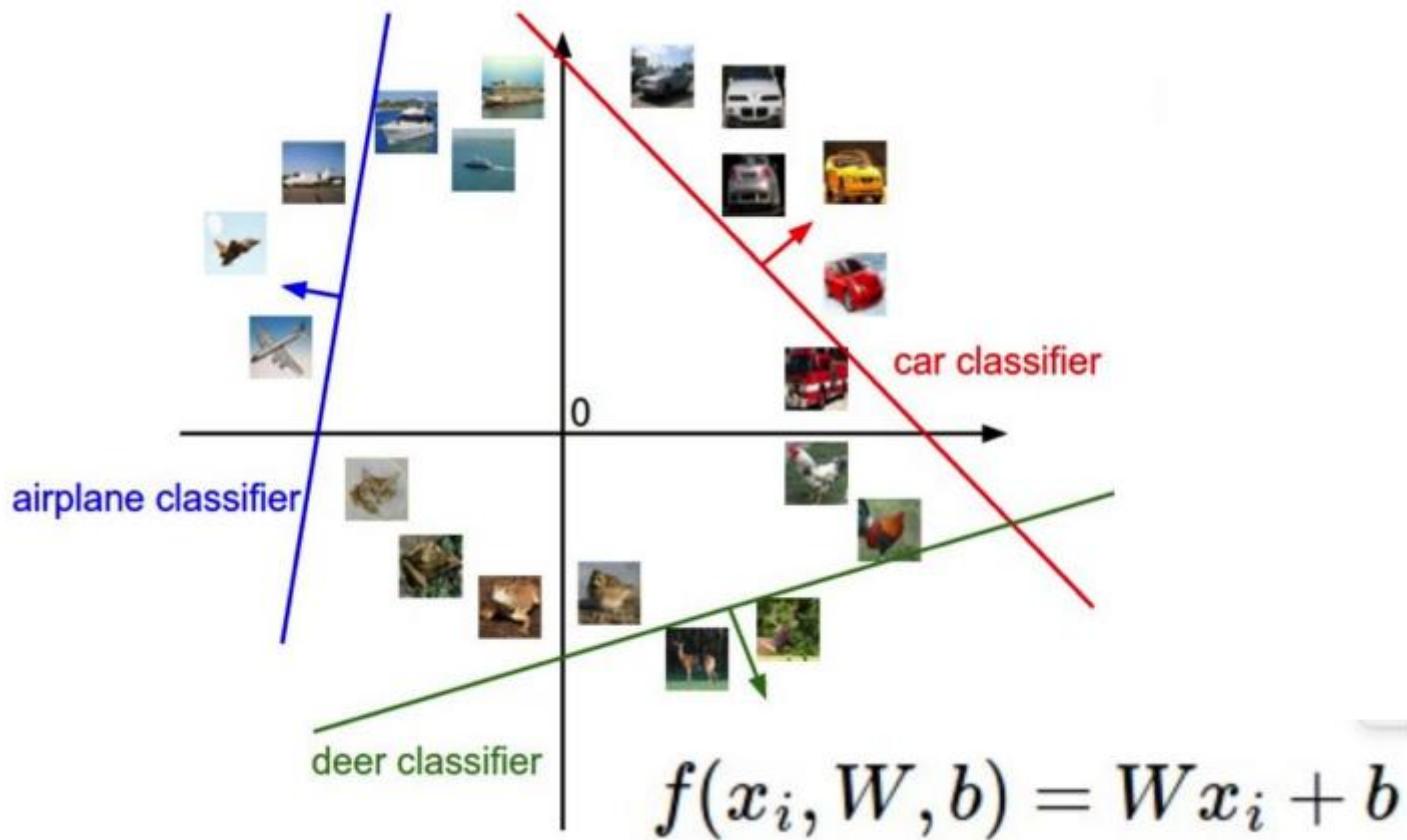


$$f(x_i, W, b) = Wx_i + b$$

Q: What would be a very hard set of classes for a linear classifier to distinguish?

(assuming x = pixels)

Linearly separable classes



Aside: Image Features

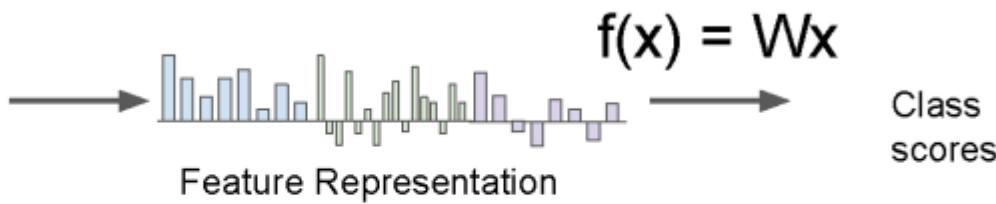
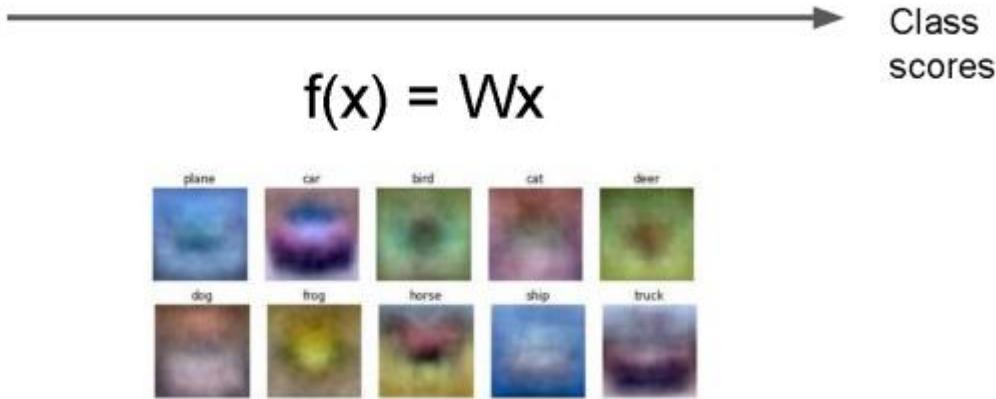
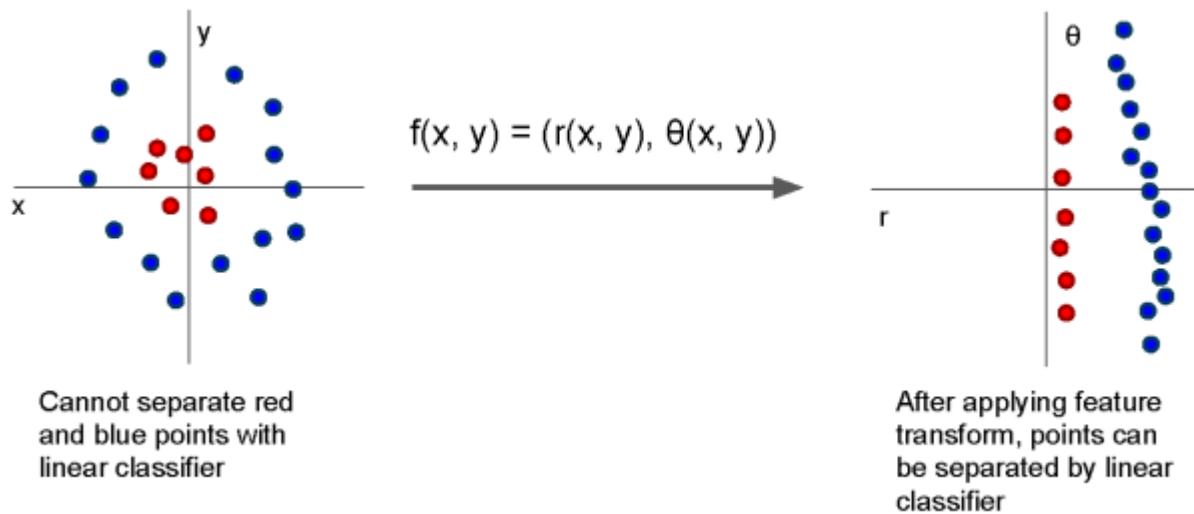


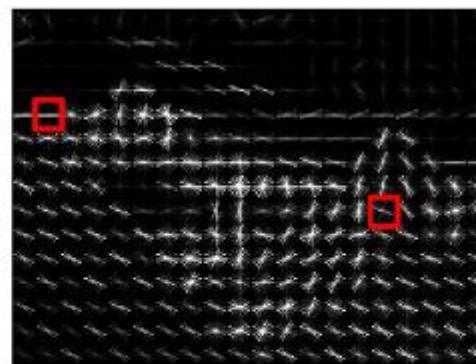
Image Features: Motivation



Example: Histogram of Oriented Gradients (HoG)

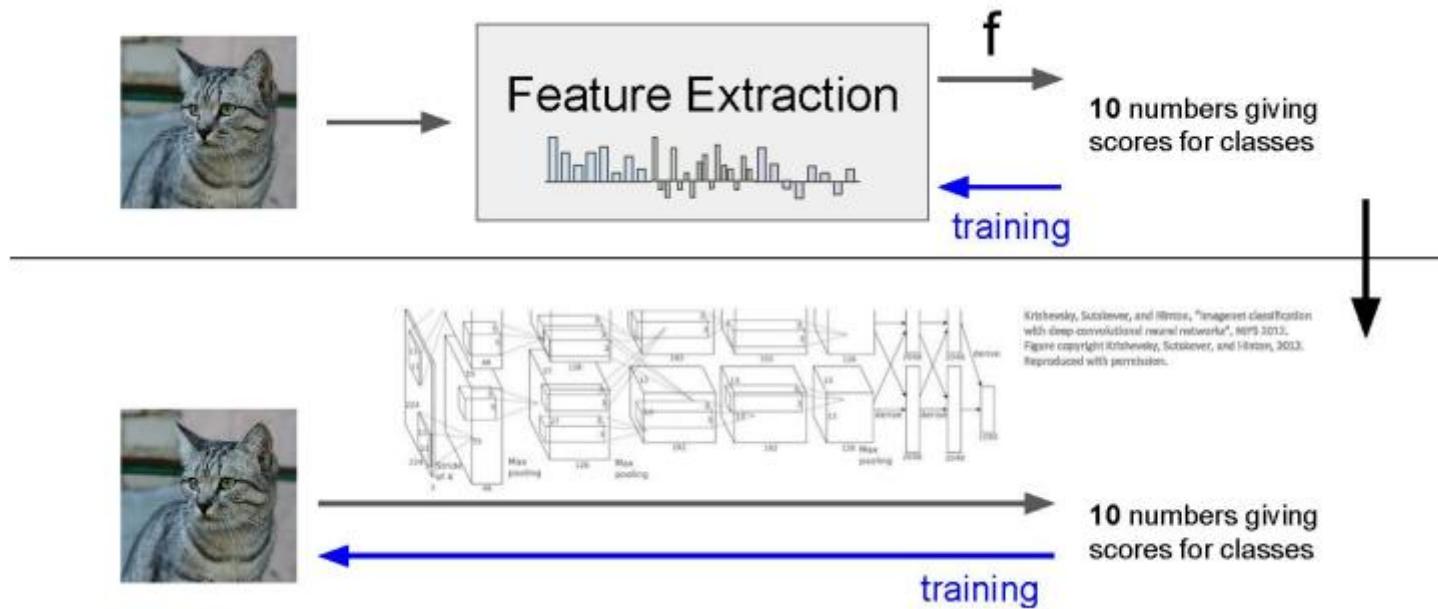


Divide image into 8x8 pixel regions
Within each region quantize edge direction into 9 bins

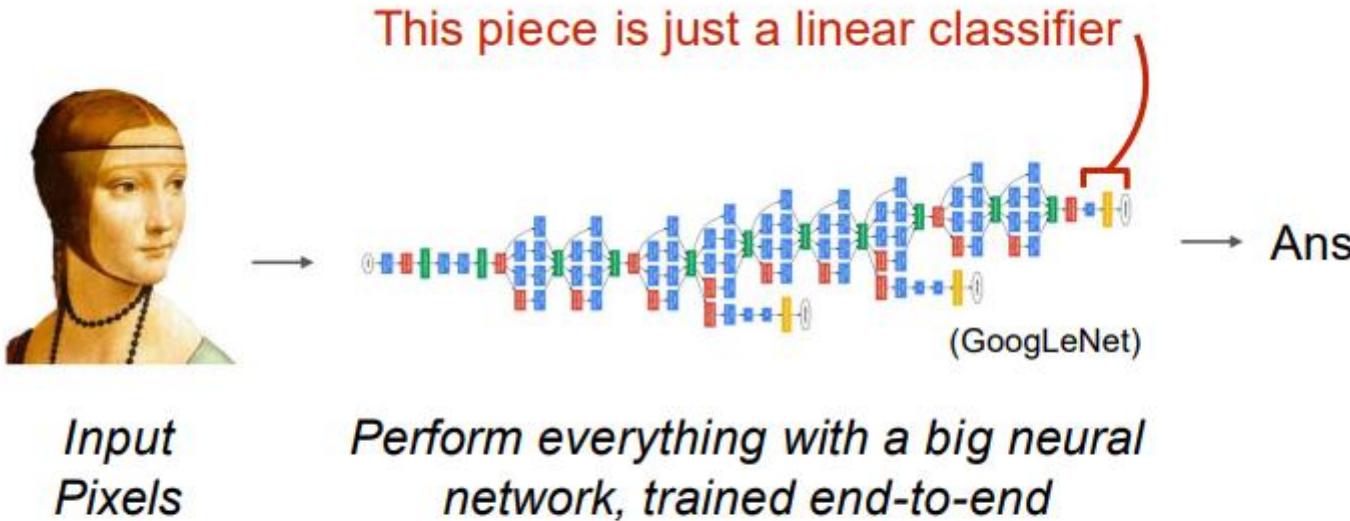


Example: 320x240 image gets divided into 40x30 bins; in each bin there are 9 numbers so feature vector has $30 \times 40 \times 9 = 10,800$ numbers

Image features vs ConvNets

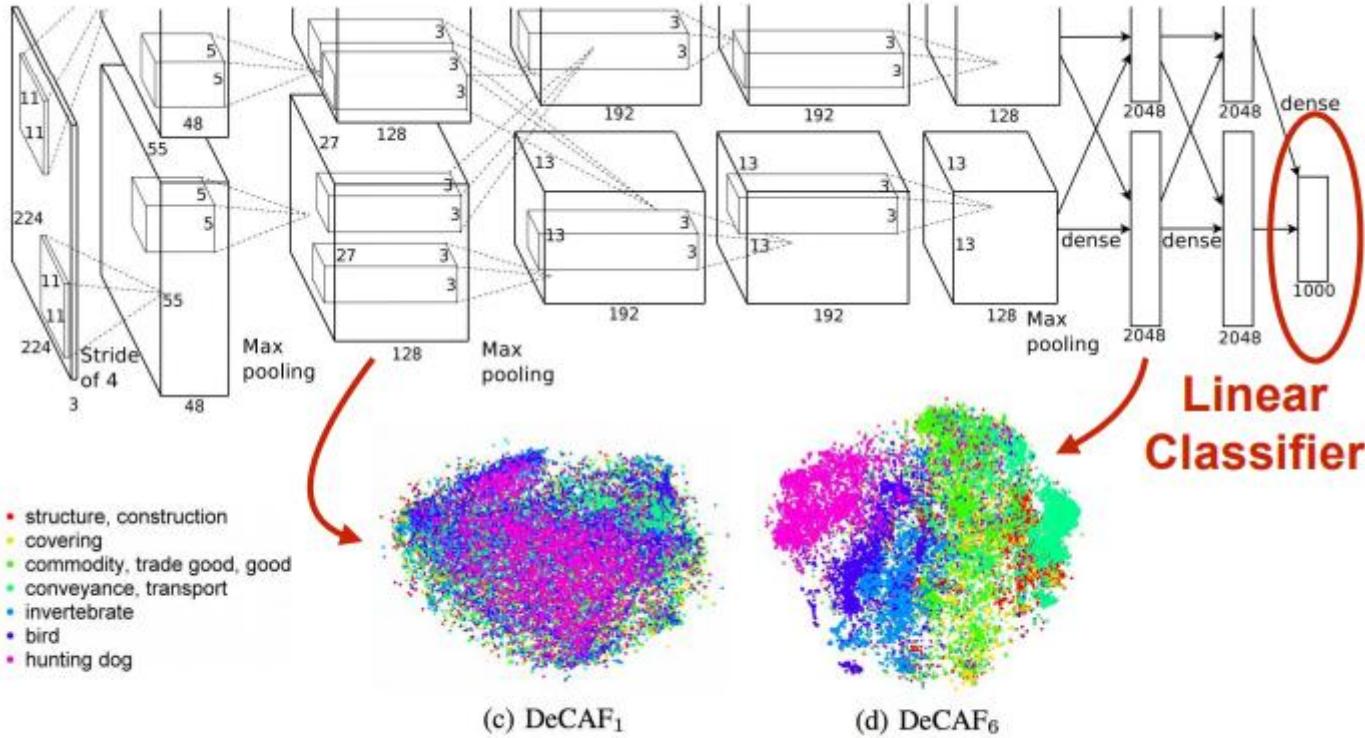


Last layer of most CNNs is a linear classifier



Key: perform enough processing so that by the time you get to the end of the network, the classes are linearly separable

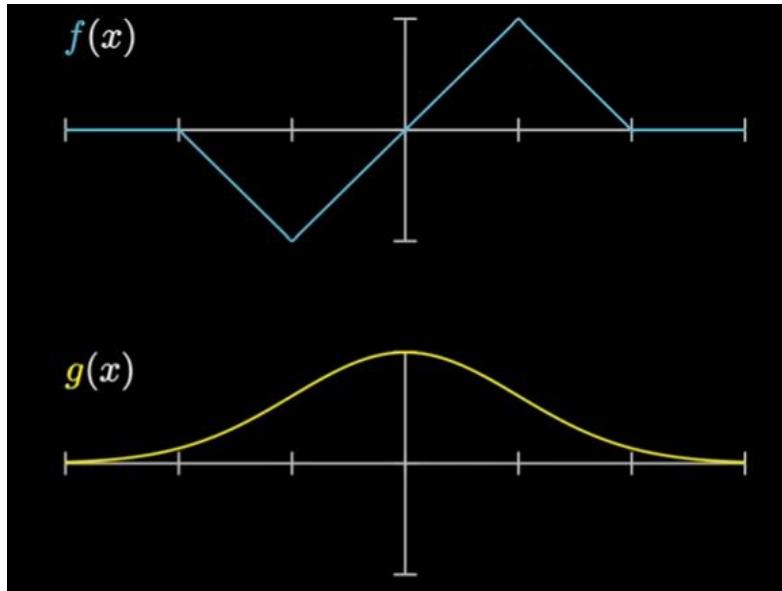
Visualizing AlexNet in 2D with t-SNE



(2D visualization using t-SNE)

