

**CS2109S: Introduction to AI and Machine Learning**

# Lecture 10: Introduction to Deep Learning

3 November 2023

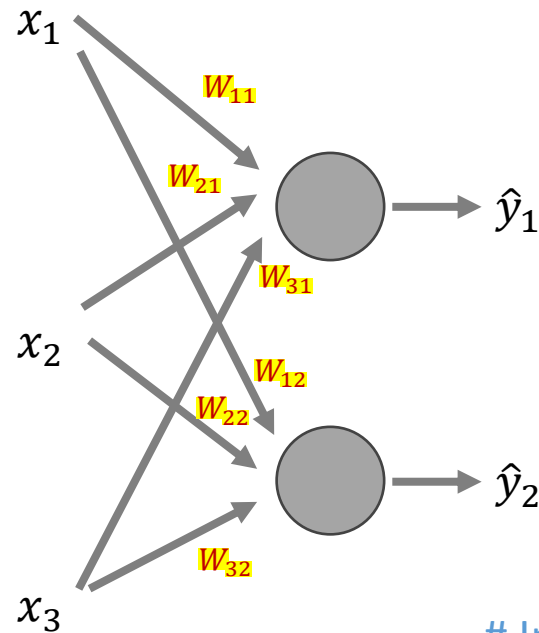
# Announcement

- Midterm grading is still on progress, sorry :(
  - Q1 is quite painful to grade
  - Probably will be done by **next week**...
- We'll release the mockup final assessment next week (hopefully)

# Recap

- Backpropagation
  - Backpropagation on different scenarios:
    - Path, branches, many features, and many samples (sum the gradients)
  - Biological plausibility of backpropagation – believed to be **not feasible**
- Automatic Differentiation
  - Reverse mode automatic differentiation – **backprop is a special case**:  $\mathbb{R}^N \rightarrow \mathbb{R}$
  - ~~• Comparison with other methods: **symbolic** and **numerical** differentiation~~
- Introduction to PyTorch
  - Tensors
    - **n-dimensional array representation** with GPU support
    - Maintain **computational graph**
  - Modules & Functions: Linear (linear), ReLU (relu), etc – they are equivalent
  - Loss function & Optimizers

# Neural Networks and Matrix Multiplication (1)



# Input (number of weights per neuron / input variables)

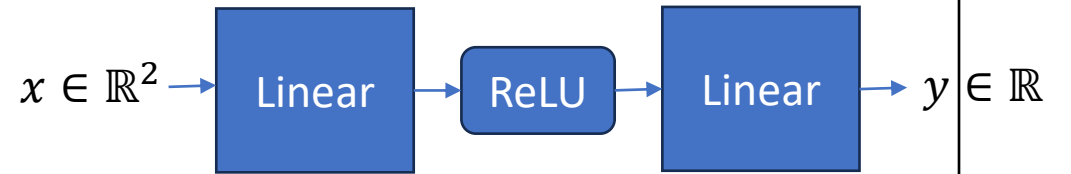
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \\ W_{31} & W_{32} \end{bmatrix} \quad \hat{\mathbf{y}} = g(\mathbf{W}^T \mathbf{x}) = g \left( \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \\ W_{31} & W_{32} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = g \left( \begin{bmatrix} W_{11} & W_{21} & W_{31} \\ W_{12} & W_{22} & W_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

# Output (number of layer's neurons / output variables)

# Modules and Functions API: Example

```
class NeuralNetRegressor(torch.nn.Module):
    def __init__(self, input_size, hidden_size):
        super().__init__()
        self.linear1 = torch.nn.Linear(input_size, hidden_size)
        self.linear2 = torch.nn.Linear(hidden_size, 1)
        self.relu = torch.nn.ReLU()

    def forward(self, x):
        f1 = self.linear1(x)
        a1 = self.relu(f1)
        f2 = self.linear2(a1)
        return f2
```



```
w1 = torch.tensor(8, 2, requires_grad=True)
w2 = torch.tensor(1, 8, requires_grad=True)

def neural_net_regressor(x): # also the same
    f1 = torch.nn.functional.linear(x, w1)
    a1 = torch.nn.functional.relu(f1)
    return torch.nn.functional.linear(a1, w2)
```

```
model1 = NeuralNetClassifier(2, 8) # 2 features, 8 hidden neurons
model2 = torch.nn.Sequential(torch.nn.Linear(2, 8), torch.nn.ReLU(), torch.nn.Linear(8, 1)) # same
```

# Gradient Descent

- Start at some  $w$
- Pick a nearby  $w$  that reduces  $J(w)$
- Repeat until minimum is reached

$$w_j \leftarrow w_j - \gamma \frac{\partial J(w_0, w_1, \dots)}{\partial w_j}$$

**Single-layer** Neural Networks with Sigmoid

$$\gamma(\hat{y} - y)\hat{y}(1 - \hat{y})x_i$$

Derive manually!

**Multi-layer** Neural Networks

Backpropagation!

## Example 2: Training NN using modular API

```
model = NeuralNetRegressor()

optimizer = torch.optim.Adam(model.parameters(), lr=0.01)
loss_function = torch.nn.MSELoss()

for epoch in range(num_epochs):
    optimizer.zero_grad()
    y_pred = model(x)
    loss = loss_function(y_pred, y)
    loss.backward()
    optimizer.step()
```

### Gradient Descent

- Start at some  $w$
- Pick a nearby  $w$  that reduces  $J(w)$   
$$w_j \leftarrow w_j - \gamma \frac{\partial J(w_0, w_1, \dots)}{\partial w_j}$$
- Repeat until minimum is reached

# Outline

- Deep Neural Networks
- Convolution Neural Networks
  - Motivation: handling spatial structure
  - Convolution, Pooling Layer, and Common Architectures
  - Applications
- Recurrent Neural Networks
  - Motivation: handling sequential data
  - Recurrent Neural Networks and Variants
  - Applications
- Attention, Transformers, GPT, and ChatGPT (if time permits)
- Issues with Deep Learning



# Outline

- **Deep Neural Networks**
- Convolution Neural Networks
  - Motivation: handling spatial structure
  - Convolution, Pooling Layer, and Common Architectures
  - Applications
- Recurrent Neural Networks
  - Motivation: handling sequential data
  - Recurrent Neural Networks and Variants
  - Applications
- Attention, Transformers, GPT, and ChatGPT (if time permits)
- Issues with Deep Learning

Looking Back...

# What is AI?



# ChatGPT

Deep Learning System



Credit: IEEE Spectrum



Credit: Guardian



Credit: NYTimes

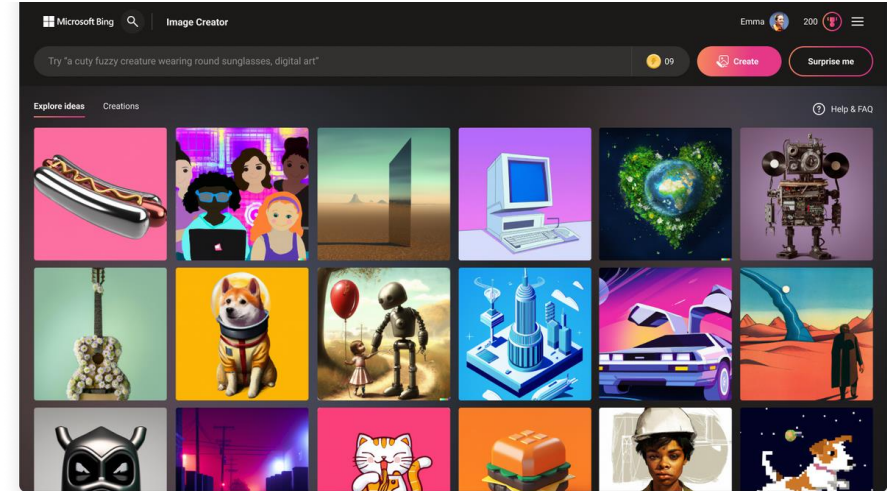
Looking Back...