[Donahue, "DeCAF: DeCAF: A Deep Convolutional ...", arXiv 2013]

Convolutional operator



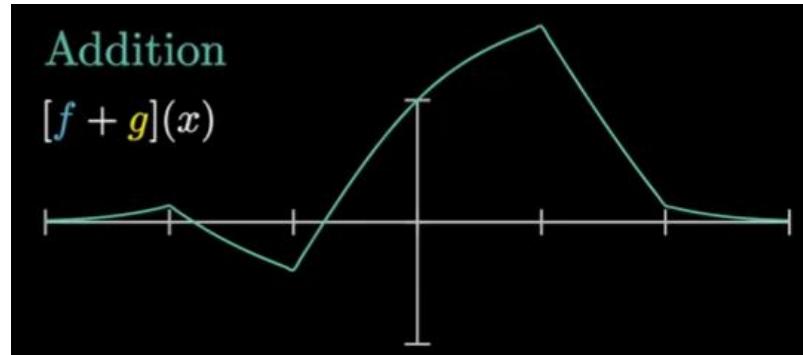
$$a = [1, 2, 3, 4]$$

$$b = [5, 6, 7, 8]$$

$$a+b = [6, 8, 10, 12]$$

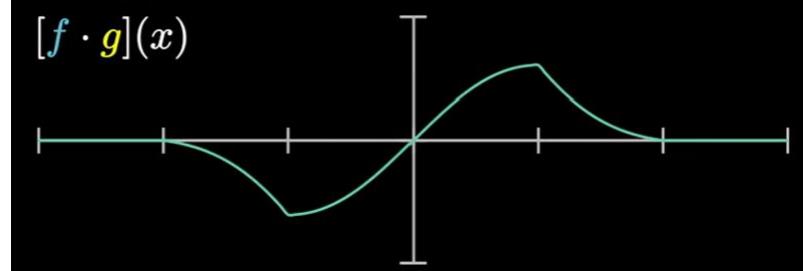
$$a \cdot b = [5, 12, 21, 32]$$

$$a * b = [5, 16, 34, 60, 61, 52, 32]$$



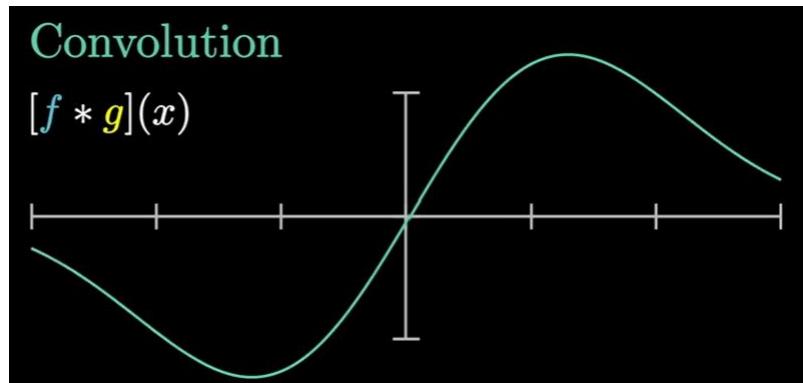
Addition

$$[f + g](x)$$



Multiplication

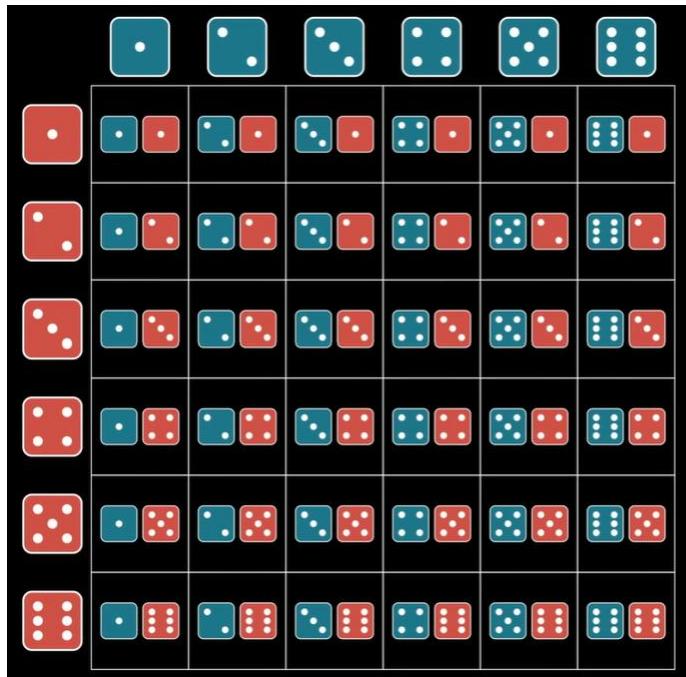
$$[f \cdot g](x)$$



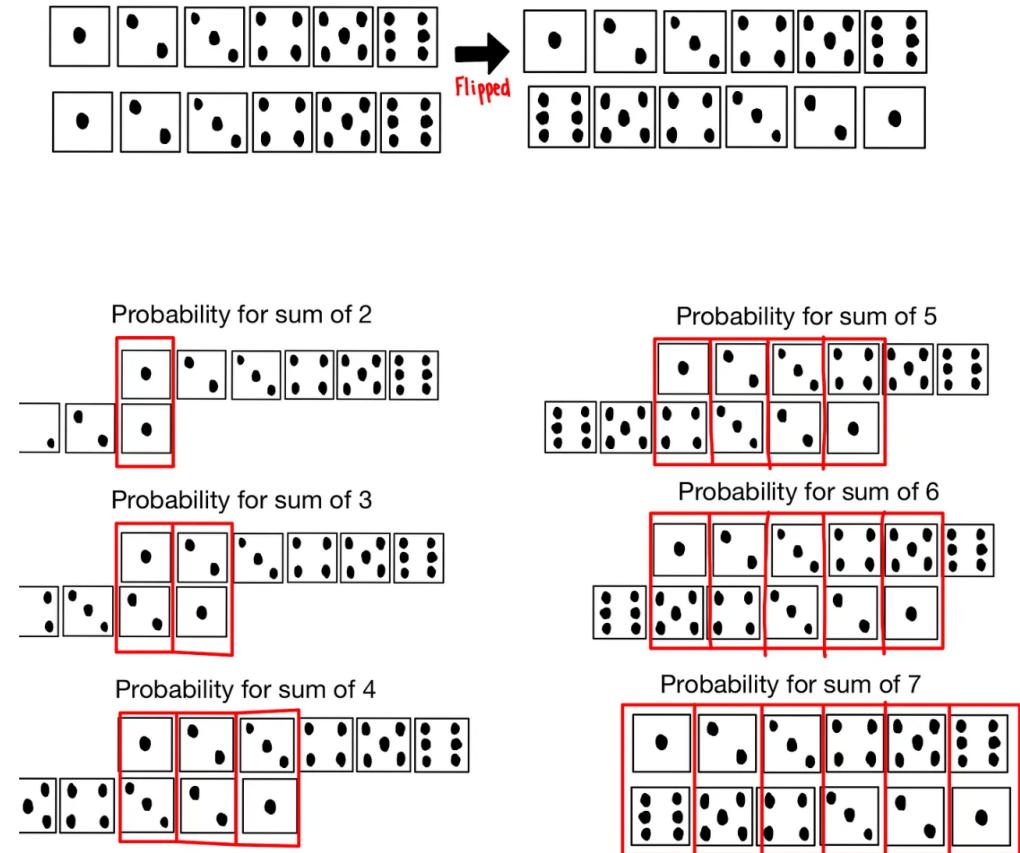
Convolution

$$[f * g](x)$$

Convolution with dice



$$6^2 = 36 \text{ Combinations}$$



Convolution of a_i and b_i

$$(a * b)_n = \sum_{\substack{i,j \\ i+j=n}} a_i \cdot b_j$$

	a_1	a_2	a_3	a_4	a_5	a_6
b_1	$a_1 \cdot b_1$	$a_2 \cdot b_1$	$a_3 \cdot b_1$	$a_4 \cdot b_1$	$a_5 \cdot b_1$	$a_6 \cdot b_1$
b_2	$a_1 \cdot b_2$	$a_2 \cdot b_2$	$a_3 \cdot b_2$	$a_4 \cdot b_2$	$a_5 \cdot b_2$	$a_6 \cdot b_2$
b_3	$a_1 \cdot b_3$	$a_2 \cdot b_3$	$a_3 \cdot b_3$	$a_4 \cdot b_3$	$a_5 \cdot b_3$	$a_6 \cdot b_3$
b_4	$a_1 \cdot b_4$	$a_2 \cdot b_4$	$a_3 \cdot b_4$	$a_4 \cdot b_4$	$a_5 \cdot b_4$	$a_6 \cdot b_4$
b_5	$a_1 \cdot b_5$	$a_2 \cdot b_5$	$a_3 \cdot b_5$	$a_4 \cdot b_5$	$a_5 \cdot b_5$	$a_6 \cdot b_5$
b_6	$a_1 \cdot b_6$	$a_2 \cdot b_6$	$a_3 \cdot b_6$	$a_4 \cdot b_6$	$a_5 \cdot b_6$	$a_6 \cdot b_6$

$$P(\square + \blacksquare = 2) = a_1 \cdot b_1$$

$$P(\square + \blacksquare = 3) = a_1 \cdot b_2 + a_2 \cdot b_1$$

$$P(\square + \blacksquare = 4) = a_1 \cdot b_3 + a_2 \cdot b_2 + a_3 \cdot b_1$$

$$P(\square + \blacksquare = 5) = a_1 \cdot b_4 + a_2 \cdot b_3 + a_3 \cdot b_2 + a_4 \cdot b_1$$

$$P(\square + \blacksquare = 6) = a_1 \cdot b_5 + a_2 \cdot b_4 + a_3 \cdot b_3 + a_4 \cdot b_2 + a_5 \cdot b_1$$

$$P(\square + \blacksquare = 7) = a_1 \cdot b_6 + a_2 \cdot b_5 + a_3 \cdot b_4 + a_4 \cdot b_3 + a_5 \cdot b_2 + a_6 \cdot b_1$$

$$P(\square + \blacksquare = 8) = a_2 \cdot b_6 + a_3 \cdot b_5 + a_4 \cdot b_4 + a_5 \cdot b_3 + a_6 \cdot b_2$$

$$P(\square + \blacksquare = 9) = a_3 \cdot b_6 + a_4 \cdot b_5 + a_5 \cdot b_4 + a_6 \cdot b_3$$

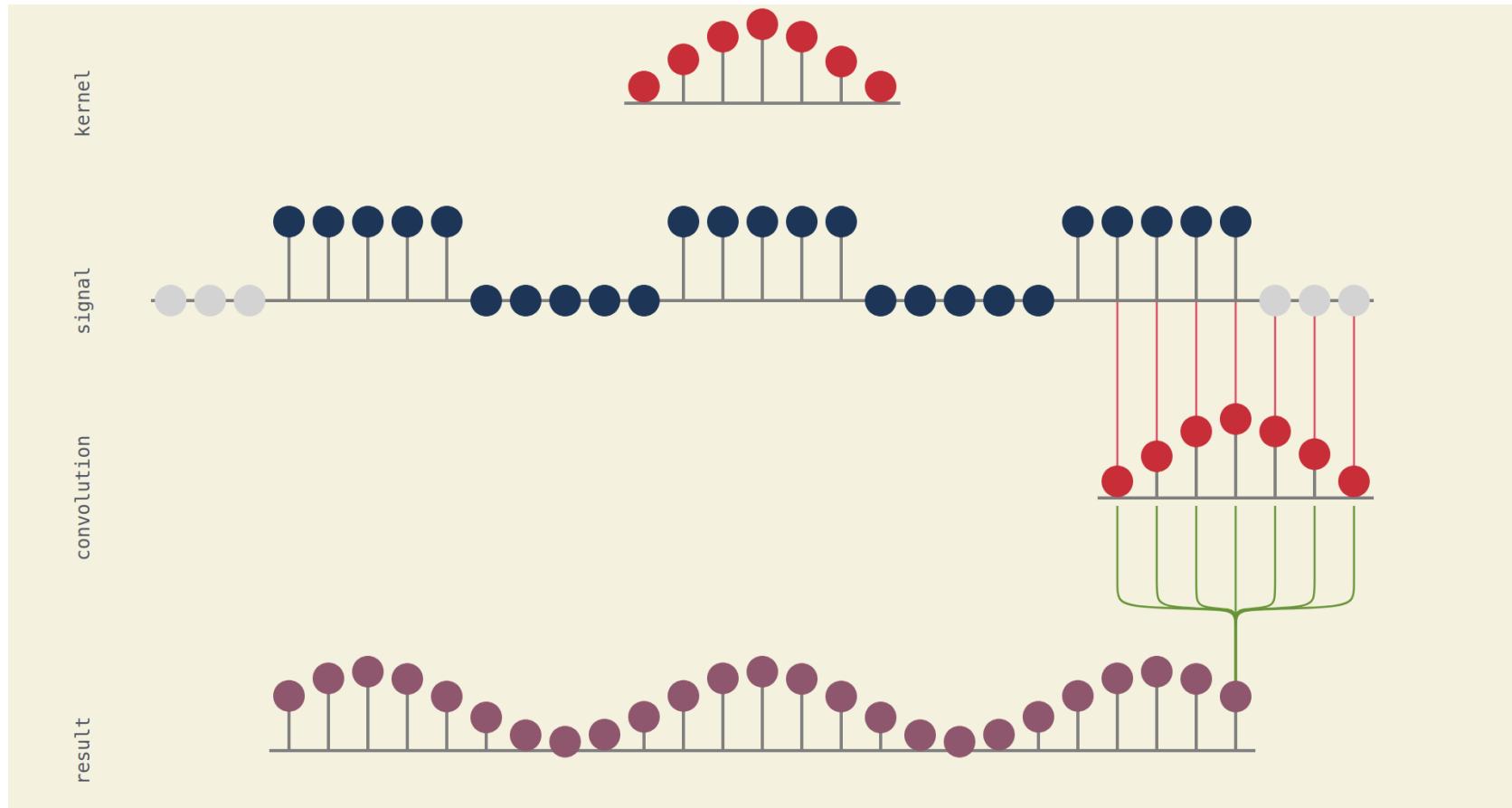
$$P(\square + \blacksquare = 10) = a_4 \cdot b_6 + a_5 \cdot b_5 + a_6 \cdot b_4$$

$$P(\square + \blacksquare = 11) = a_5 \cdot b_6 + a_6 \cdot b_5$$

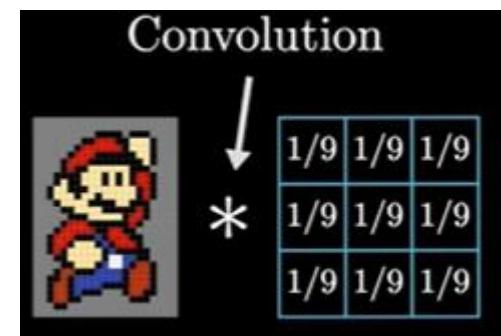
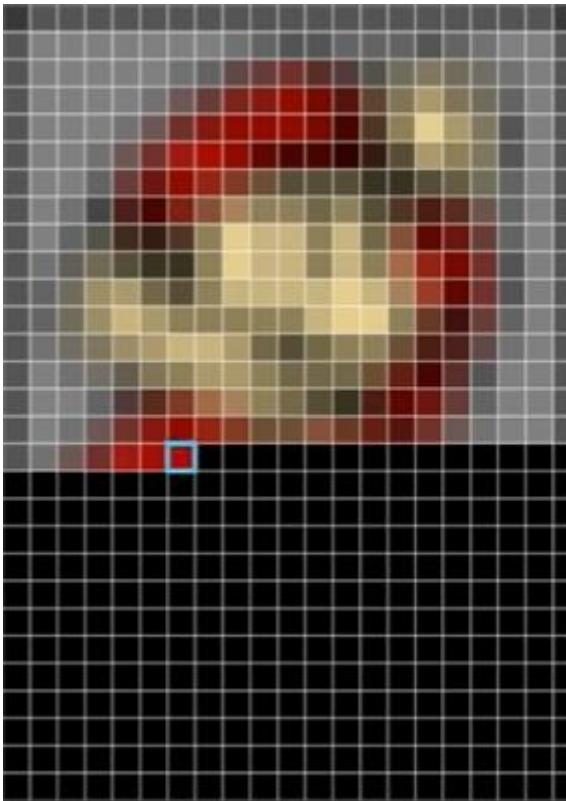
$$P(\square + \blacksquare = 12) = a_6 \cdot b_6$$

Example: $(1,2,3) * (4,5,6)$

Example: Moving average



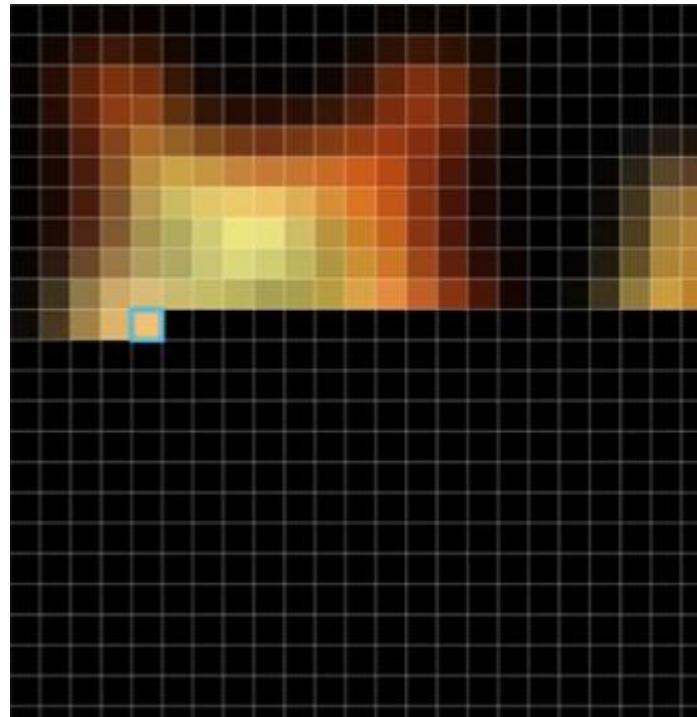
Example: 2D



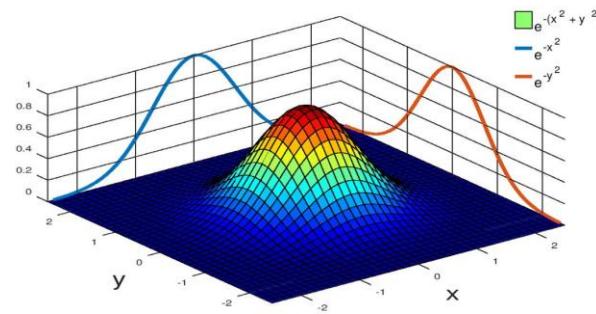
1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

$$\frac{1}{9} \begin{bmatrix} 0.5 \\ 0.5 \\ 0.5 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.8 \\ 0.0 \\ 0.0 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.8 \\ 0.0 \\ 0.0 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.8 \\ 0.0 \\ 0.0 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix} + \frac{1}{9} \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