# What is AI?

## Deep Learning System



Credit: Tesla

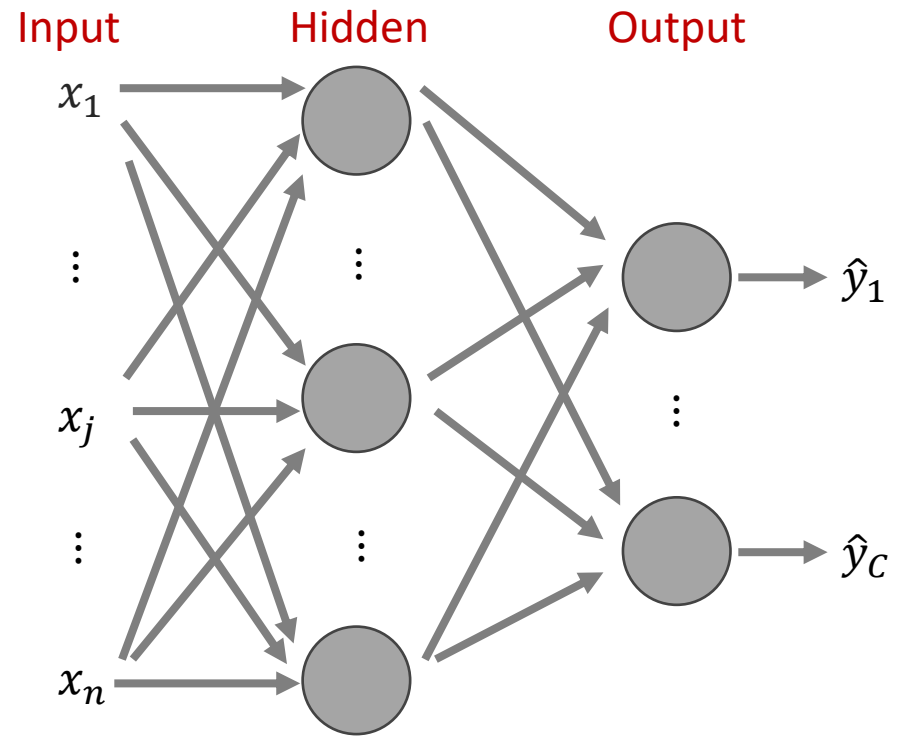


Credit: Eden AI

Credit: Microsoft

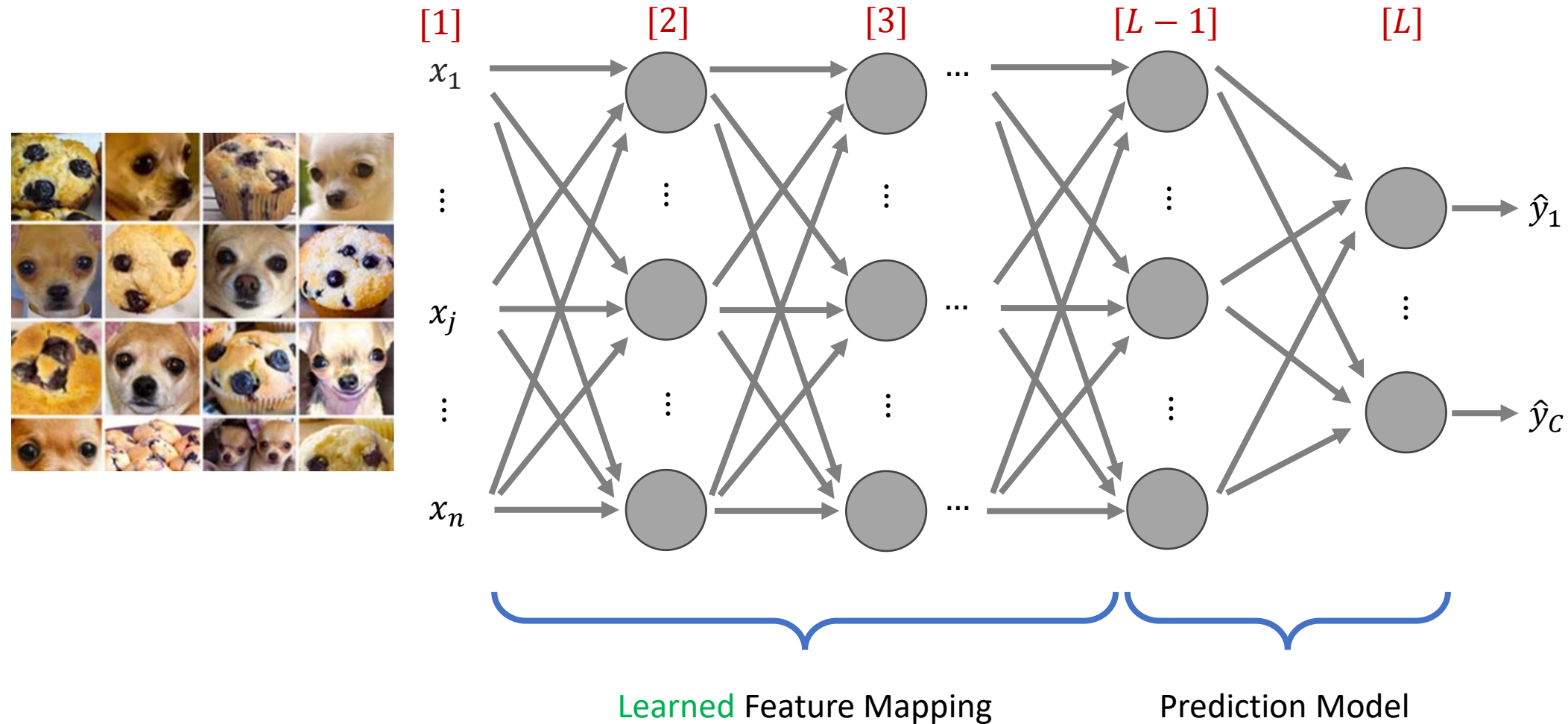


# Shallow Neural Networks



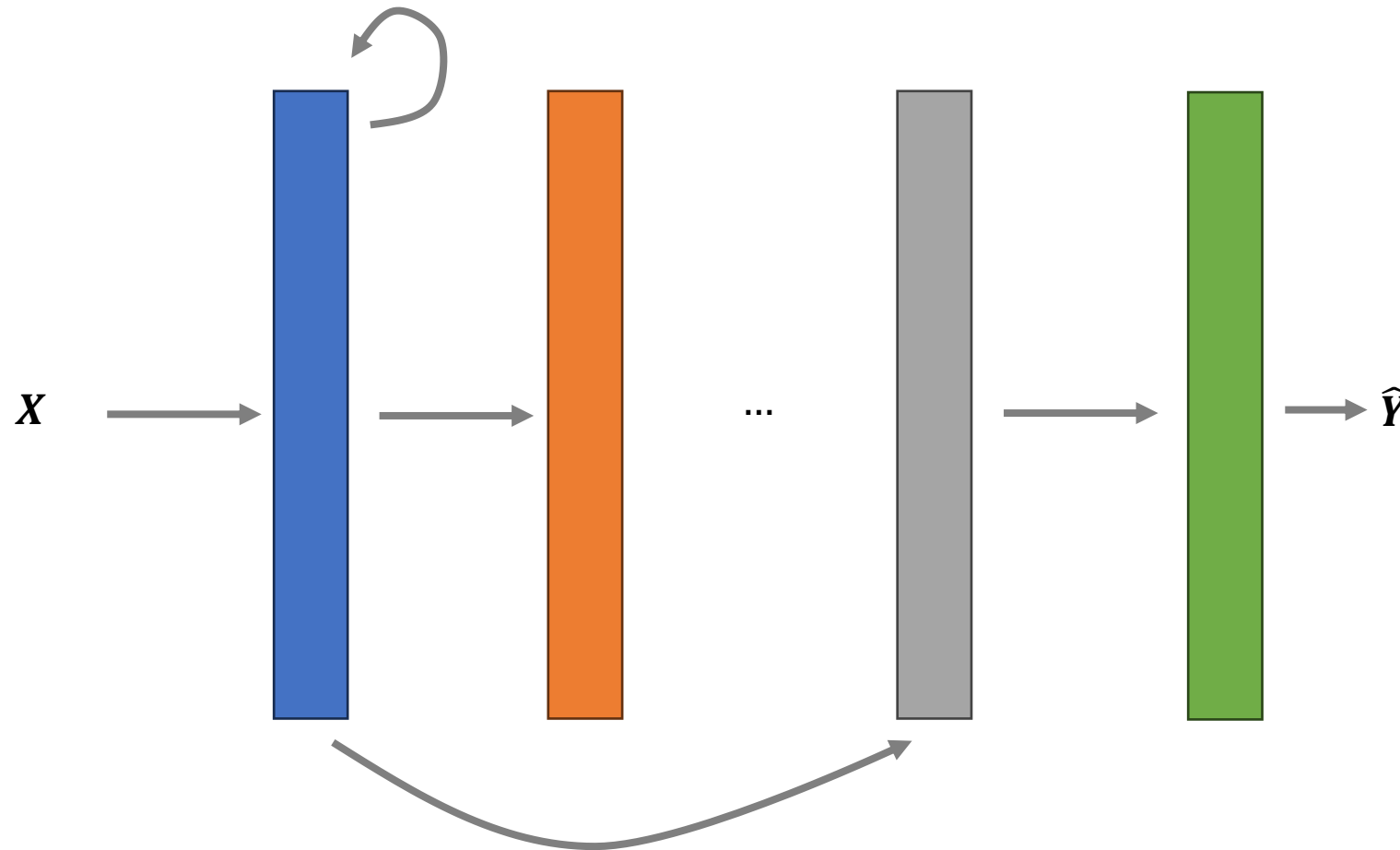
# Deep Neural Networks

Neural Network with the number of layers  $L > 3$



# Deep Neural Networks

Neural Network with the number of layers  $L > 3$



How do we decide  
the functions and  
the compositions?

**Arbitrary function compositions:** can have any\* functions and any\* compositions!

XNOR, XOR

Differentiable

T&C applies

# Outline

- Deep Neural Networks
- **Convolution Neural Networks**
  - Motivation: handling spatial structure
  - Convolution, Pooling Layer, and Common Architectures
  - Applications
- Recurrent Neural Networks
  - Motivation: handling sequential data
  - Recurrent Neural Networks and Variants
  - Applications
- Attention, Transformers, GPT, and ChatGPT (if time permits)
- Issues with Deep Learning