Example: 2D



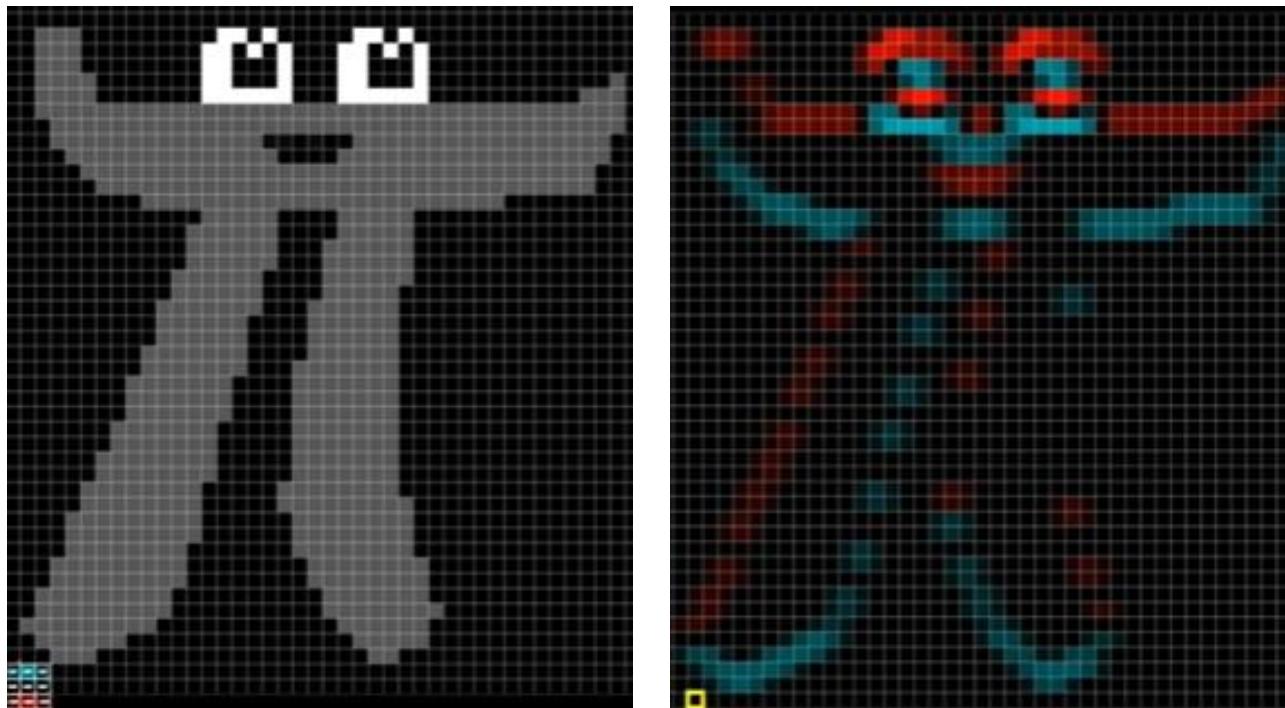
0.003	0.013	0.022	0.013	0.003
0.013	0.060	0.098	0.060	0.013
0.022	0.098	0.162	0.098	0.022
0.013	0.060	0.098	0.060	0.013
0.003	0.013	0.022	0.013	0.003



Example: 2D

kernel

0.25	0.00	-0.25
0.50	0.00	-0.50
0.25	0.00	-0.25



Example: (3,1,4,1,5,9) * (7,7,5)

- Classical signal processing:
 - mean kernel → spatial smoothing
 - Gaussian kernel → physically meaningful smoothing
 - Zero-sum kernel → edge detection
- However, real-world tasks are different:
 - Typhoon structure detection
 - Atmospheric pattern identification
 - Temperature field prediction

What is the correct kernel for a given task?

CNN: Learning the Convolutional operator

- Any linear and shift-invariant operator can be written as a convolution:

$$T(f) = f * w,$$

where w is the kernel.

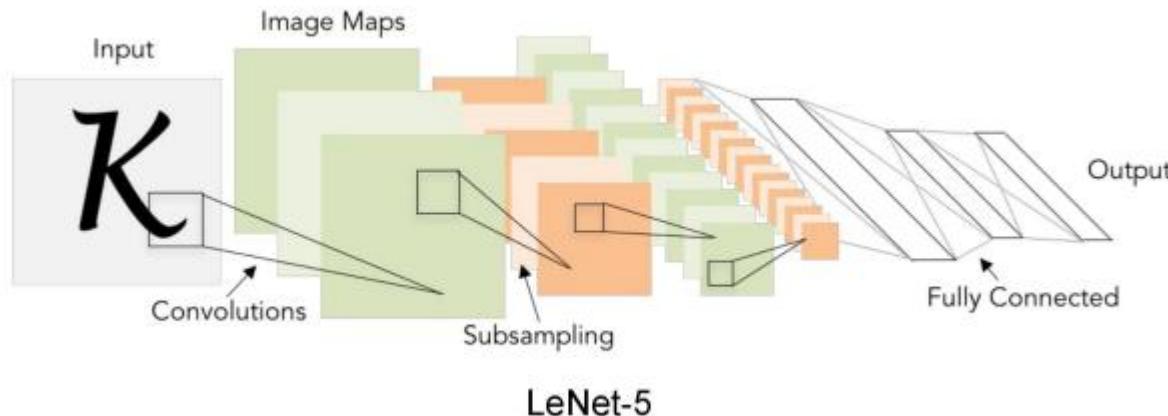
- In a convolutional neural network:
 - A *convolution layer* implements a linear and shift-invariant operator

$$(a * b)[n] = \sum a[k]b[n - k]$$

- The kernel w is treated as *learnable parameters*
- Training *optimizes the kernel* for the task

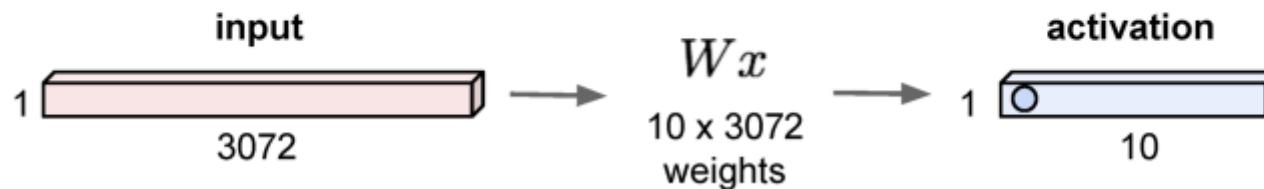
Convolutional neural networks

- Layer types:
 - Fully-connected layer
 - Convolutional layer
 - Pooling layer

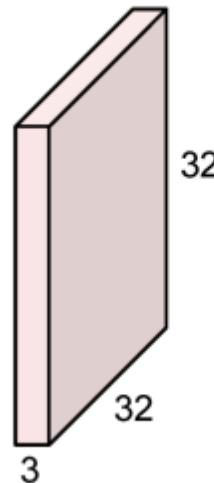


Fully Connected Layer & Convolution Layer

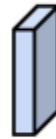
32x32x3 image -> stretch to 3072 x 1



32x32x3 image



5x5x3 filter



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

- Fully connected (MLP):

$$y = Wx$$

- Each spatial location has independent weights
- Not shift-invariant and ignores locality

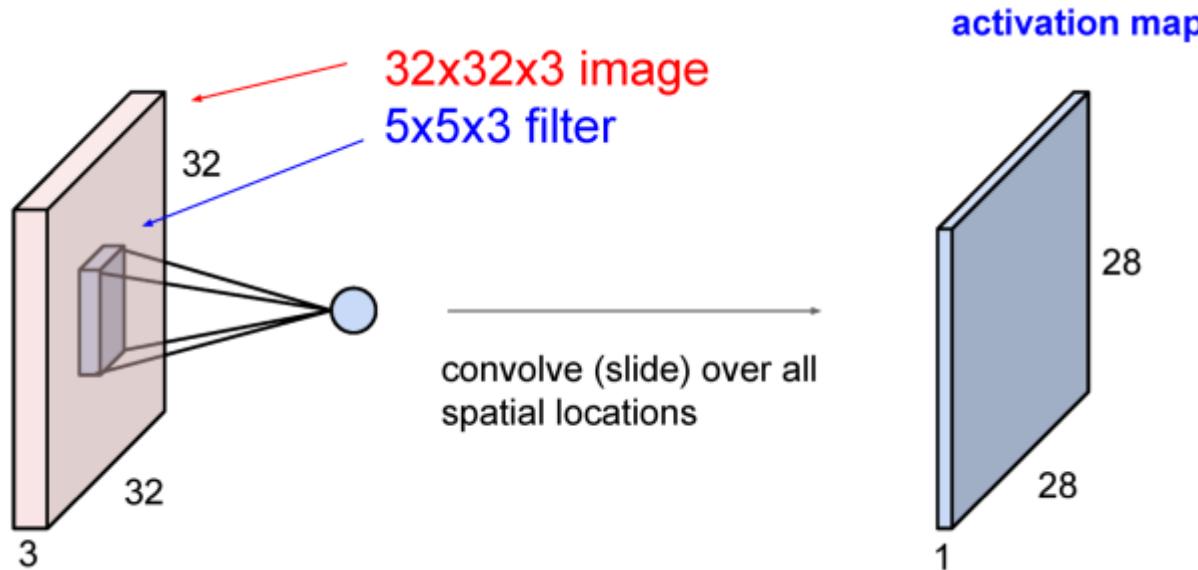
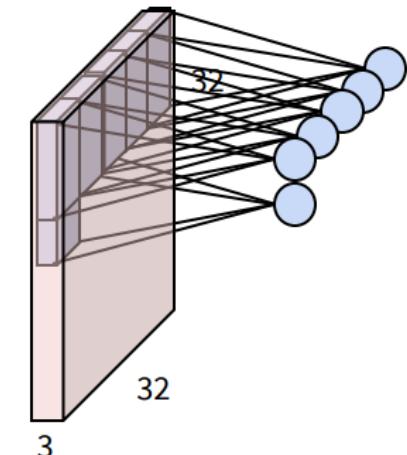
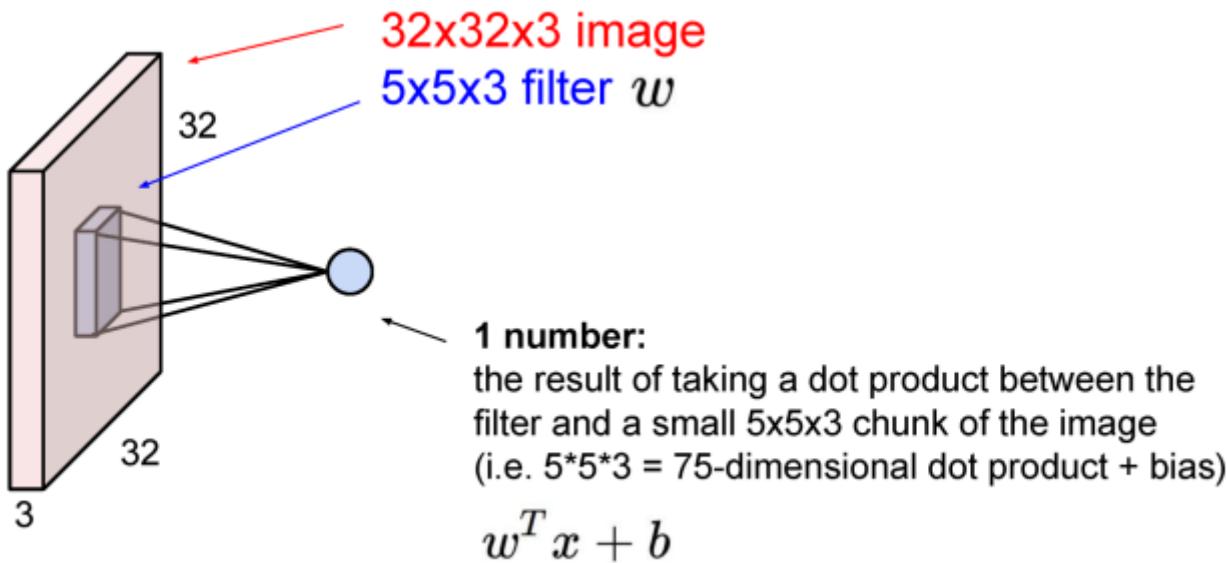
- Convolution:

$$(T * K)(x, y) = \sum_{i,j} T(i, j)K(x - i, y - j)$$

- Same kernel applied at every location (weight sharing)
- Each output is a weighted sum of neighboring values

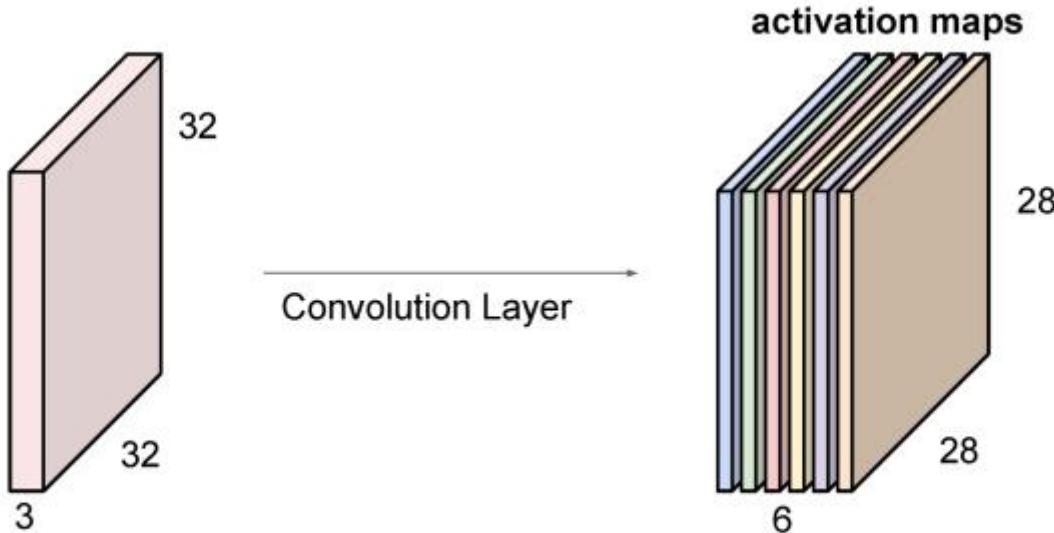
*For spatially continuous satellite fields,
convolution preserves structures.*

Convolution Layer



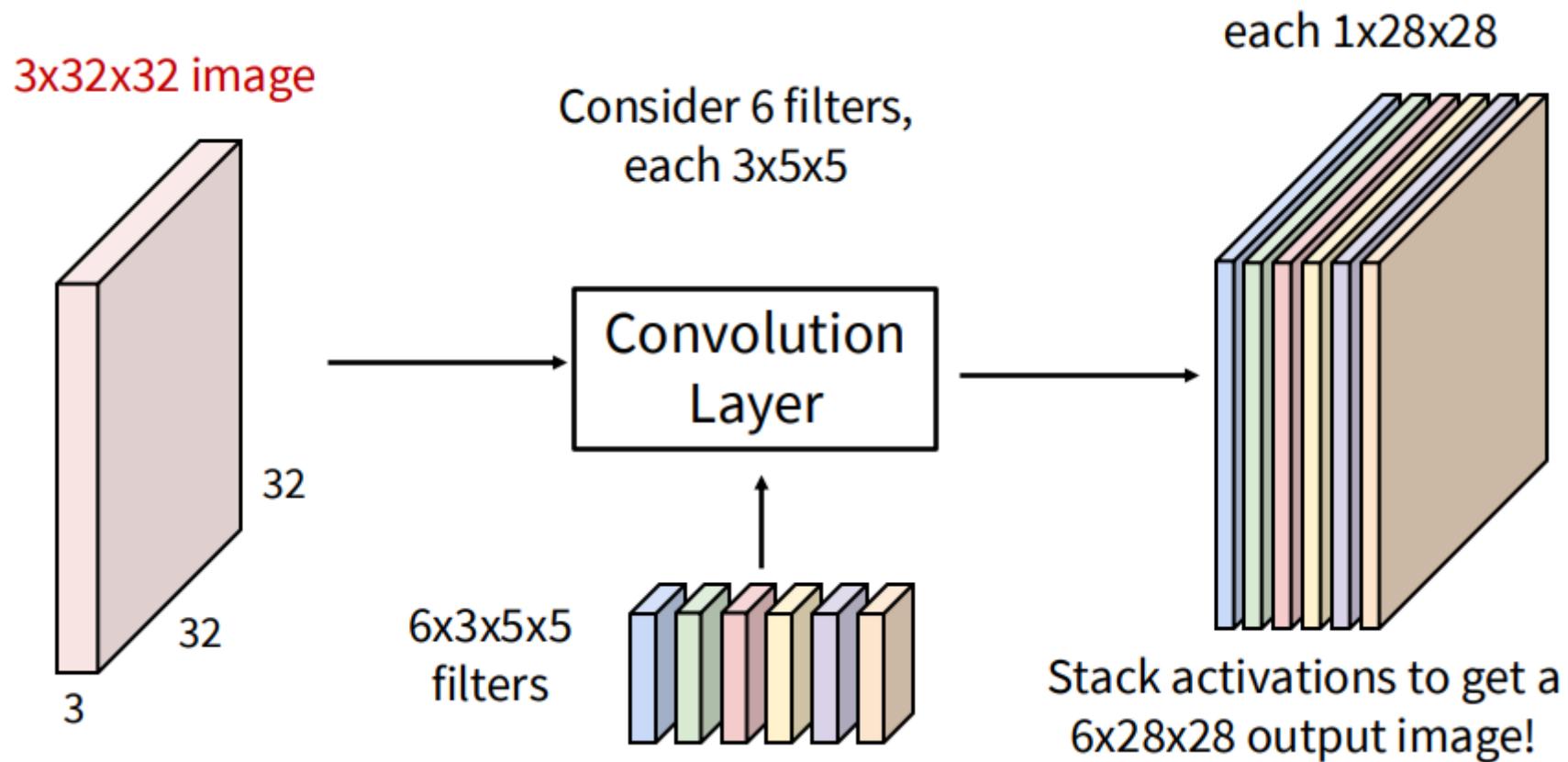
Convolution Layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