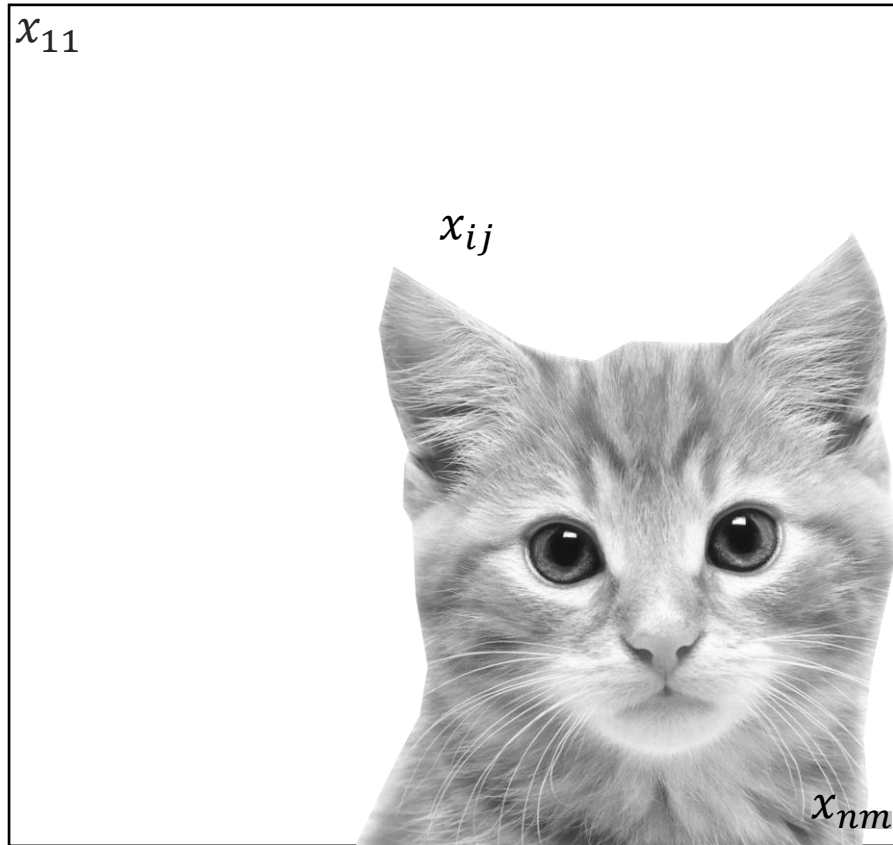
# Computer Vision Problem



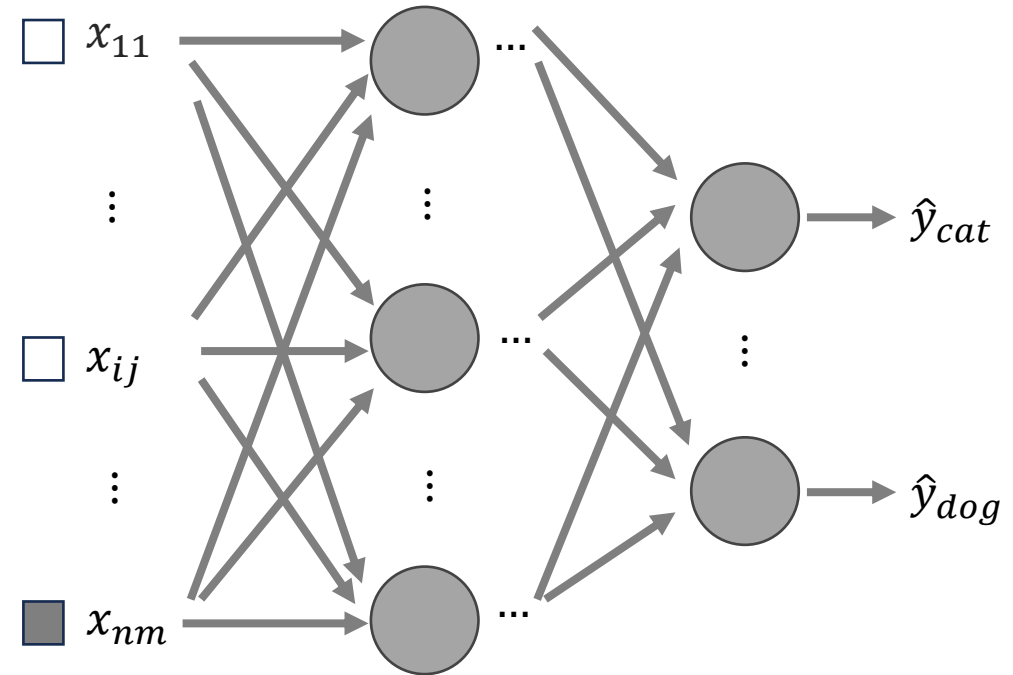
**Cat or Dog?**



# Computer Vision Problem: A Naïve Attempt



Input: Pixels



# Computer Vision Problem: A Naïve Attempt

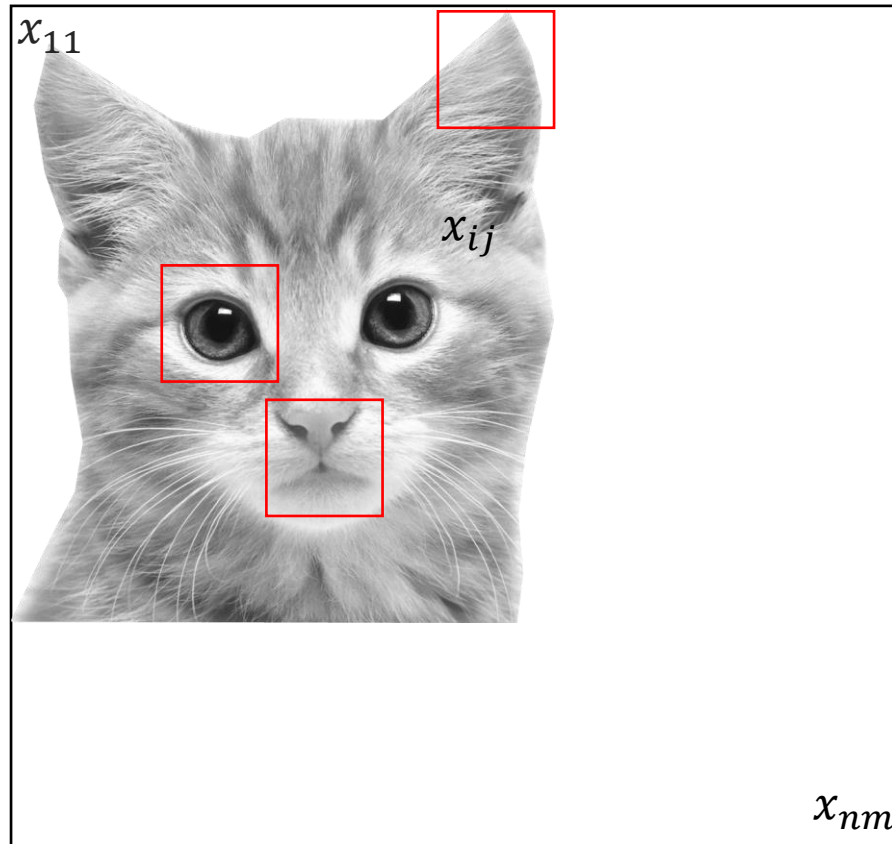
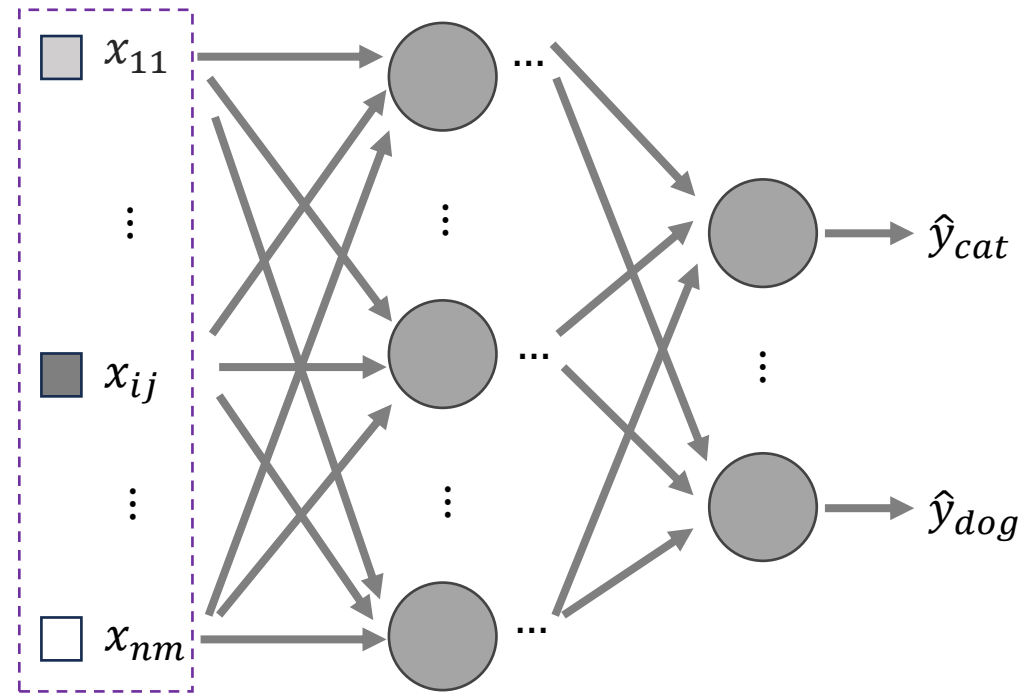


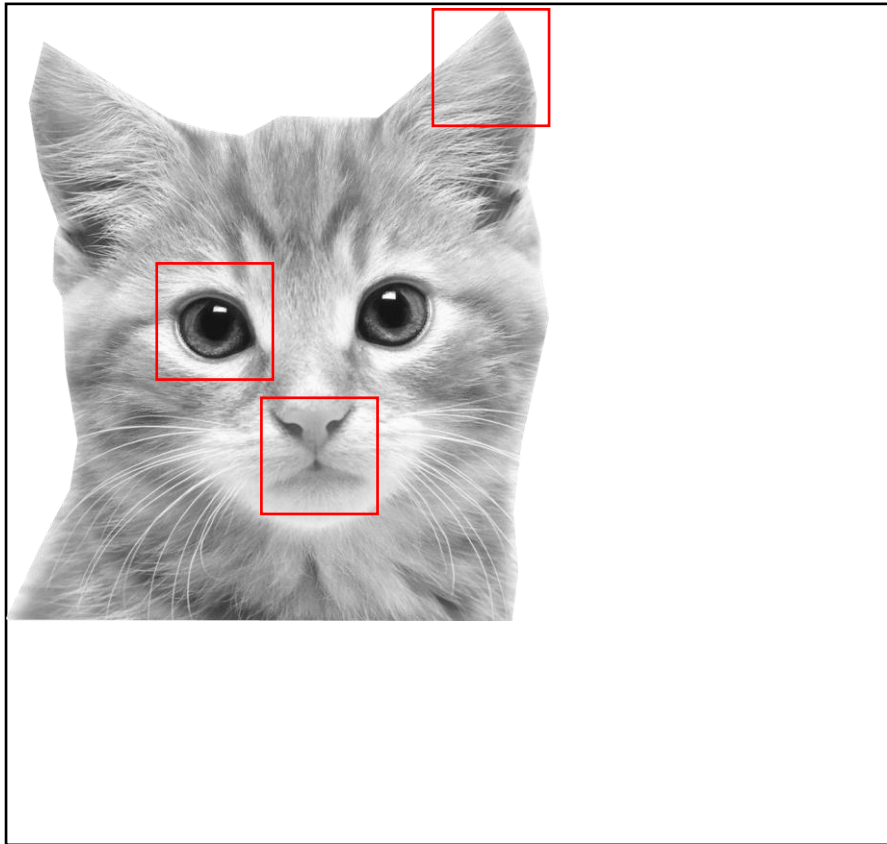
Image has *spatial structures*, which are **ignored**

Input: Pixels

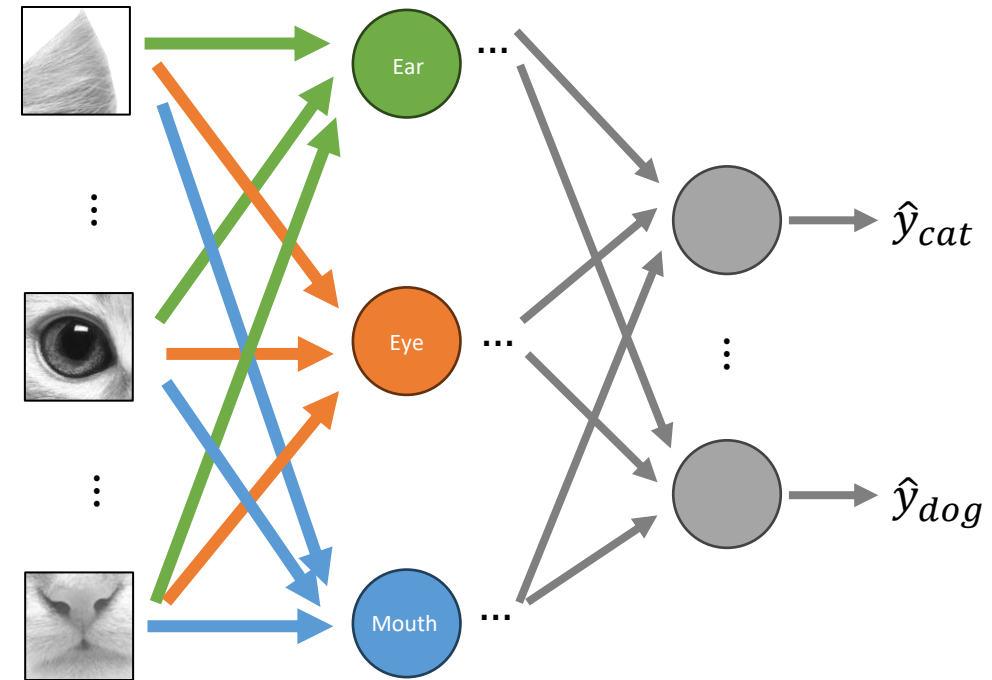


Same image, shifted → **dramatic change** in input of NN

# Computer Vision Problem: A Better Idea



Input: A group of Pixels



**Note:** same colour → same function

# Computer Vision Problem: A Better Idea

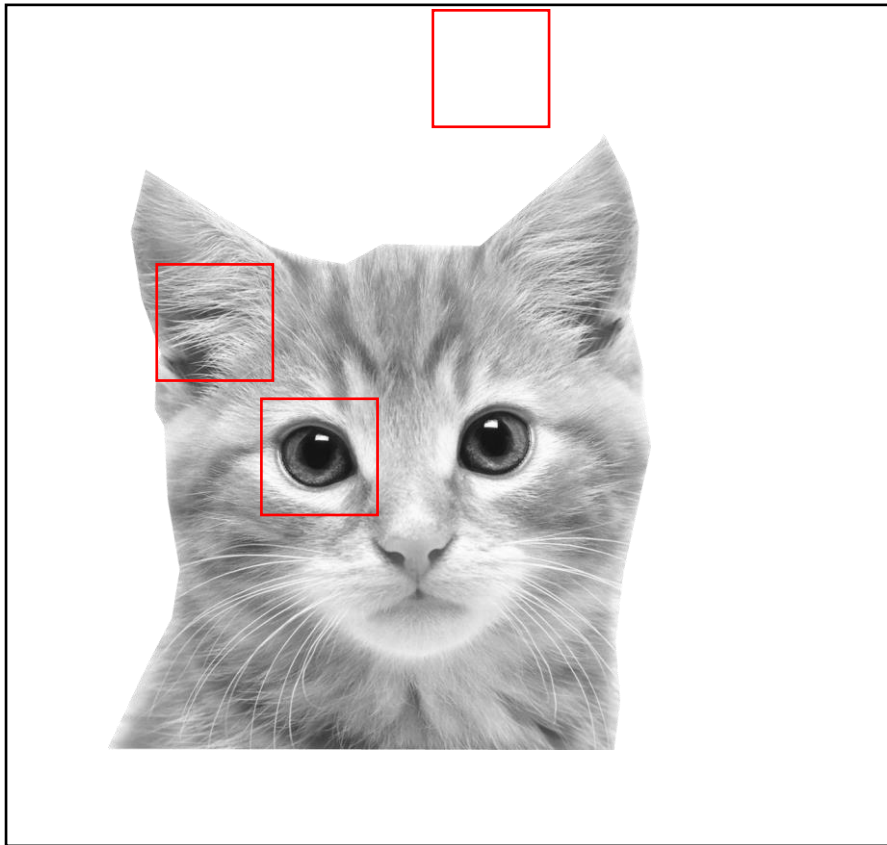
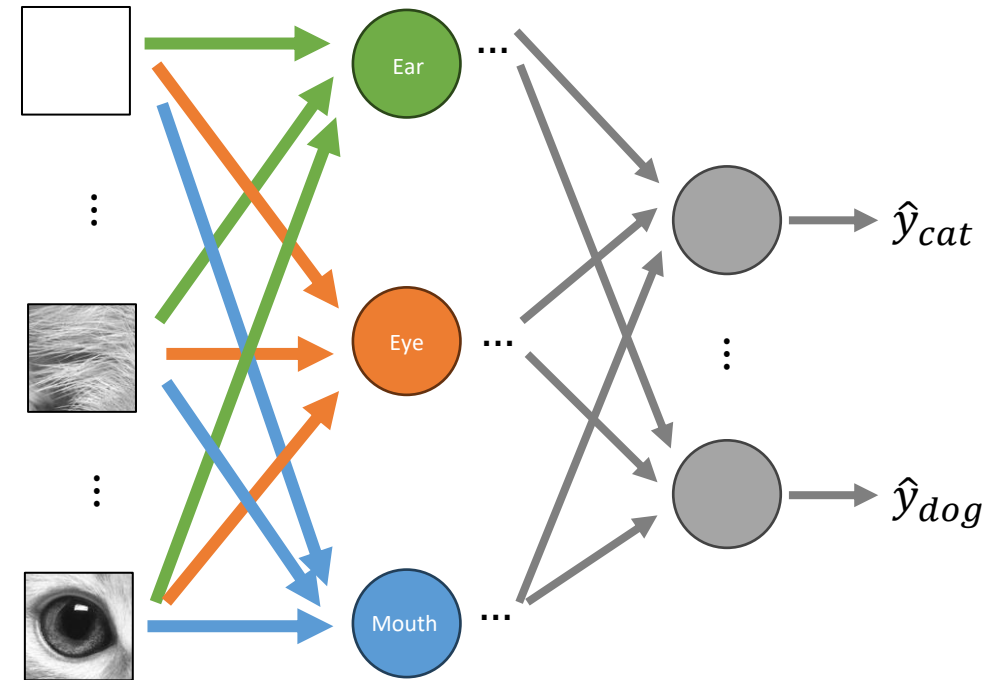