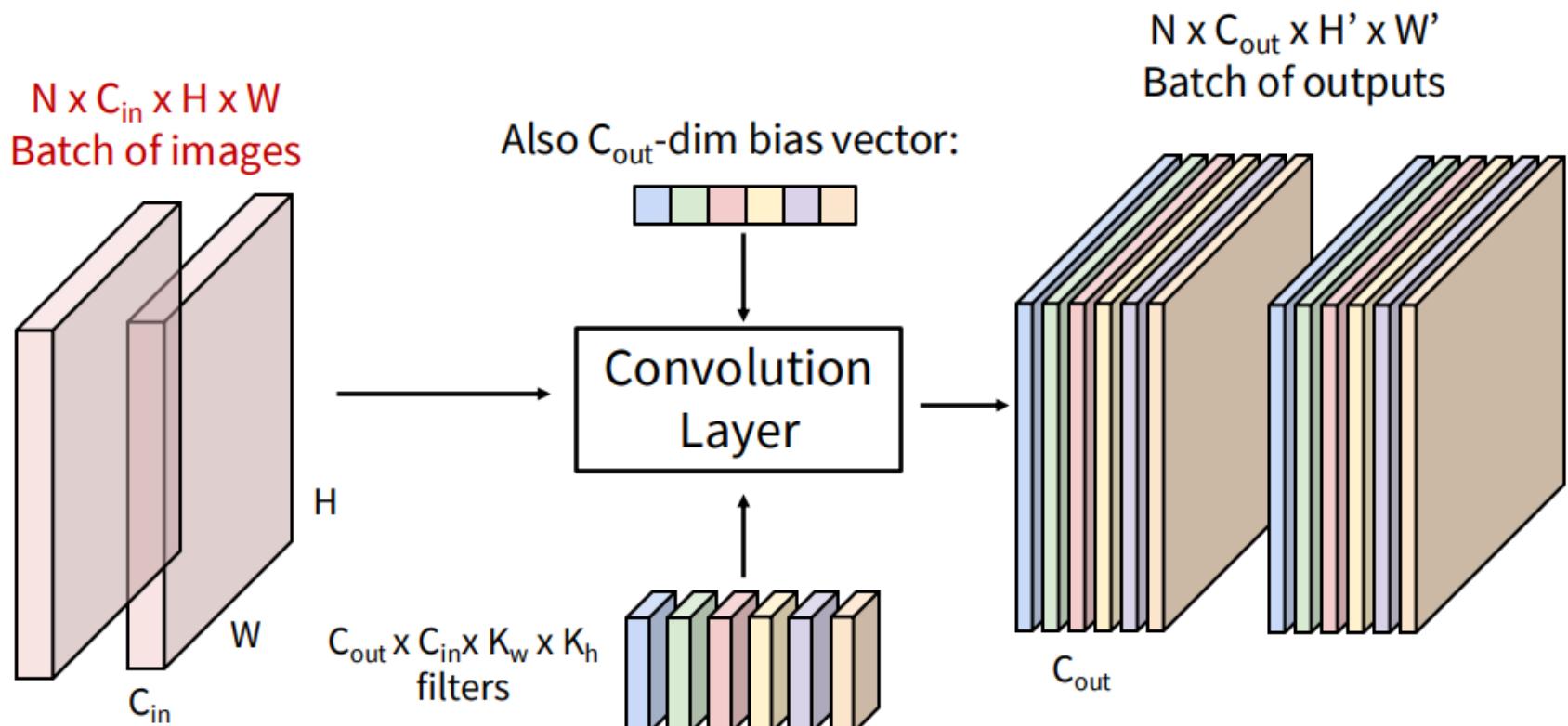


We stack these up to get a “new image” of size 28x28x6!

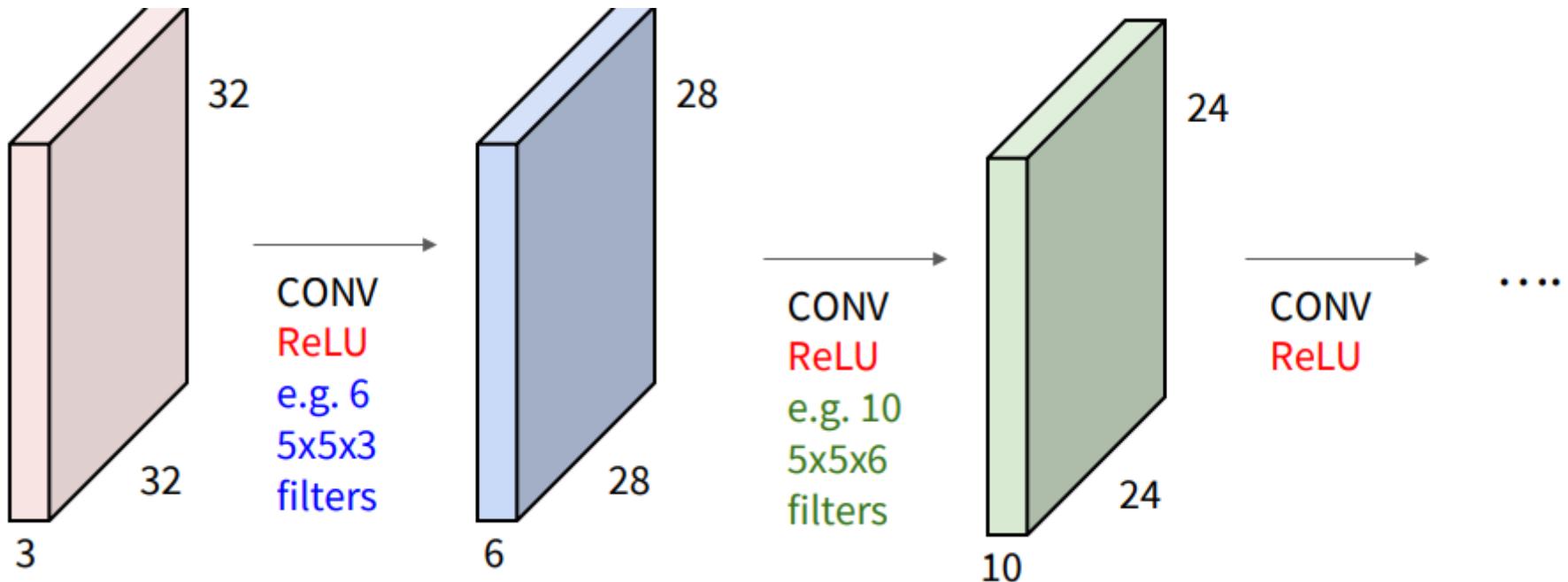
Convolution Layer



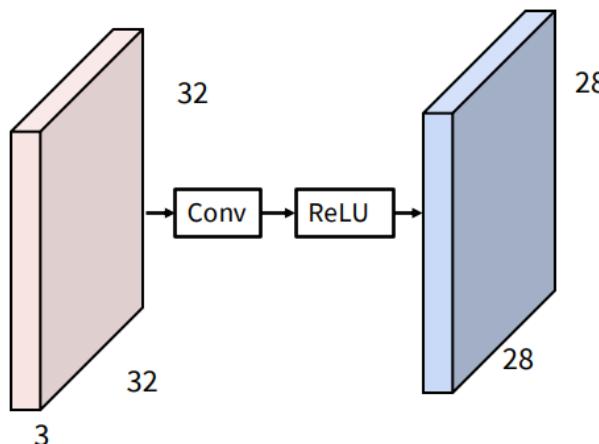
Convolution Layer



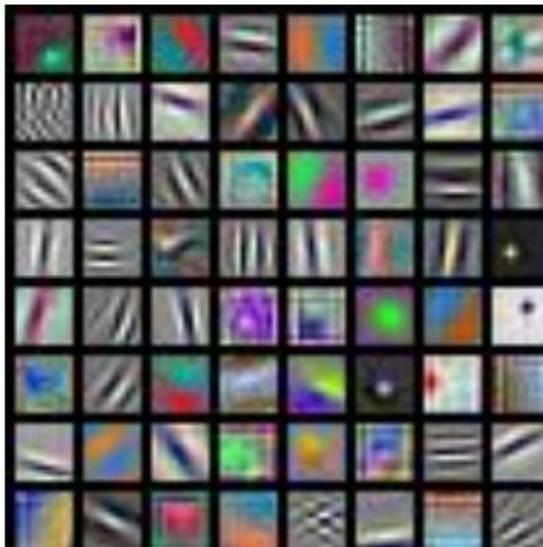
Convolution Layer



What do Conv filters learn?



First-layer conv filters: local image templates
(Often learns oriented edges, opposing colors)



AlexNet: 64 filters, each $3 \times 11 \times 11$

Deeper conv layers: Harder to visualize
Tend to learn larger structures e.g. eyes, letters



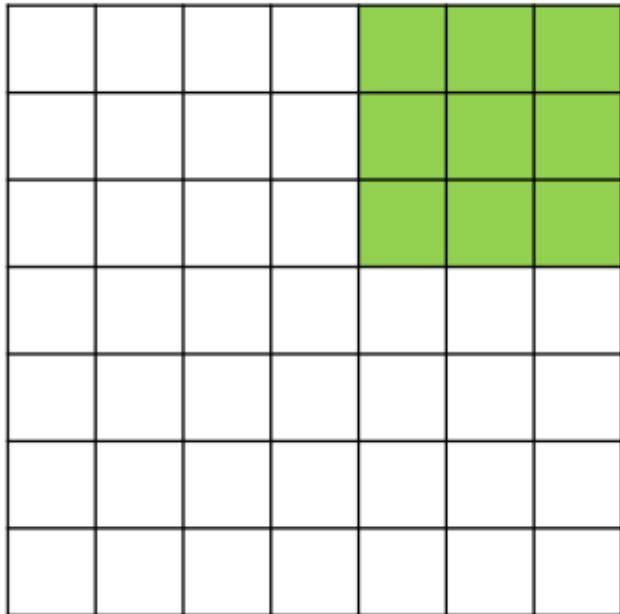
6th layer conv layer from an ImageNet model

Visualization from [Springenberg et al, ICLR 2015]

Convolution: Spatial Dimensions

7

7



Input: 7x7
Filter: 3x3
Output: 5x5

Problem: Feature maps shrink with each layer!

In general
Input: W
Filter: K
Output: $W - K + 1$

Convolution: Spatial Dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7

Filter: 3x3

Output: 5x5

Problem: Feature maps shrink with each layer!

In general

Input: W

Filter: K

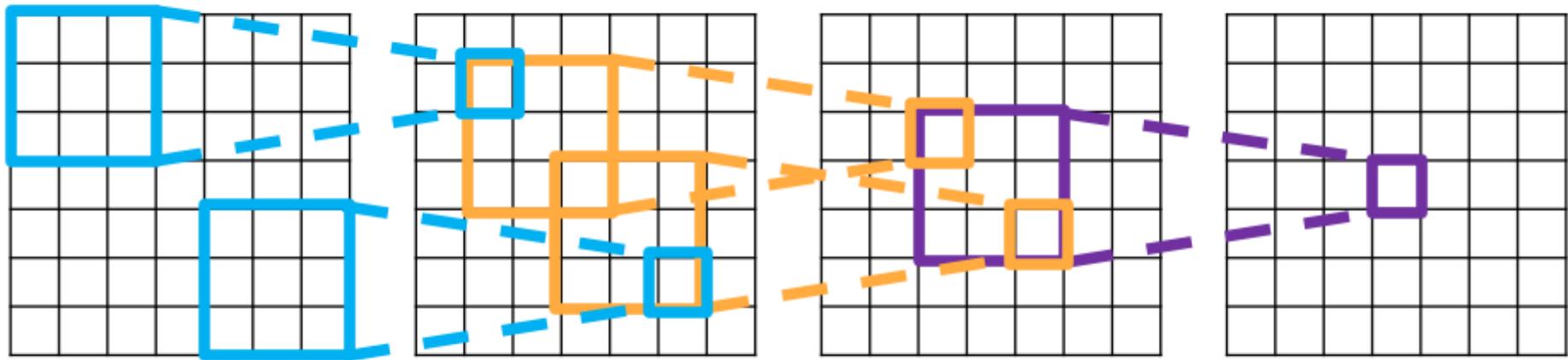
Padding: P

Output: $W - K + 1 + 2P$

Solution: Add **padding** around the input before sliding the filter

Receptive fields

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$



Input

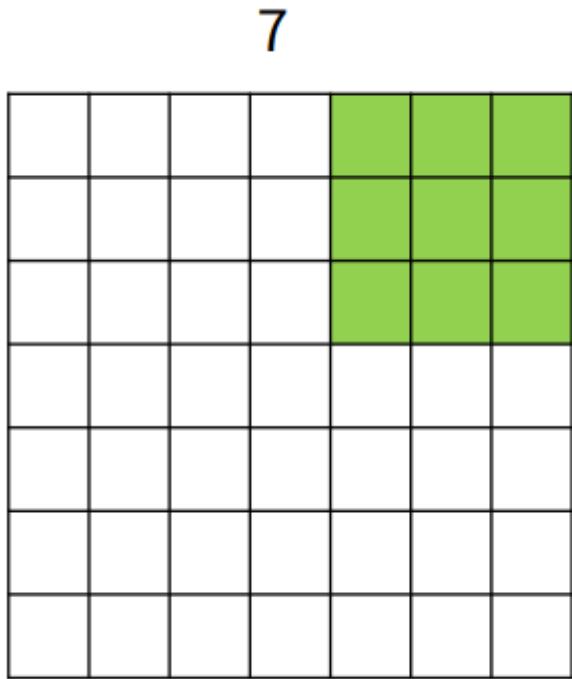
Problem: For large images we need many layers for each output to “see” the whole image

Solution: Downsample inside the network

Output

Slide inspiration: Justin Johnson

Strided Convolution



Input: 7x7
Filter: 3x3
Stride: 2
Output: 3x3

7

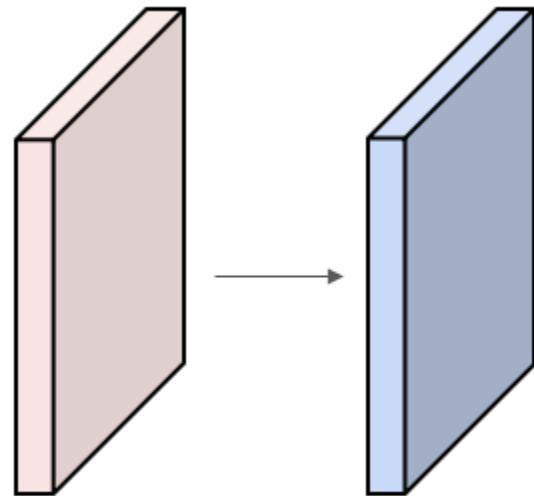
In general:
Input: W
Filter: K
Padding: P
Stride: S

Output:
$$(W - K + 2P) / S + 1$$

Convolution example

Input volume: $3 \times 32 \times 32$
10 5x5 filters with stride 1, pad 2

Output volume size: ?



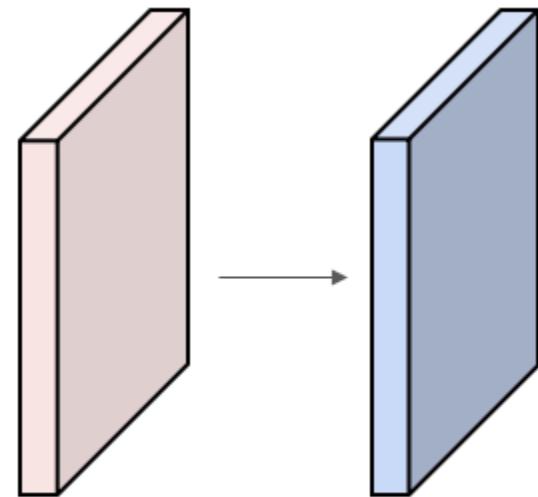
Convolution example

Input volume: $3 \times 32 \times 32$

$10 \times 5 \times 5$ filters with stride 1 , pad 2

Output volume size: $10 \times 32 \times 32$

$$32 = (32 + 2 * 2 - 5) / 1 + 1$$



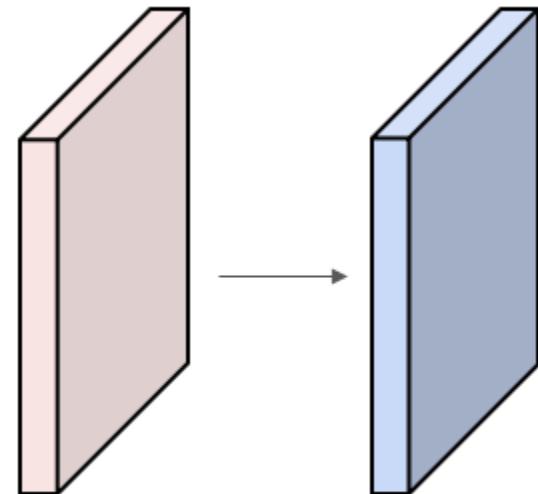
Convolution example

Input volume: **3** x 32 x 32

10 **5x5** filters with stride 1, pad 2

Output volume size: 10 x 32 x 32

Number of learnable parameters: ?



Number of learnable parameters: 760

Parameters per filter: **3*5*5 + 1** (for bias) = **76**

10 filters, so total is **10 * 76 = 760**

Convolution example

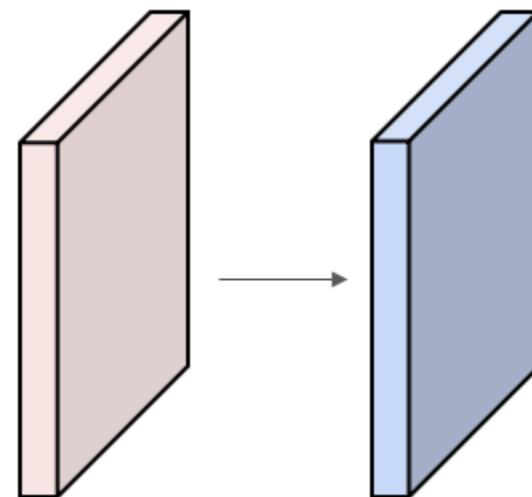
Input volume: **3** x 32 x 32

10 **5x5** filters with stride 1, pad 2

Output volume size: 10 x 32 x 32

Number of learnable parameters: 760

Number of multiply-add operations?



Number of learnable parameters: 760

Number of multiply-add operations: **768,000**

10*32*32 = 10,240 outputs

Each output is the inner product of two **3x5x5** tensors (75 elems)

Total = $75 * 10240 = \mathbf{768K}$

PyTorch Convolution Layer

```
CLASS torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0,  
dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) \[SOURCE\]
```

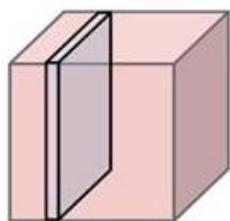
Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{\text{out}}, H_{\text{out}}, W_{\text{out}})$ can be precisely described as:

$$\text{out}(N_i, C_{\text{out}_j}) = \text{bias}(C_{\text{out}_j}) + \sum_{k=0}^{C_{\text{in}}-1} \text{weight}(C_{\text{out}_j}, k) \star \text{input}(N_i, k)$$

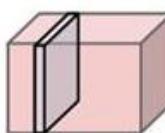
Pooling layer

$64 \times 112 \times 112$



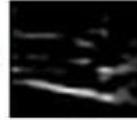
$64 \times 224 \times 224$

pool



Given an input $C \times H \times W$,
downsample each $1 \times H \times W$ plane

224



downsampling

112



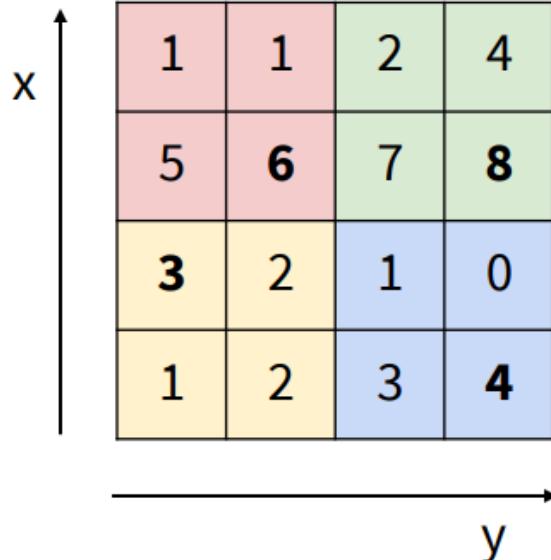
Hyperparameters:

Kernel Size

Stride

Pooling function

Single depth slice



$64 \times 224 \times 224$



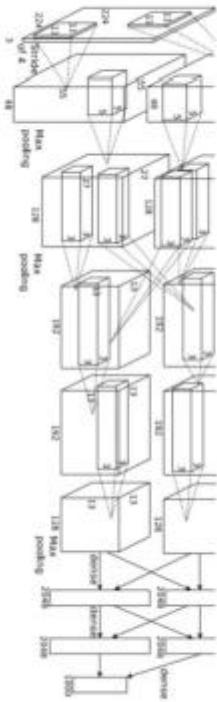
Max pooling with 2x2
kernel size and stride 2



Gives **invariance** to small spatial
shifts. No learnable parameters.

CNN Architectures

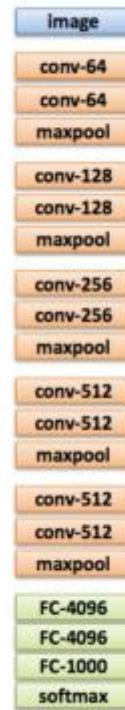
“AlexNet”



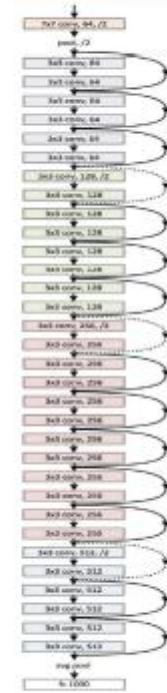
“GoogLeNet”



“VGG Net”



“ResNet”



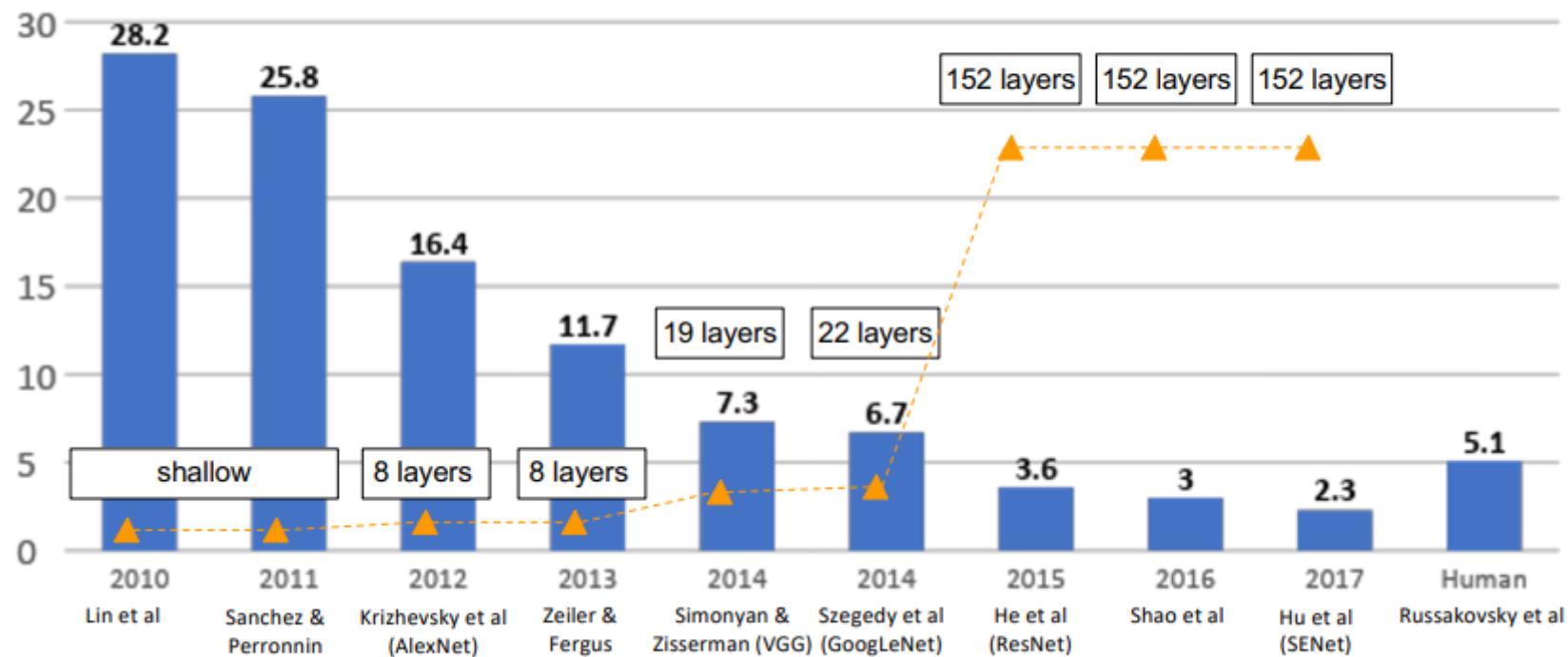
[Krizhevsky et al. NIPS 2012]

[Szegedy et al. CVPR 2015]

[Simonyan & Zisserman,
ICLR 2015]

[He et al. CVPR 2016]

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



High-level Idea: Learn parameters that let us **scale / shift the input data**

1. Normalize input data
2. Scale / shift using learned parameters

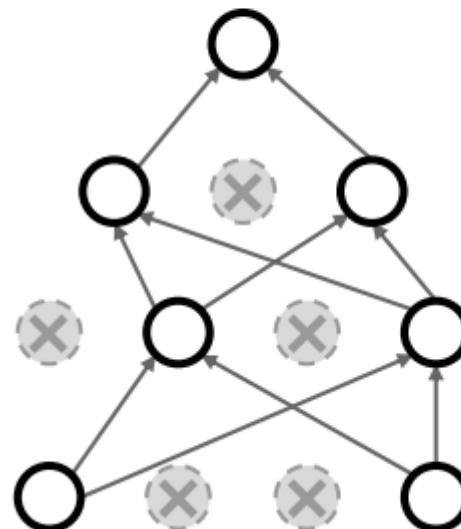
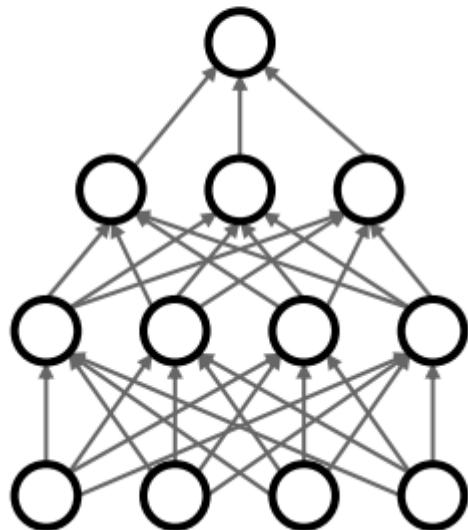
Statistics calculated per batch →

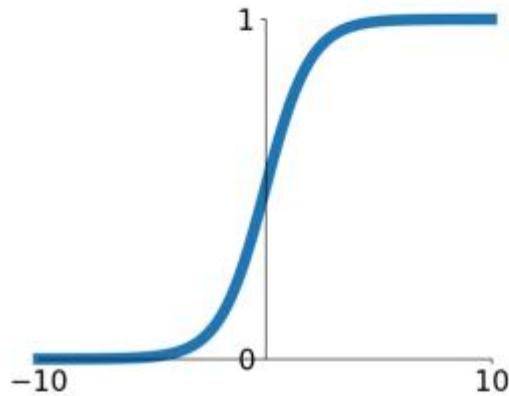
Learned parameters applied to each sample →

$$\begin{array}{c} \mathbf{x} : N \times D \\ \boxed{\text{Normalize}} \\ \boldsymbol{\mu}, \boldsymbol{\sigma} : N \times 1 \\ \mathbf{y}, \boldsymbol{\beta} : 1 \times D \\ \mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta} \end{array}$$

Regularization: Dropout

In each forward pass, randomly set some neurons to zero
Probability of dropping is a hyperparameter; 0.5 is common





Sigmoid

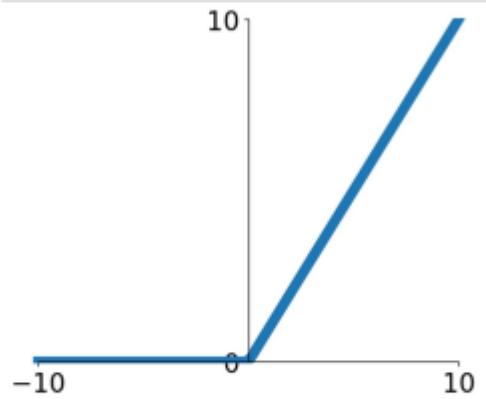
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

Key problem:

Many layers of sigmoids → smaller and smaller gradients.

Q: In which regions does sigmoid have a small gradient?



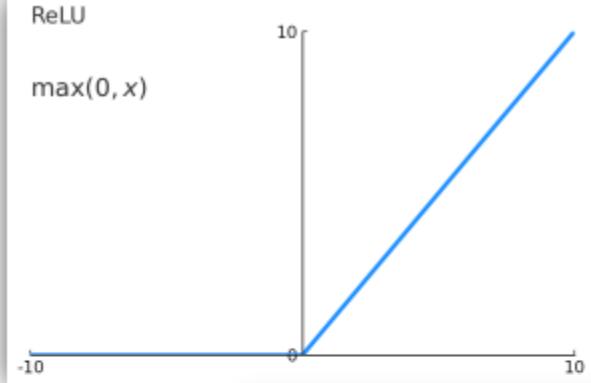
- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid in practice (e.g. 6x)

ReLU
(Rectified Linear Unit)

Activation Functions

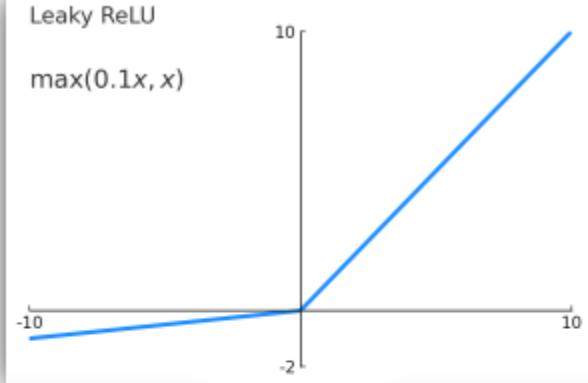
ReLU

$$\max(0, x)$$



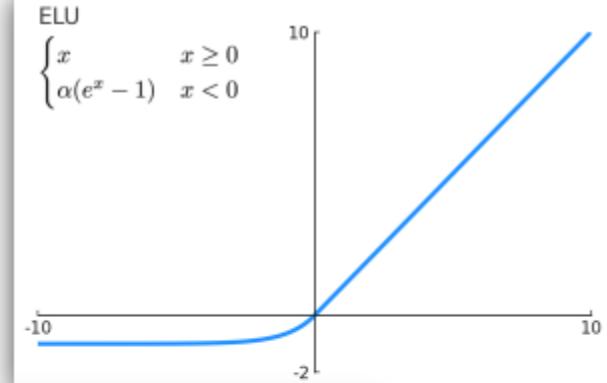
Leaky ReLU

$$\max(0.1x, x)$$



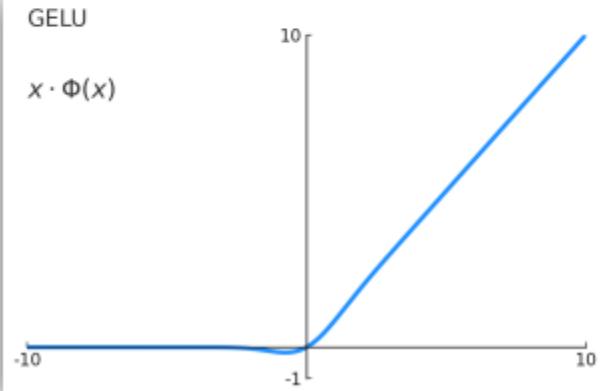
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



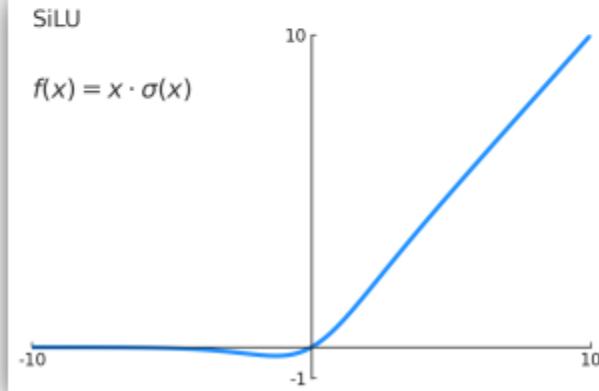
GELU

$$x \cdot \Phi(x)$$



SiLU

$$f(x) = x \cdot \sigma(x)$$



Case study: VGGNet

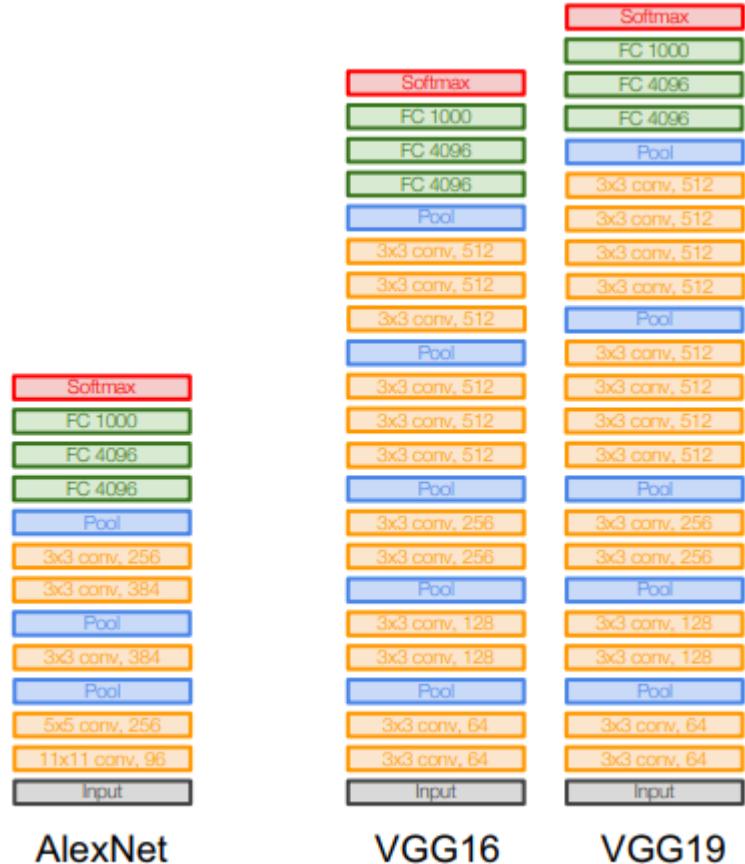
Small filters, Deeper networks

8 layers (AlexNet)

-> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13
(ZFNet)
-> 7.3% top 5 error in ILSVRC'14



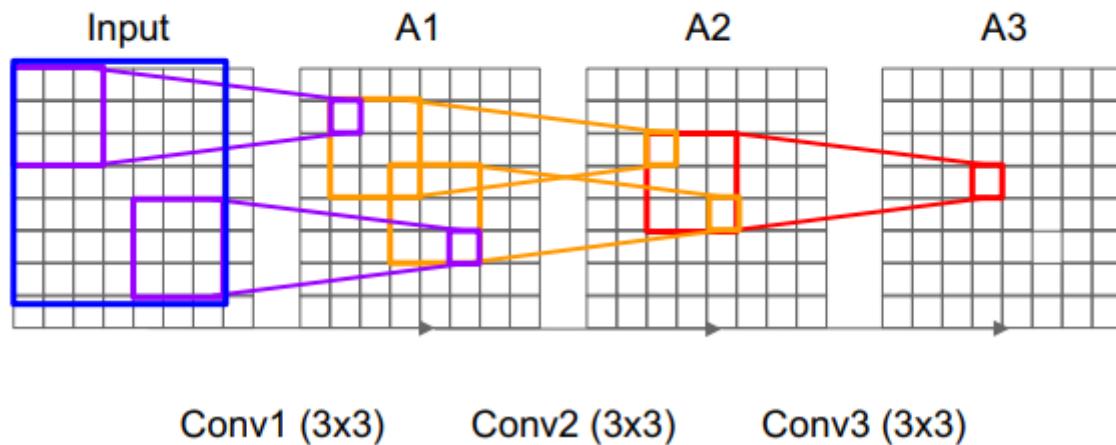
AlexNet

VGG16

VGG19

Case study: VGGNet

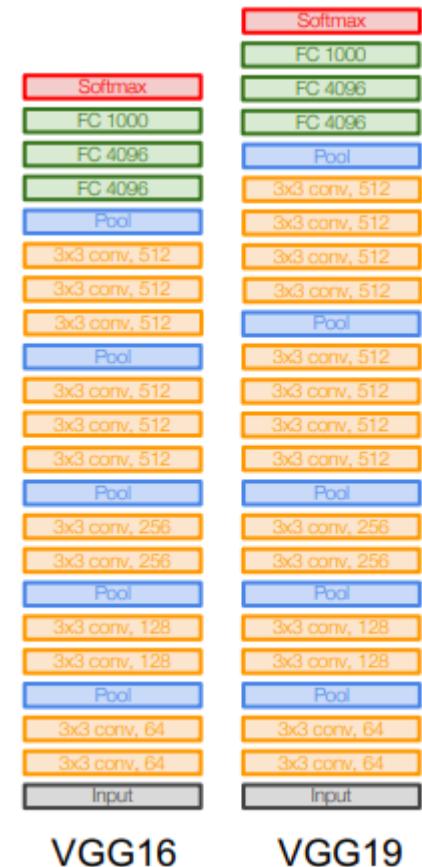
Q: What is the effective receptive field of three 3x3 conv (stride 1) layers?



Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

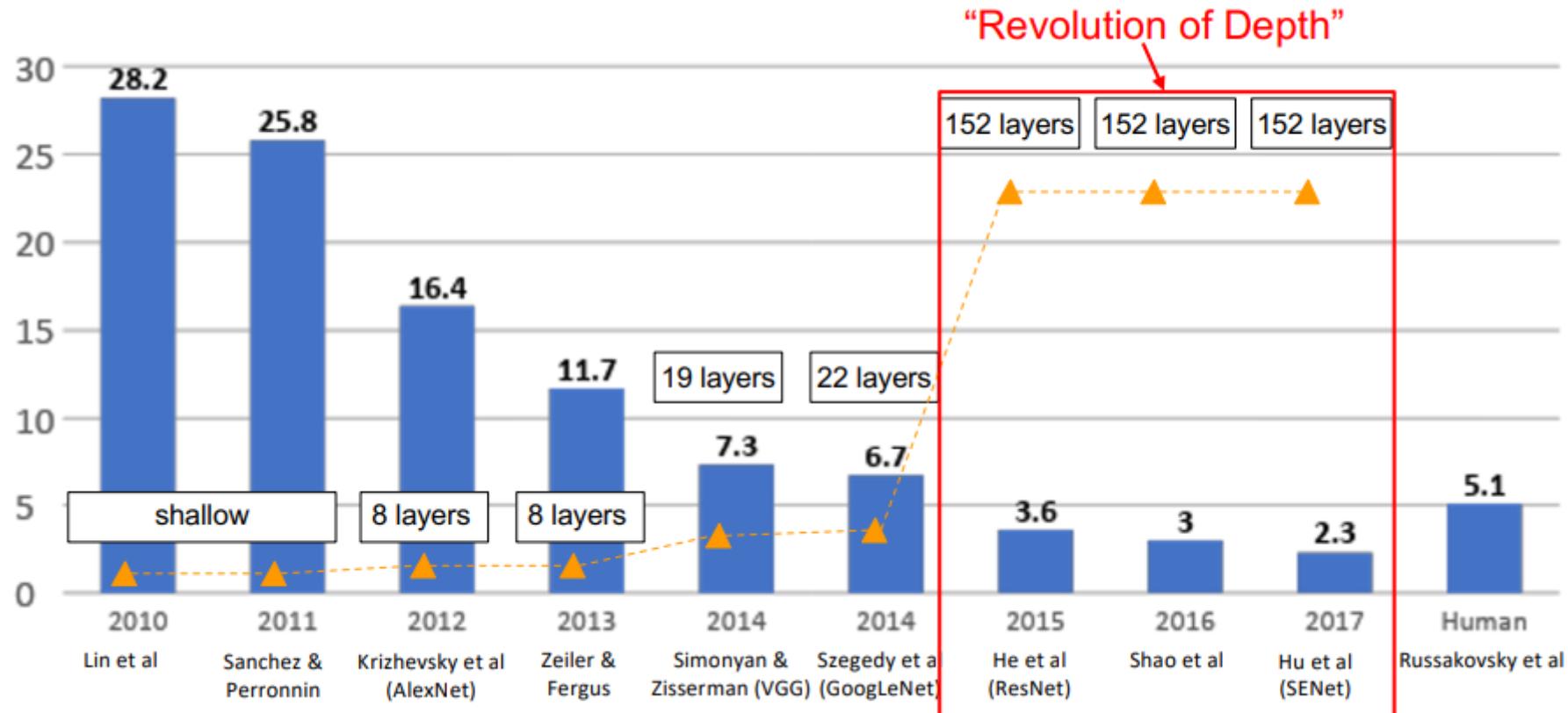
But deeper, more non-linearities

And fewer parameters: $3 * (3^2 C^2)$ vs. $7^2 C^2$ for C channels per layer



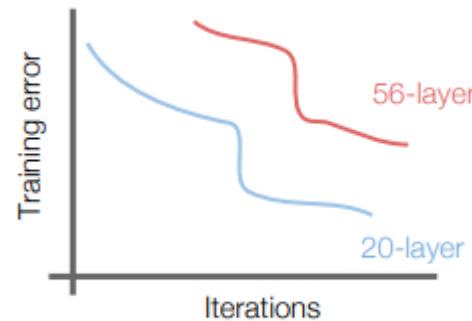
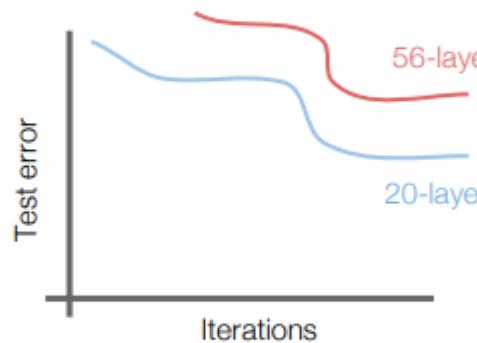
Case Study: ResNet

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Case Study: ResNet

What happens when we continue stacking deeper layers on a “plain” convolutional neural network?



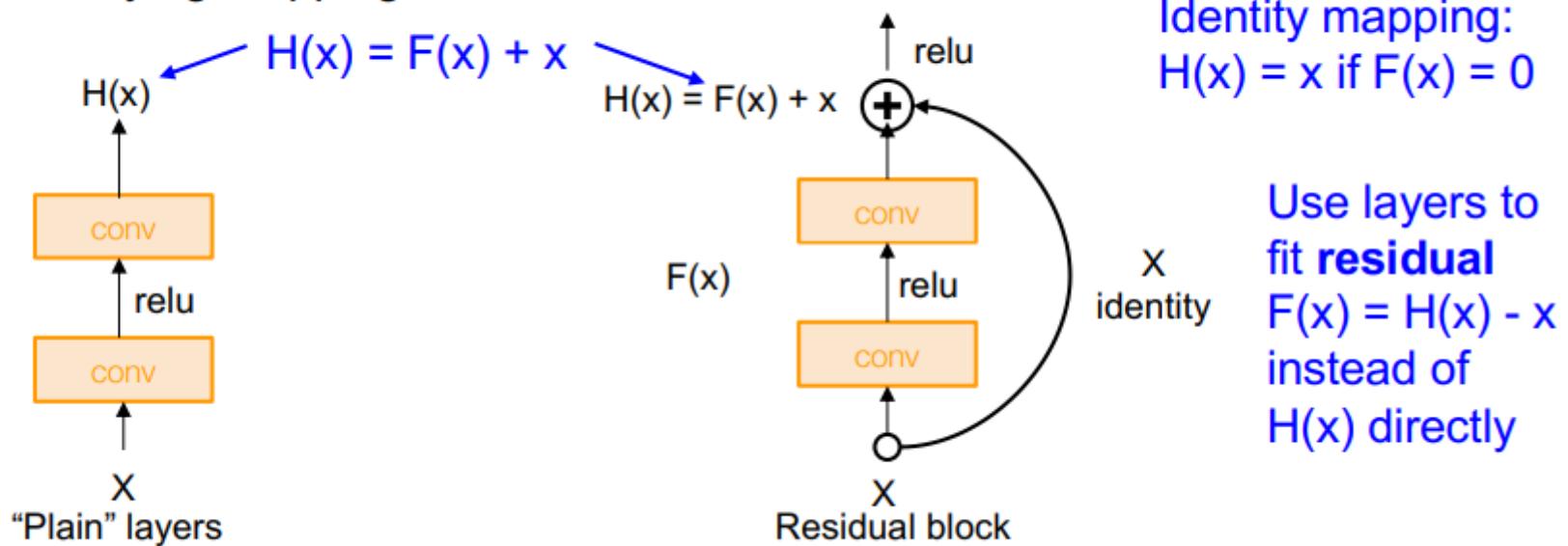
56-layer model performs worse on both test and training error
-> The deeper model performs worse, but it's **not caused by overfitting!**

Fact: Deep models have more representation power (more parameters) than shallower models.

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

Case Study: ResNet

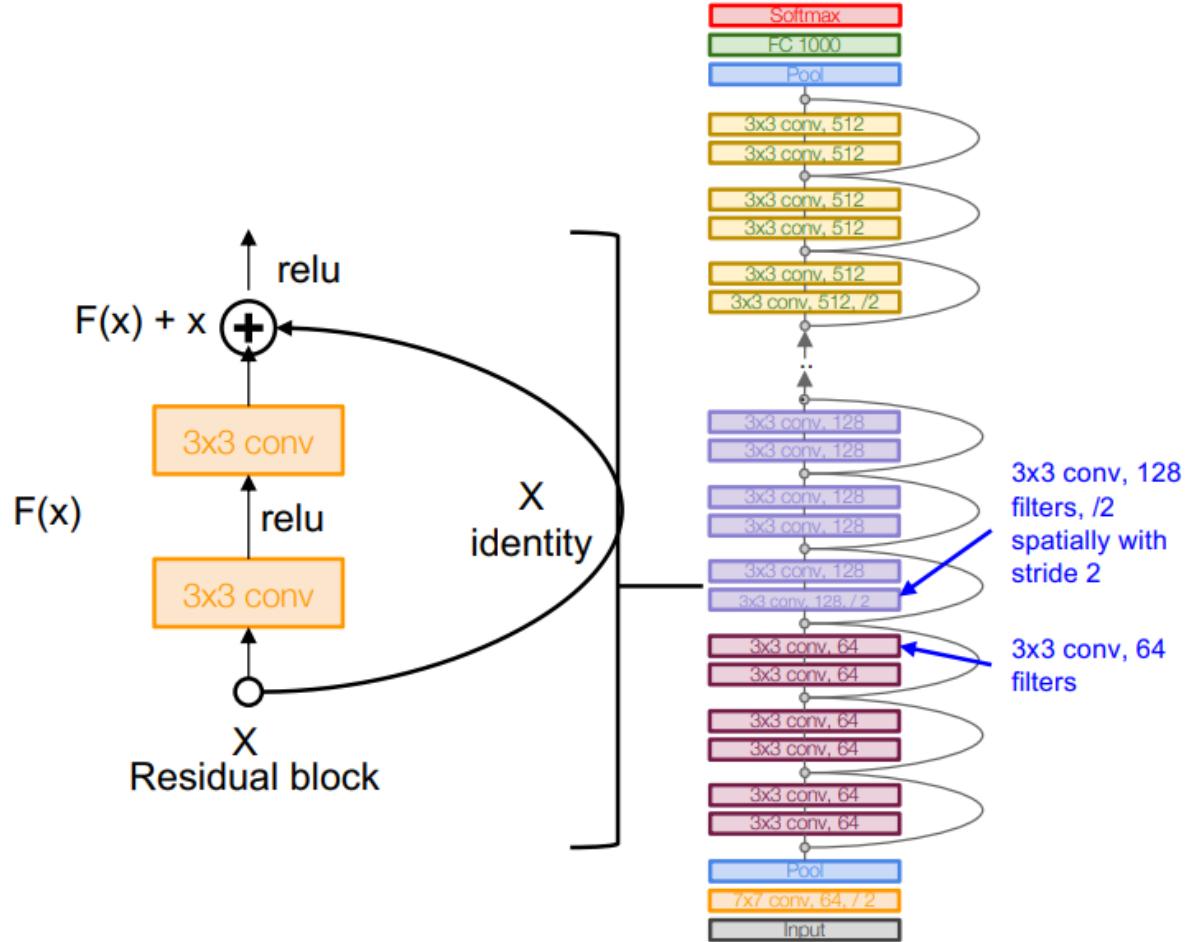
Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



Case Study: ResNet

Full ResNet architecture:

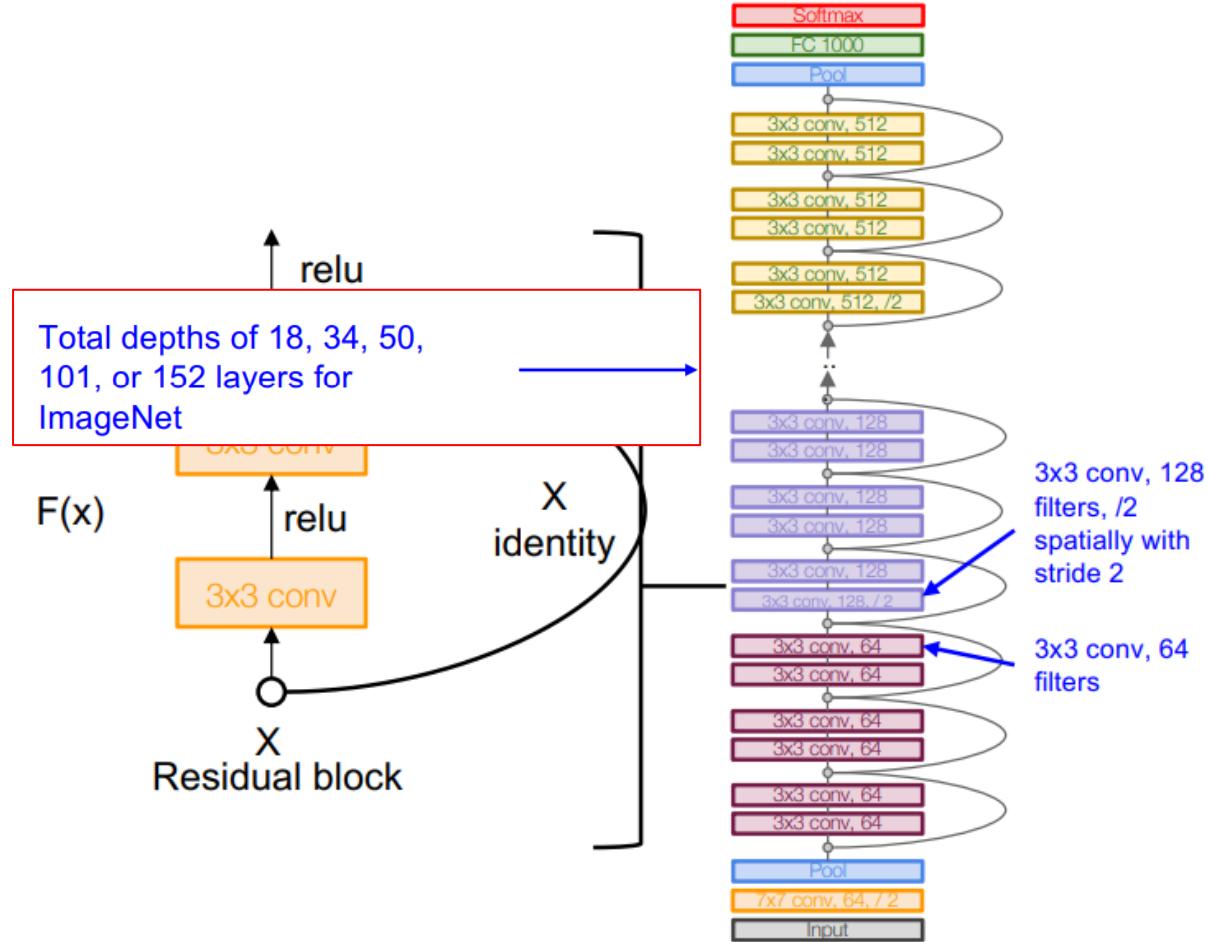
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
Reduce the activation volume by half.