Image has *spatial structures*, which are **preserved**

Input: A group of Pixels



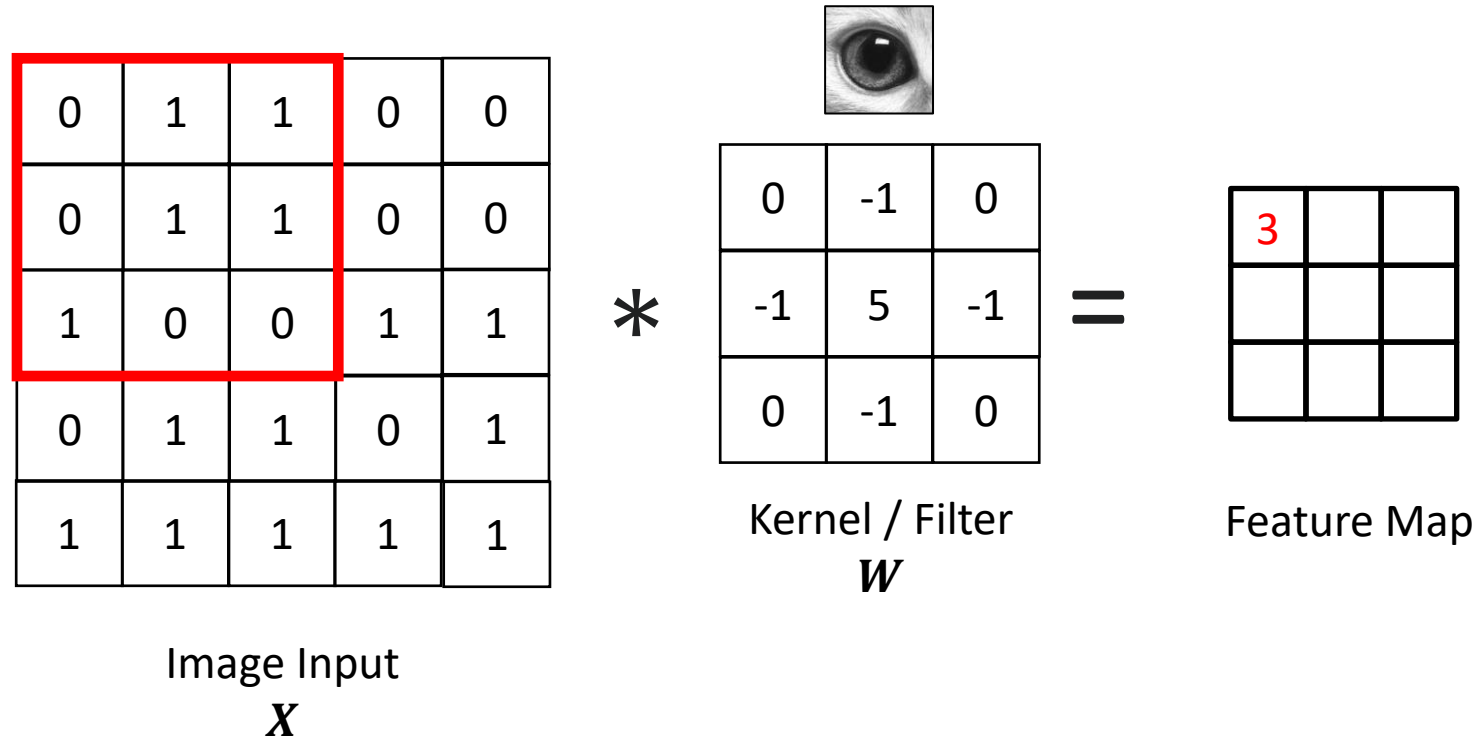
**Note:** same colour  $\rightarrow$  same function

Same image, shifted  $\rightarrow$  **same** set of inputs

How to do this?

# Convolution: 2D

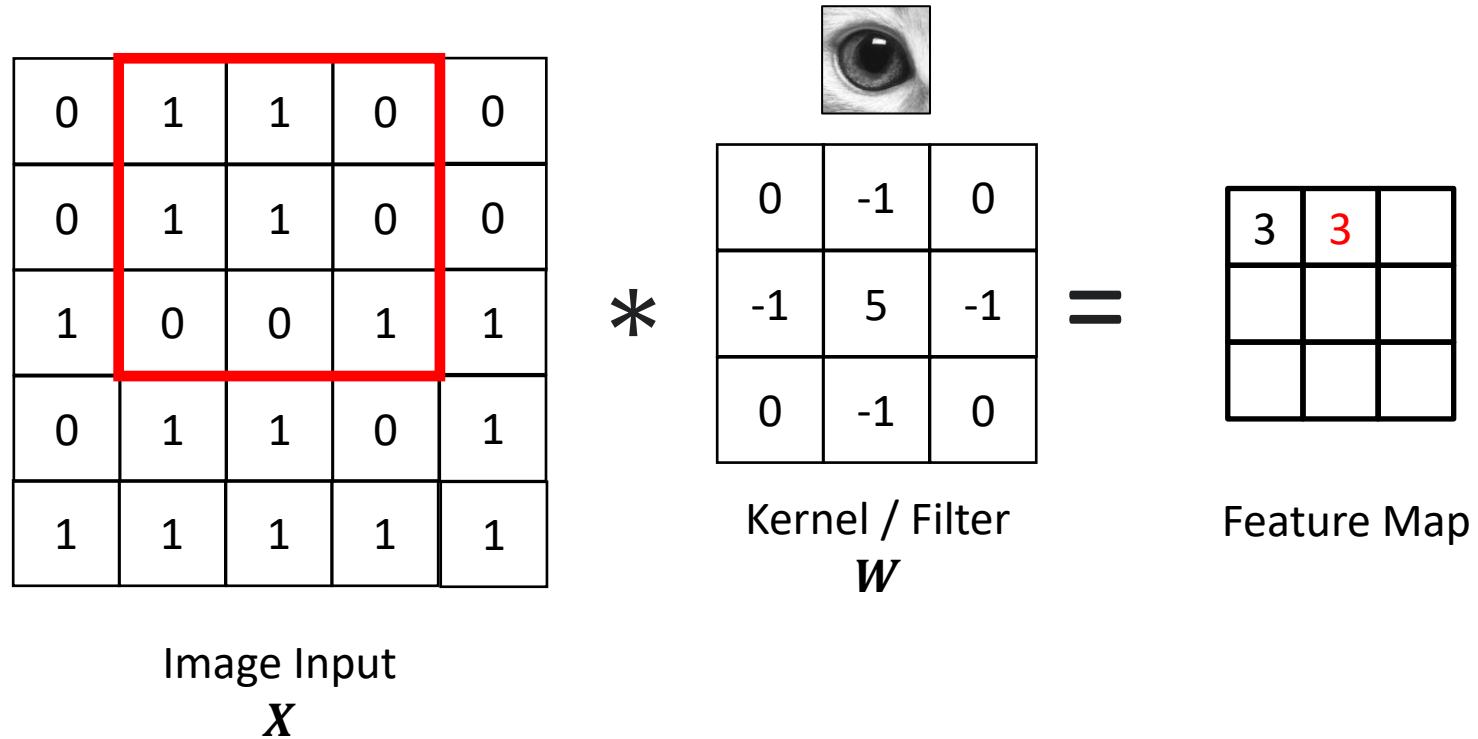
$$f_{conv}(X) = W * X$$



**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

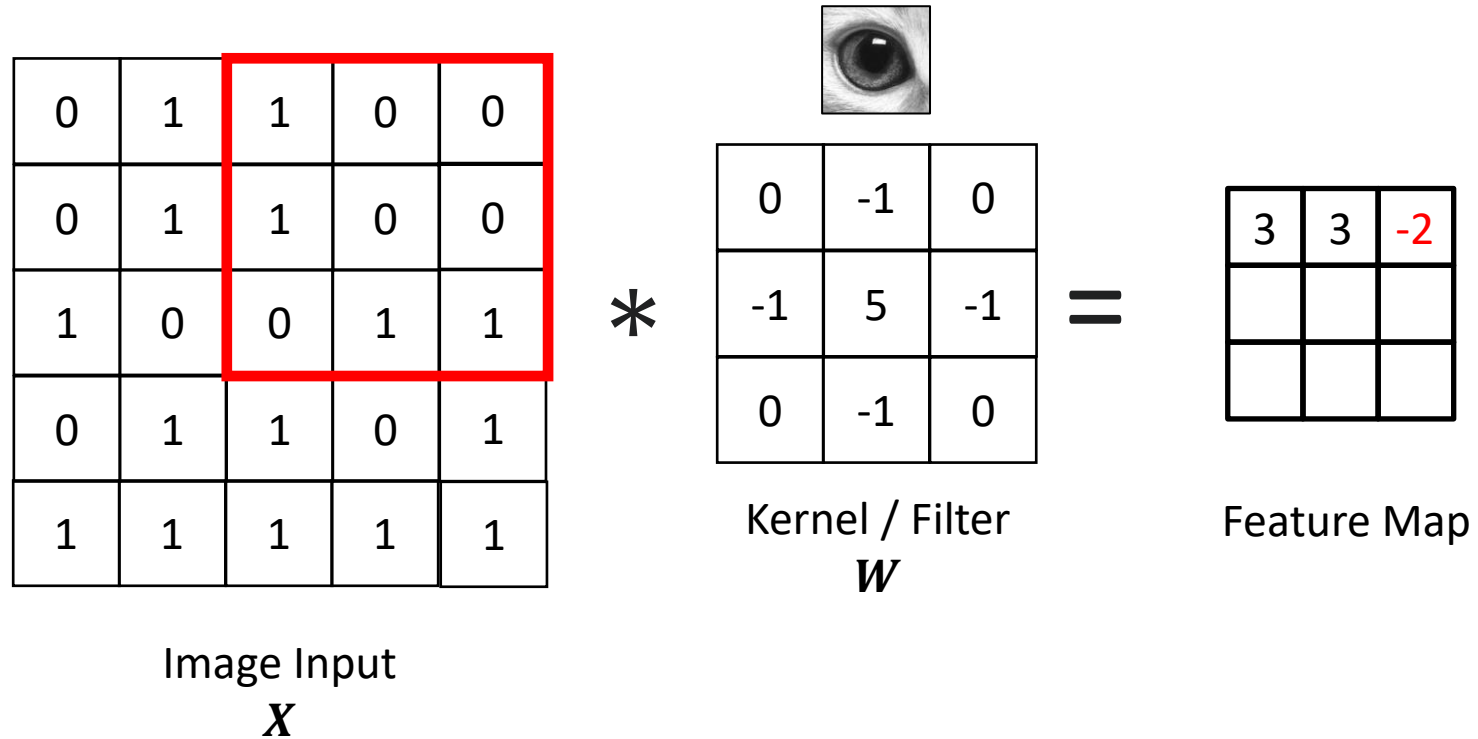
$$f_{conv}(X) = W * X$$



**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

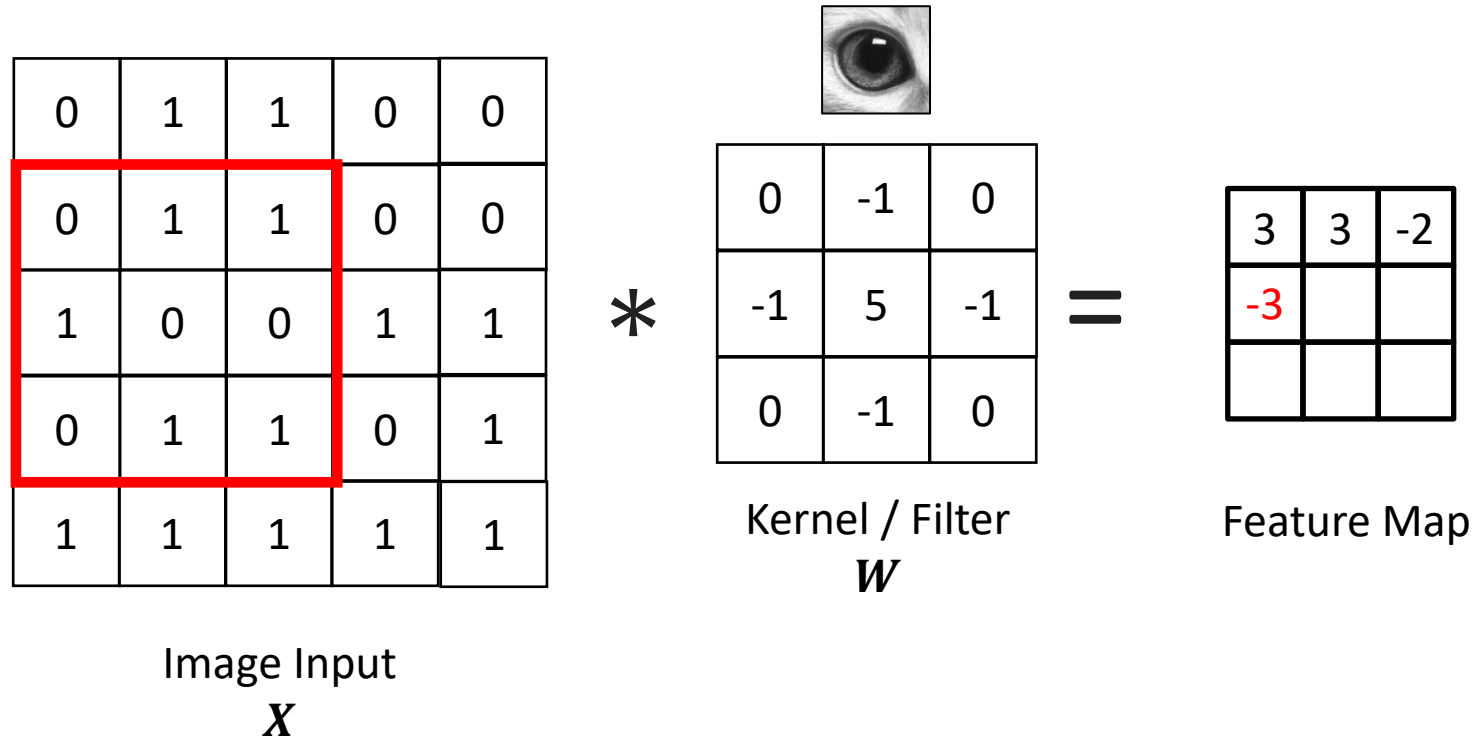