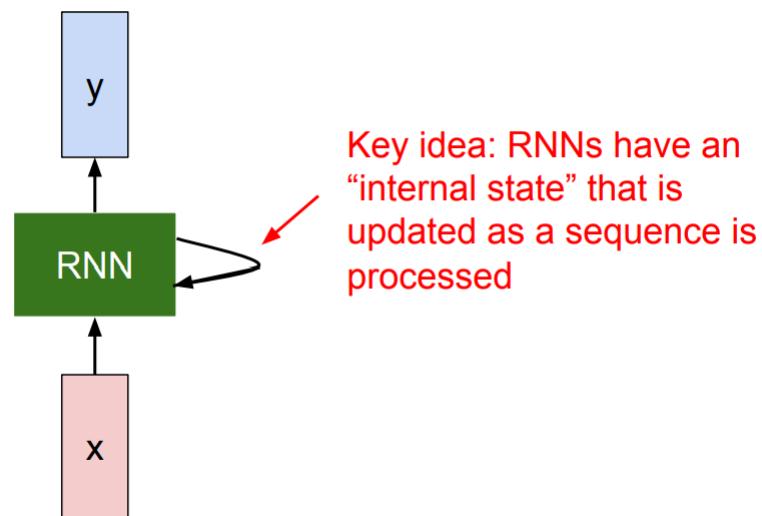
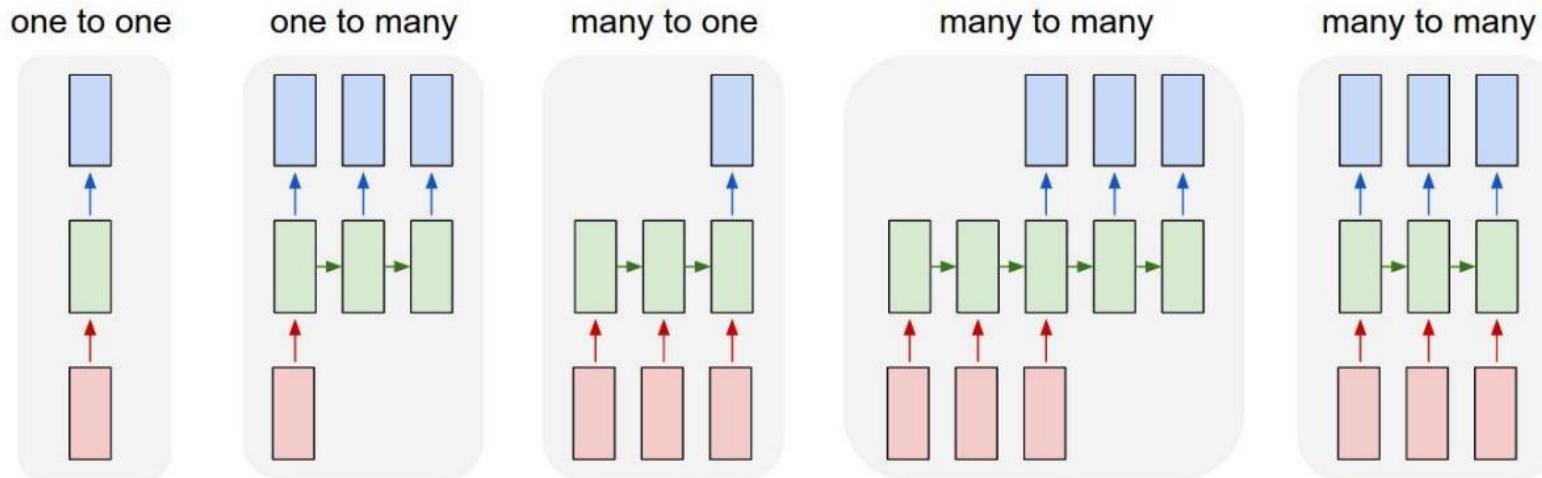


Case Study: ResNet

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
Reduce the activation volume by half.



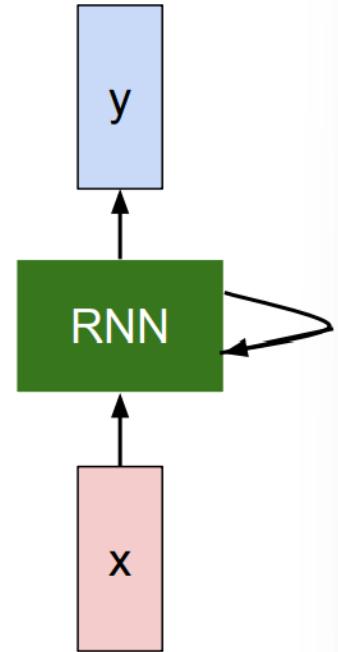


RNN hidden state update

We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state / old state input vector at
 | some function some time step
 \ with parameters W

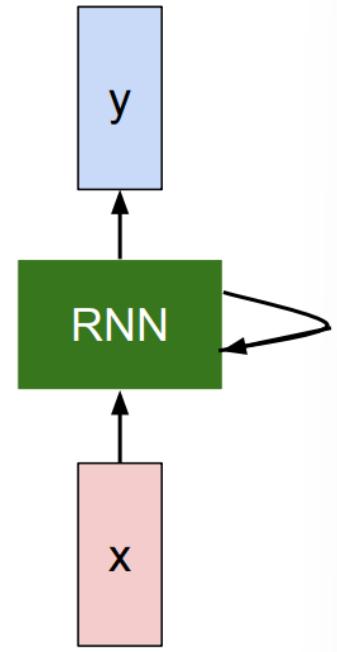


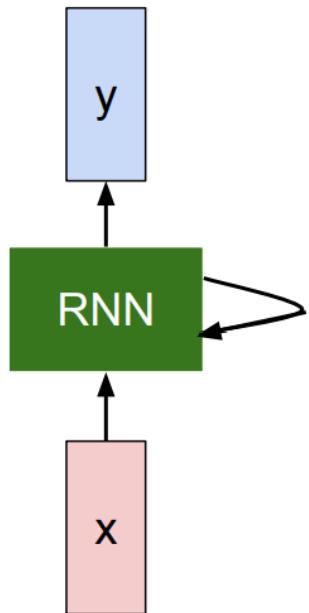
RNN output generation

We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

$$y_t = f_{W_{hy}}(h_t)$$

output new state
another function
with parameters W_o



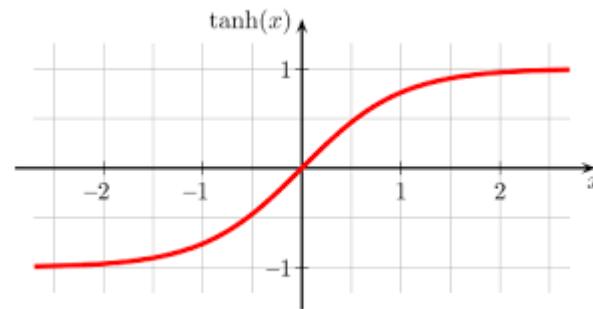


$$h_t = f_W(h_{t-1}, x_t)$$

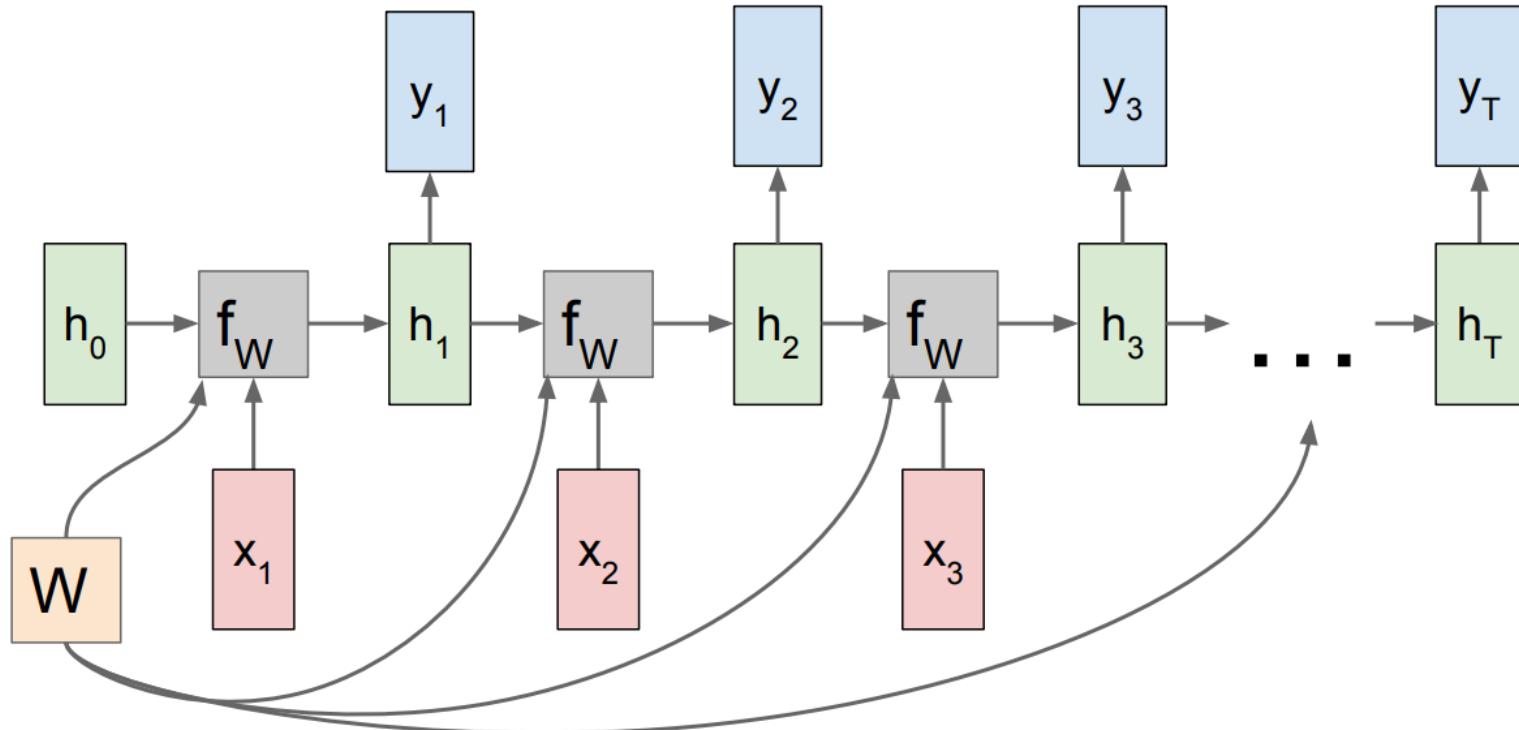
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

Sometimes called a “Vanilla RNN” or an “Elman RNN” after Prof. Jeffrey Elman



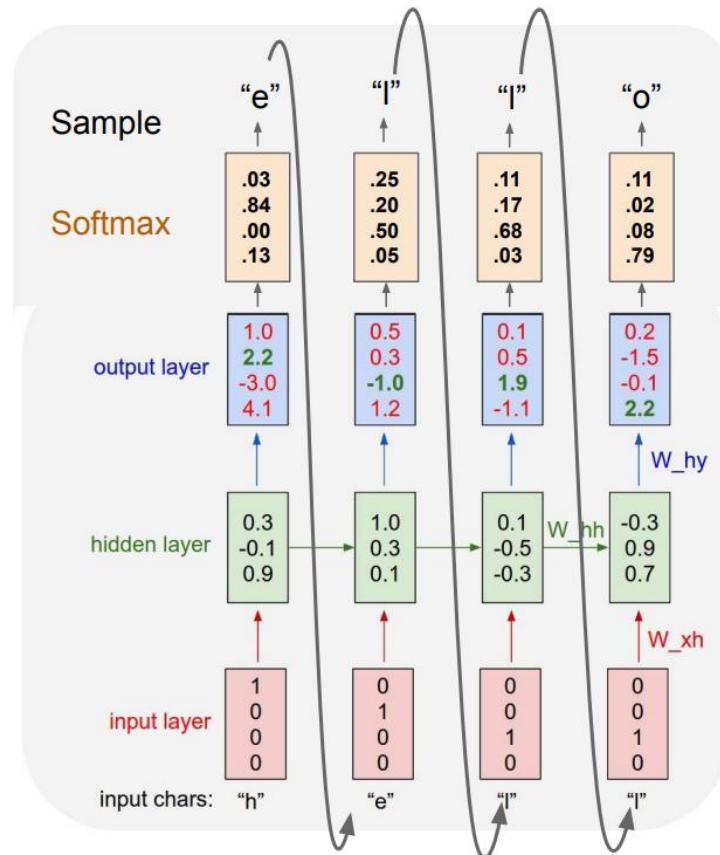
```
1   import torch
2   import torch.nn as nn
3
4   class MyRNN(nn.Module):
5       def __init__(self, input_size, hidden_size, output_size):
6           super().__init__()
7
8           self.hidden_size = hidden_size
9           self.i2h = nn.Linear(input_size + hidden_size, hidden_size)
10          self.h2o = nn.Linear(hidden_size, output_size)
11
12      def forward(self, input, hidden):
13          combined = torch.cat((input,hidden), 1)
14          hidden = torch.tanh(self.i2h(combined))
15          output = self.h2o(hidden)
16          return output, hidden
17
18      def get_hidden(self):
19          return torch.zeros(1, self.hidden_size)
20
21
22
23  rnn_model = MyRNN(input_size=4, hidden_size=1024, output_size=2)
24  hidden = rnn_model.get_hidden()
```



Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

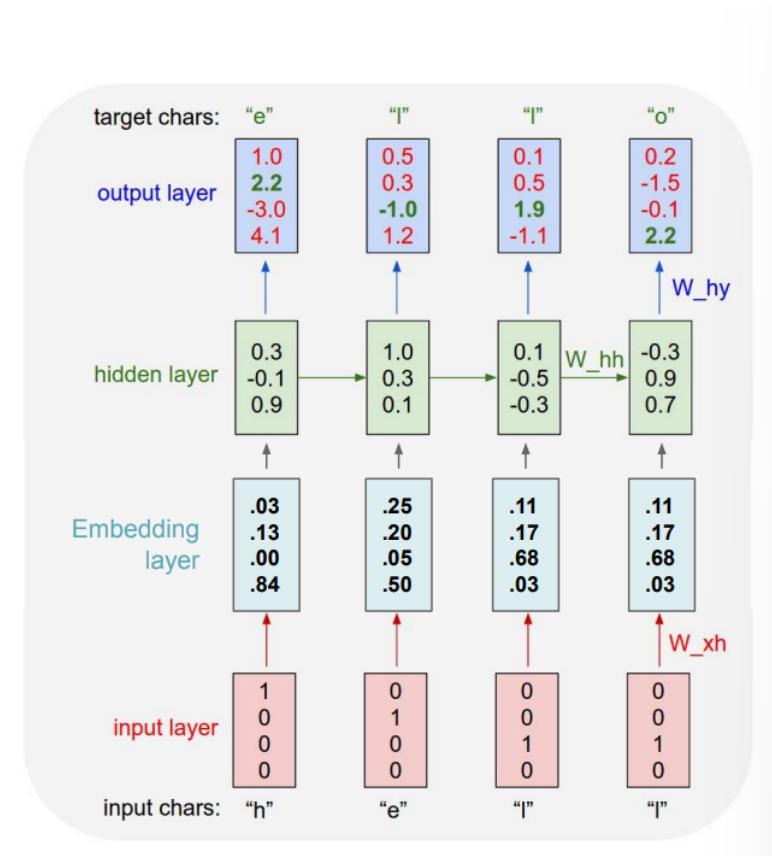
At test-time sample characters one at a time, feed back to model



Example: Character-level Language Model Sampling

$$\begin{aligned}
 & [w_{11} \ w_{12} \ w_{13} \ w_{14}] [1] = [w_{11}] \\
 & [w_{21} \ w_{22} \ w_{23} \ w_{14}] [0] = [w_{21}] \\
 & [w_{31} \ w_{32} \ w_{33} \ w_{14}] [0] = [w_{31}] \\
 & \quad [0]
 \end{aligned}$$

Matrix multiply with a one-hot vector just extracts a column from the weight matrix. We often put a separate **embedding** layer between input and hidden layers.



Word embedding

Vocabulary:
Man, woman, boy,
girl, prince,
princess, queen,
king, monarch

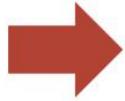


	1	2	3	4	5	6	7	8	9
man	1	0	0	0	0	0	0	0	0
woman	0	1	0	0	0	0	0	0	0
boy	0	0	1	0	0	0	0	0	0
girl	0	0	0	1	0	0	0	0	0
prince	0	0	0	0	1	0	0	0	0
princess	0	0	0	0	0	1	0	0	0
queen	0	0	0	0	0	0	1	0	0
king	0	0	0	0	0	0	0	1	0
monarch	0	0	0	0	0	0	0	0	1

Each word gets
a 1x9 vector
representation

Try to build a lower dimensional embedding

Vocabulary:
Man, woman, boy,
girl, prince,
princess, queen,
king, monarch

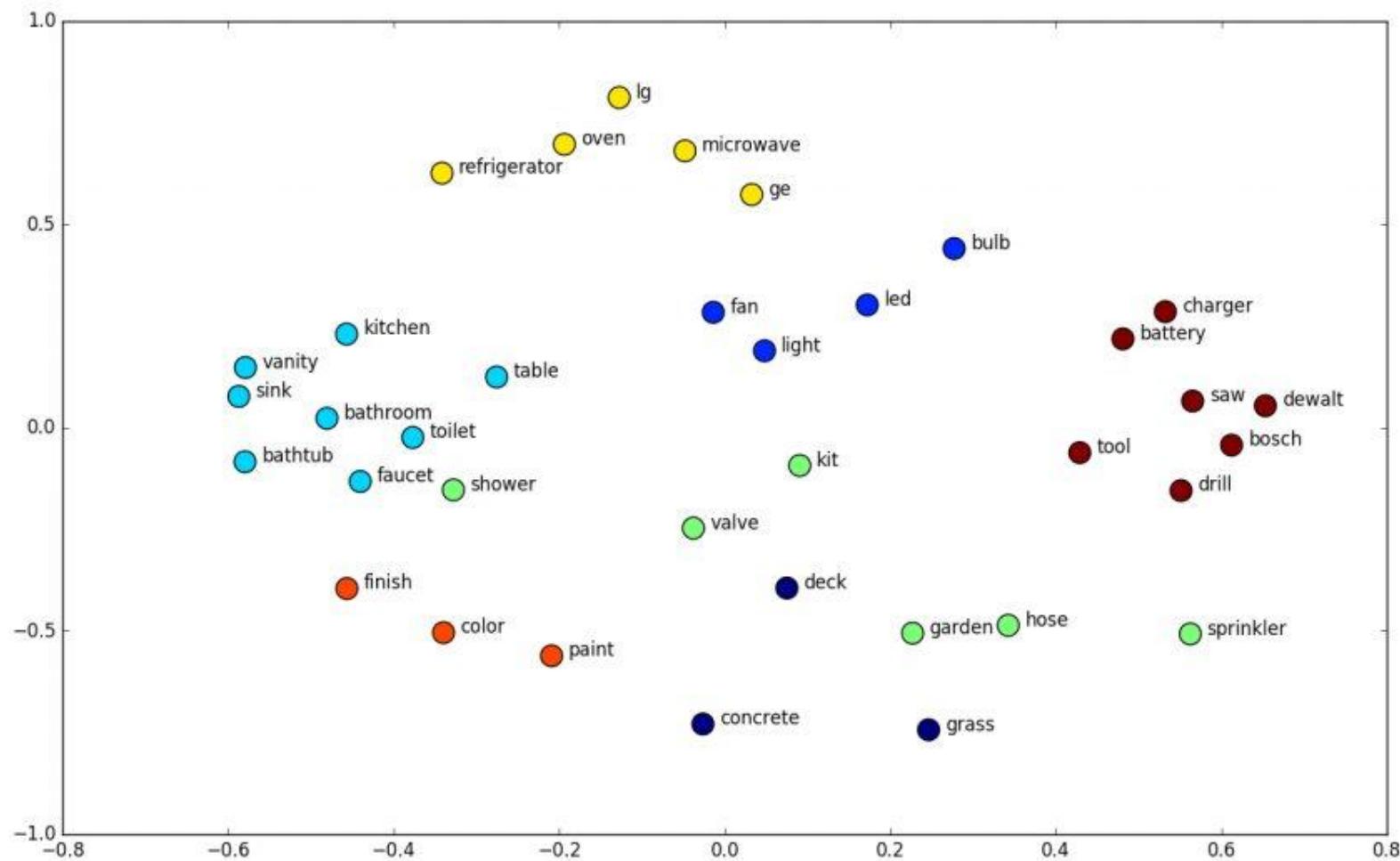


	Femininity	Youth	Royalty
Man	0	0	0
Woman	1	0	0
Boy	0	1	0
Girl	1	1	0
Prince	0	1	1
Princess	1	1	1
Queen	1	0	1
King	0	0	1
Monarch	0.5	0.5	1