$$f_{conv}(X) = W * X$$



**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

$$f_{conv}(X) = W * X$$

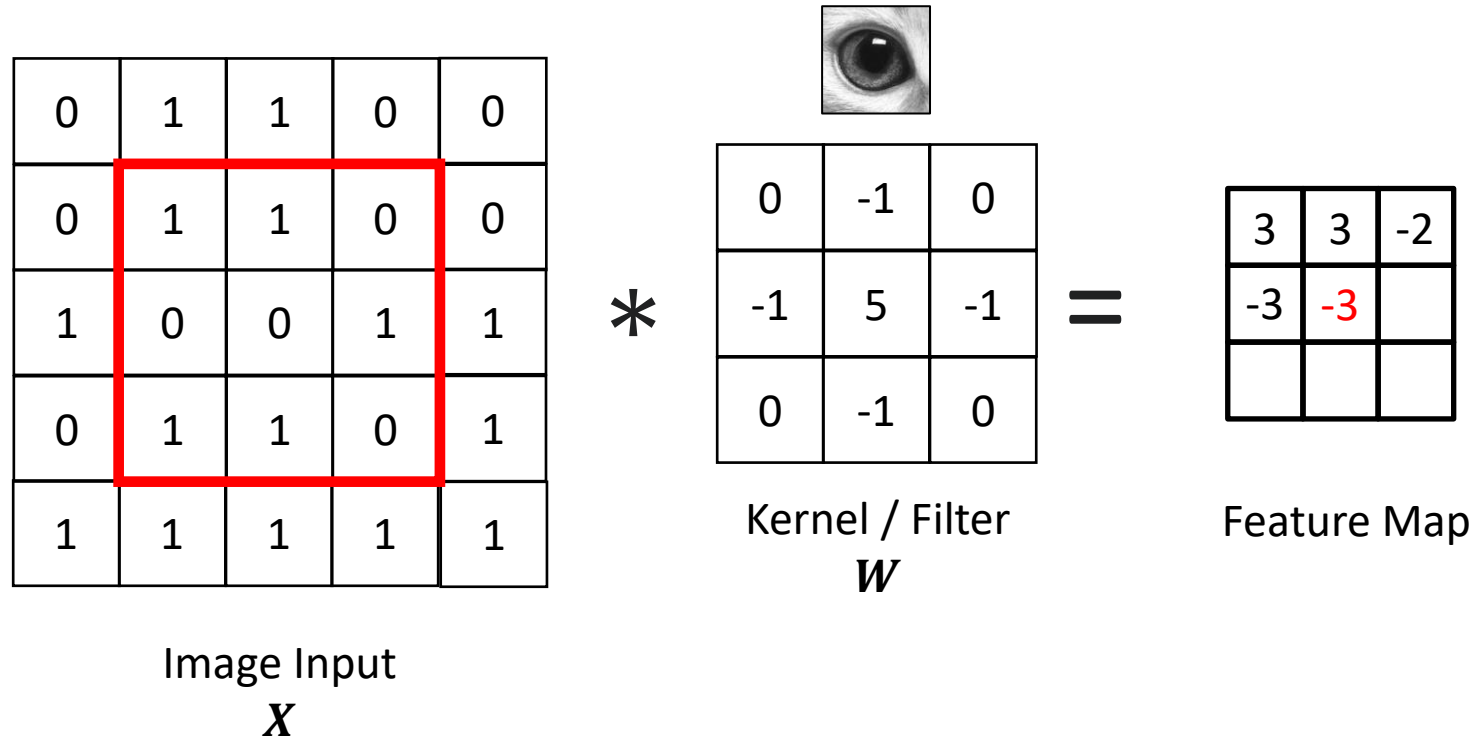


**Multiply** the sliding input window with kernel then **sum**



# Convolution: 2D

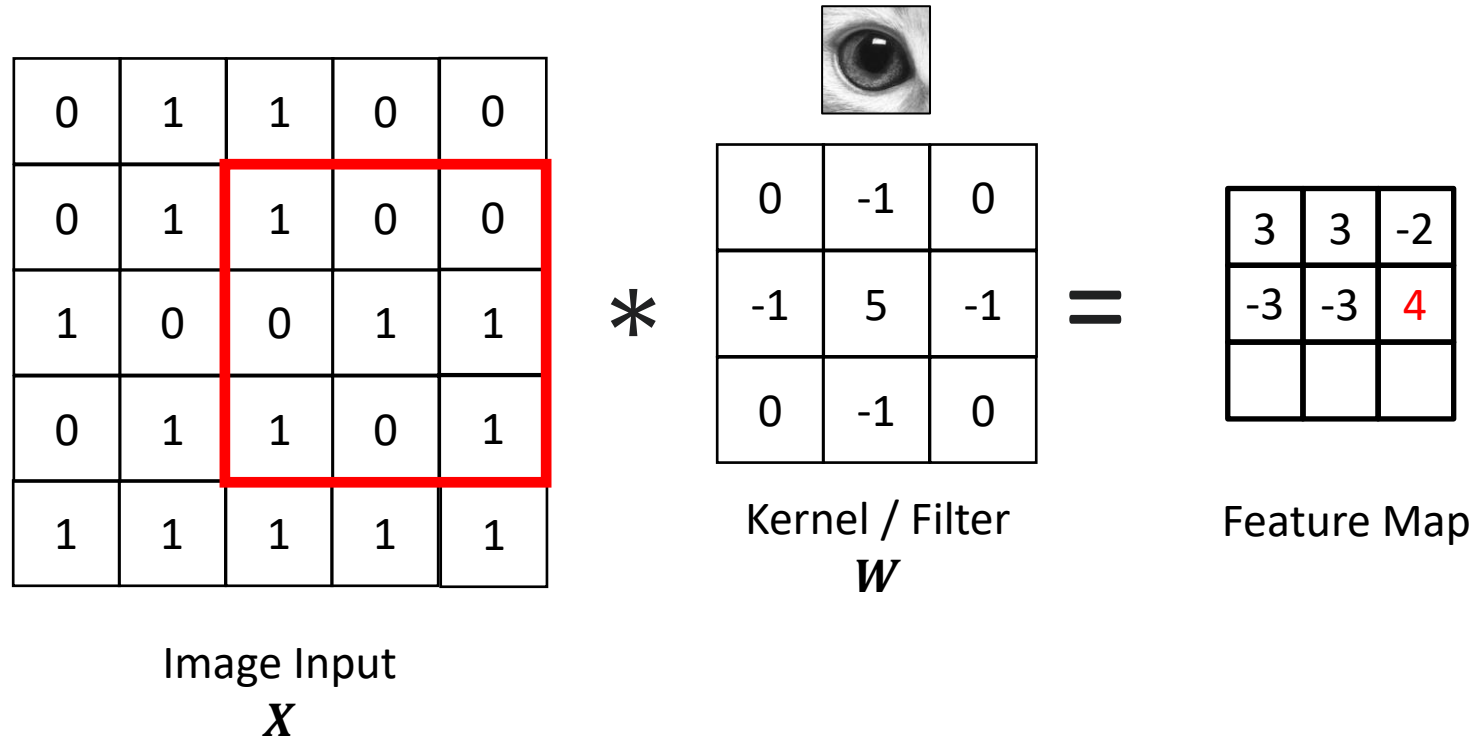