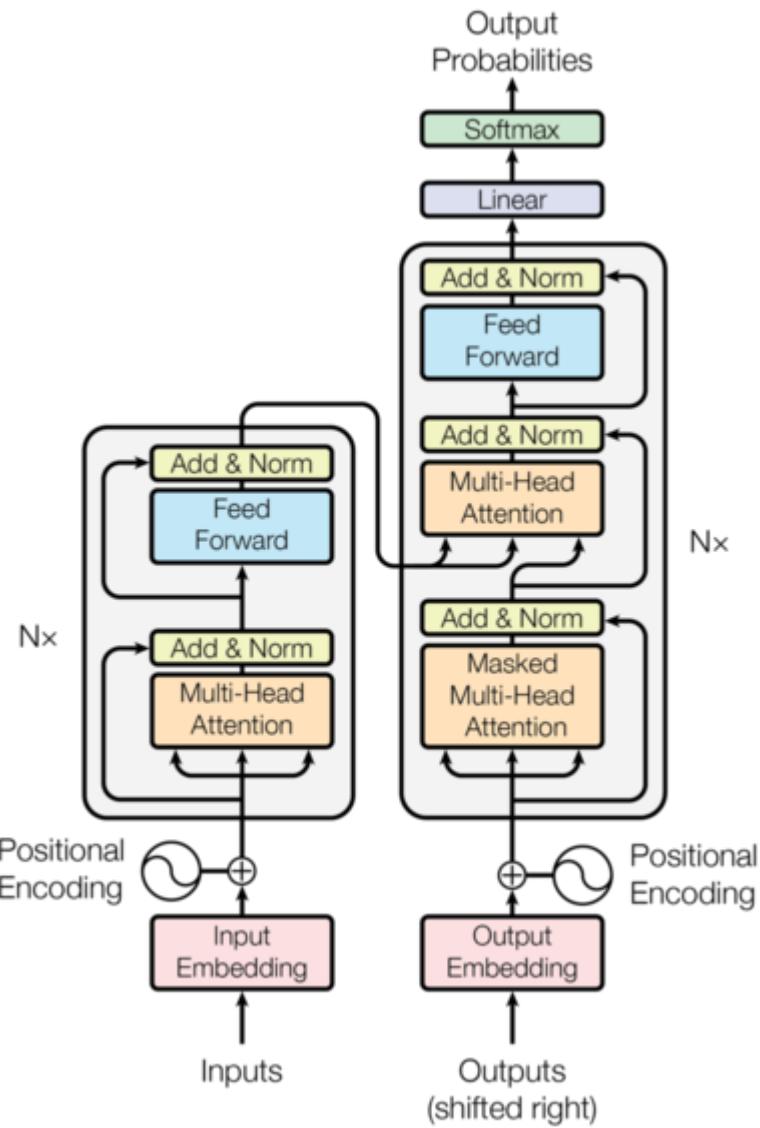
Each word gets a
1x3 vector

Similar words...
similar vectors

Word embedding



Word embedding



OpenAI Platform

language model, and the total count of tokens in this piece of text.

GPT-5.x & O1/3

GPT-4 & GPT-3.5 (legacy)

GPT-3 (legacy)

I love chicken and beer.

Clear

Show example

Tokens

6

Characters

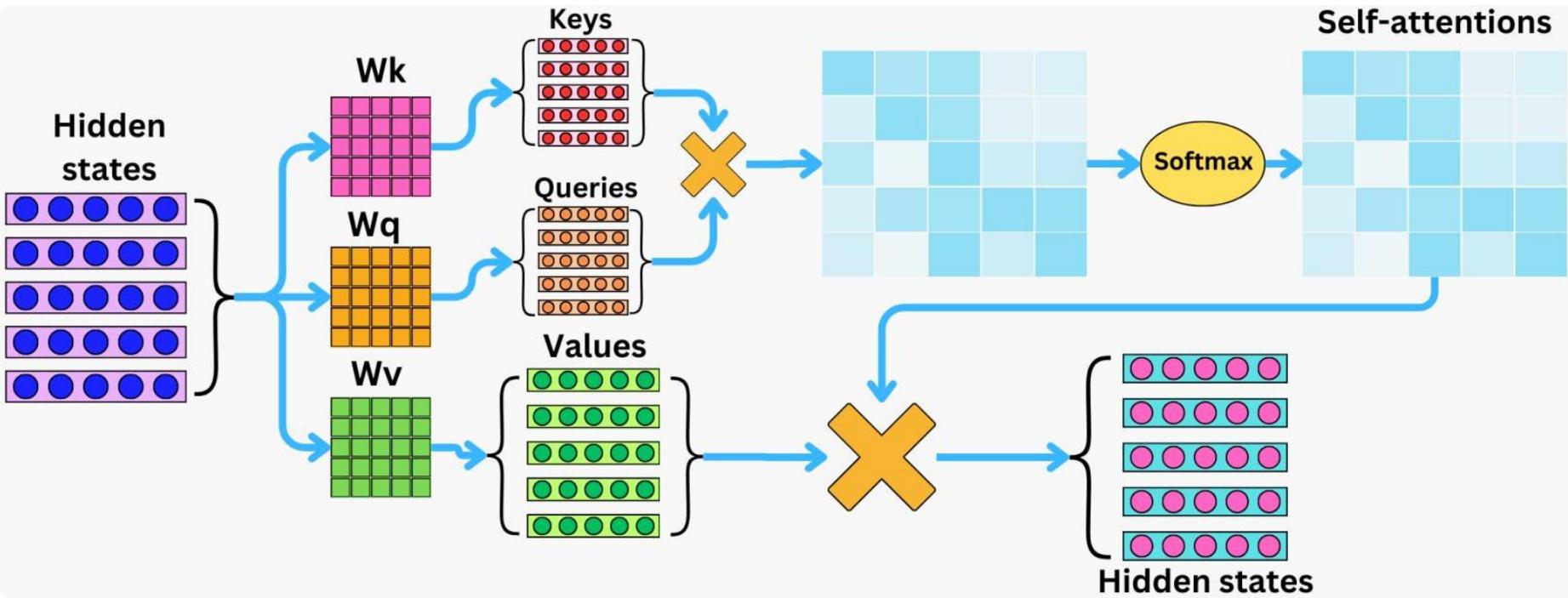
24

[40, 3047, 21663, 326, 20696, 13]

Text

Token IDs

Self attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

배를 먹었더니 배가 아프다

Self attention

```
import torch
import torch.nn as nn
import torch.nn.functional as F

class SelfAttention(nn.Module):
    def __init__(self, embed_dim, atten_dim):
        super().__init__()
        self.query = nn.Linear(embed_dim, atten_dim, bias=False)
        self.key = nn.Linear(embed_dim, atten_dim, bias=False)
        self.value = nn.Linear(embed_dim, atten_dim, bias=False)

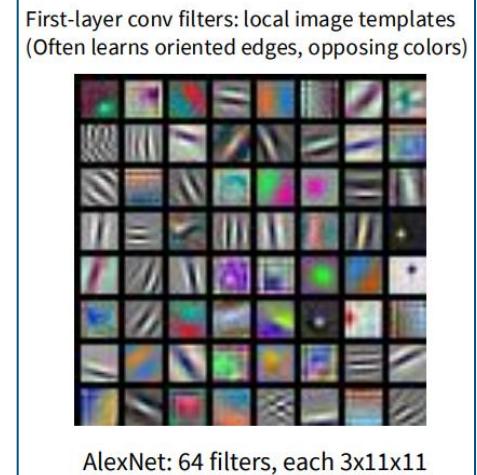
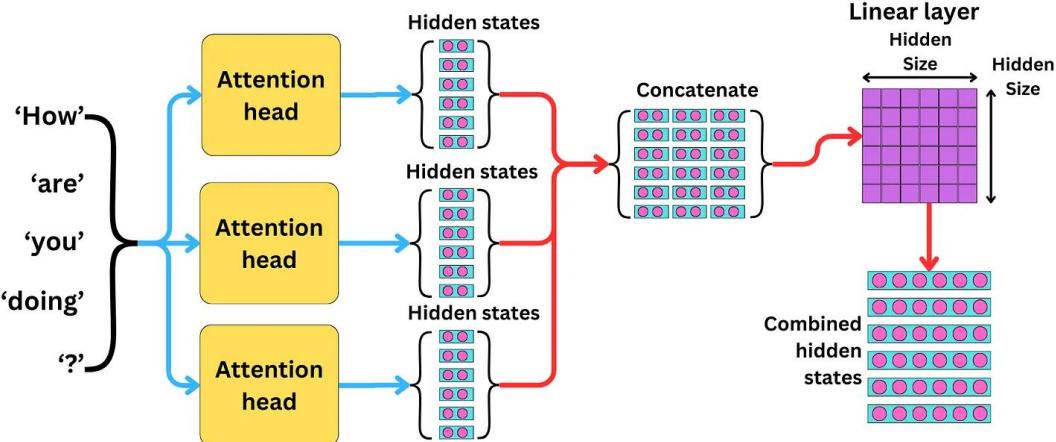
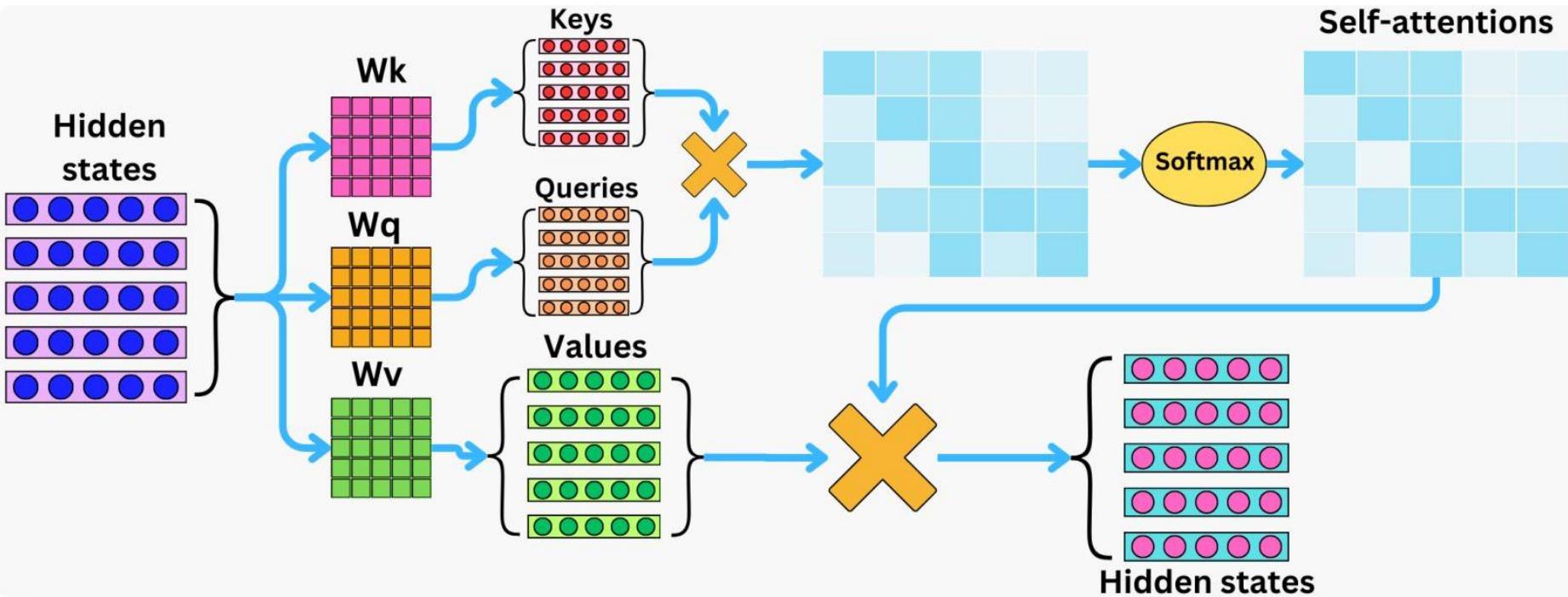
    def forward(self, x):
        query = self.query(x)
        key = self.key(x)
        value = self.value(x)

        scores = torch.matmul(query, key.transpose(-2, -1))
        scores = scores / key.size(-1)**0.5

        attention_weights = F.softmax(scores, dim=-1)
        weighted_values = torch.matmul(attention_weights, value)

        return weighted_values
```

Multi head Self attention



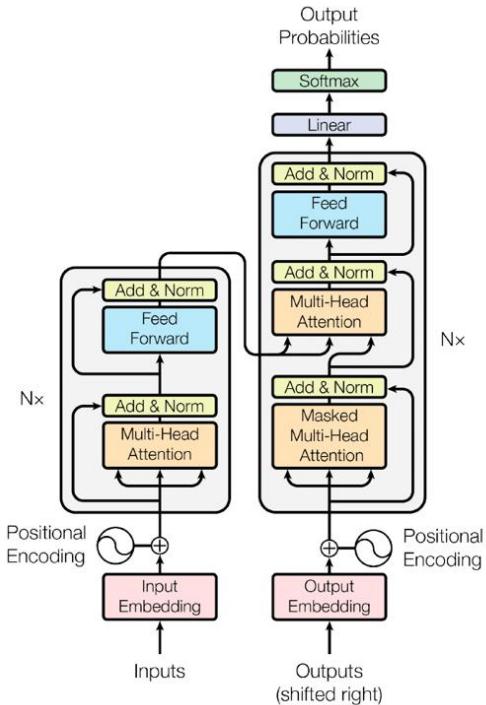
Multi head Self attention

```
class MultiheadAttention(nn.Module):
    def __init__(self, embed_dim, num_heads):
        super().__init__()
        attention_dim = embed_dim // num_heads
        self.attentions = nn.ModuleList([SelfAttention(embed_dim, attention_dim) for _ in range(num_heads)])
        self.fc = nn.Linear(embed_dim, embed_dim)

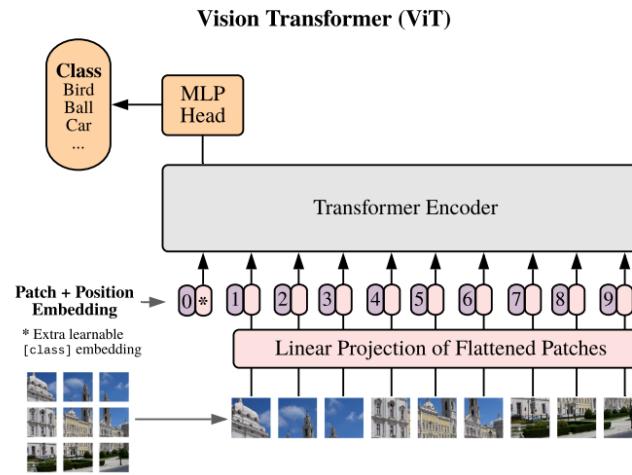
    def forward(self, x):
        head_outputs = []
        for attention in self.attentions:
            head_output = attention(x)
            head_outputs.append(head_output)

        concatenated_heads = torch.cat(head_outputs, dim=-1)
        print("concatenated_heads", concatenated_heads.shape)
        output = self.fc(concatenated_heads)
        print("output", output.shape)
        return output
```

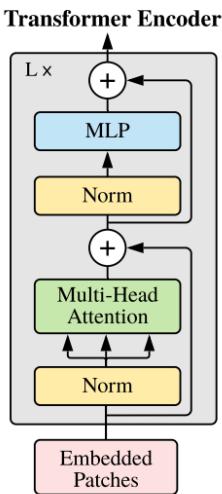
Transformer encoder



Transformer architecture



Vision Transformer architecture



Transformer encoder

```
class FeedForward(nn.Module):
    def __init__(self, embed_dim, ff_dim):
        super().__init__()
        self.net = nn.Sequential(
            nn.Linear(embed_dim, ff_dim),
            nn.ReLU(),
            nn.Linear(ff_dim, embed_dim),
        )
    def forward(self, x):
        return self.net(x)

class TransformerBlock(nn.Module):
    def __init__(self, embed_dim, n_head):
        super().__init__()
        self.layer_norm1 = nn.LayerNorm(embed_dim)
        self.multihead_attn = MultiheadAttention(embed_dim, n_head)

        self.layer_norm2 = nn.LayerNorm(embed_dim)
        self.feed_forward = FeedForward(embed_dim, 4*embed_dim)

    def forward(self, x):
        x = x + self.multihead_attn(self.layer_norm1(x))
        x = x + self.feed_forward(self.layer_norm2(x))
        return x
```

Transformer encoder

v2.10.0 (stable) ▾

[Home](#)[Install PyTorch](#)[User Guide](#)[Reference API](#)[Developer Notes](#)[Community](#)[Tutorials](#)

Search the docs ...

GRU

RNNCell

LSTMCell

GRUCell

Transformer

TransformerEncoder

TransformerDecoder

TransformerEncoderLayer

TransformerDecoderLayer

Identity

Linear

Bilinear

LazyLinear



> Reference API > torch.nn > TransformerEncoder

Rate this Page

TransformerEncoder

```
class torch.nn.TransformerEncoder(encoder_layer, num_layers, norm=None,  
enable_nested_tensor=True, mask_check=True) #
```

[\[source\]](#)

TransformerEncoder is a stack of N encoder layers.

This TransformerEncoder layer implements the original architecture described in the [Attention Is All You Need](#) paper.

The intent of this layer is as a reference implementation for foundational understanding and thus it contains only limited features relative to newer Transformer architectures. Given the fast pace of innovation in transformer-like architectures, we recommend exploring this [tutorial](#) to build efficient layers from building blocks in core or using higher level libraries from the [PyTorch Ecosystem](#).

Examples

```
>>> encoder_layer = nn.TransformerEncoderLayer(d_model=512, nhead=8)  
>>> transformer_encoder = nn.TransformerEncoder(encoder_layer, num_layers=6)  
>>> src = torch.rand(10, 32, 512)  
>>> out = transformer_encoder(src)
```

Thank you for listening 😊