$$f_{conv}(\mathbf{X}) = \mathbf{W} * \mathbf{X}$$



**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

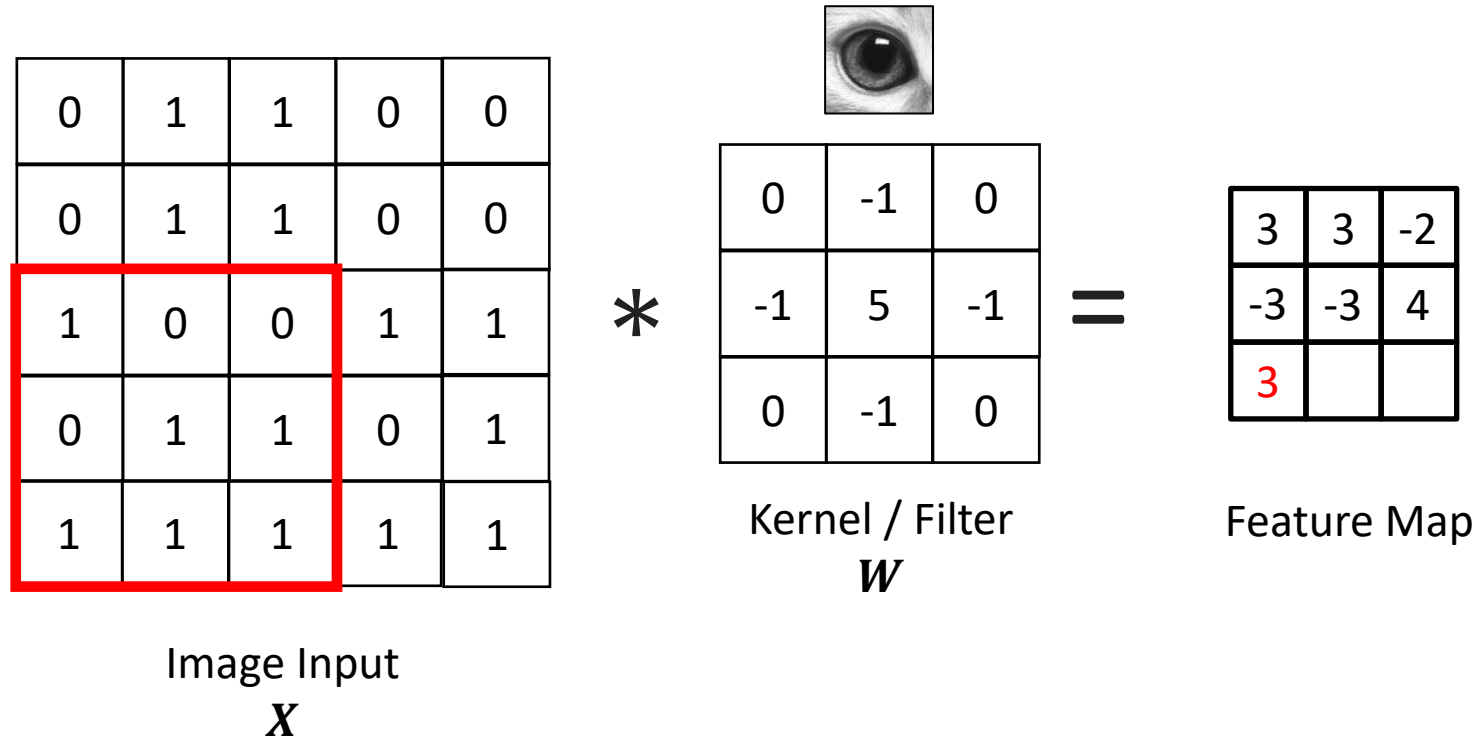
$$f_{conv}(X) = W * X$$



**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

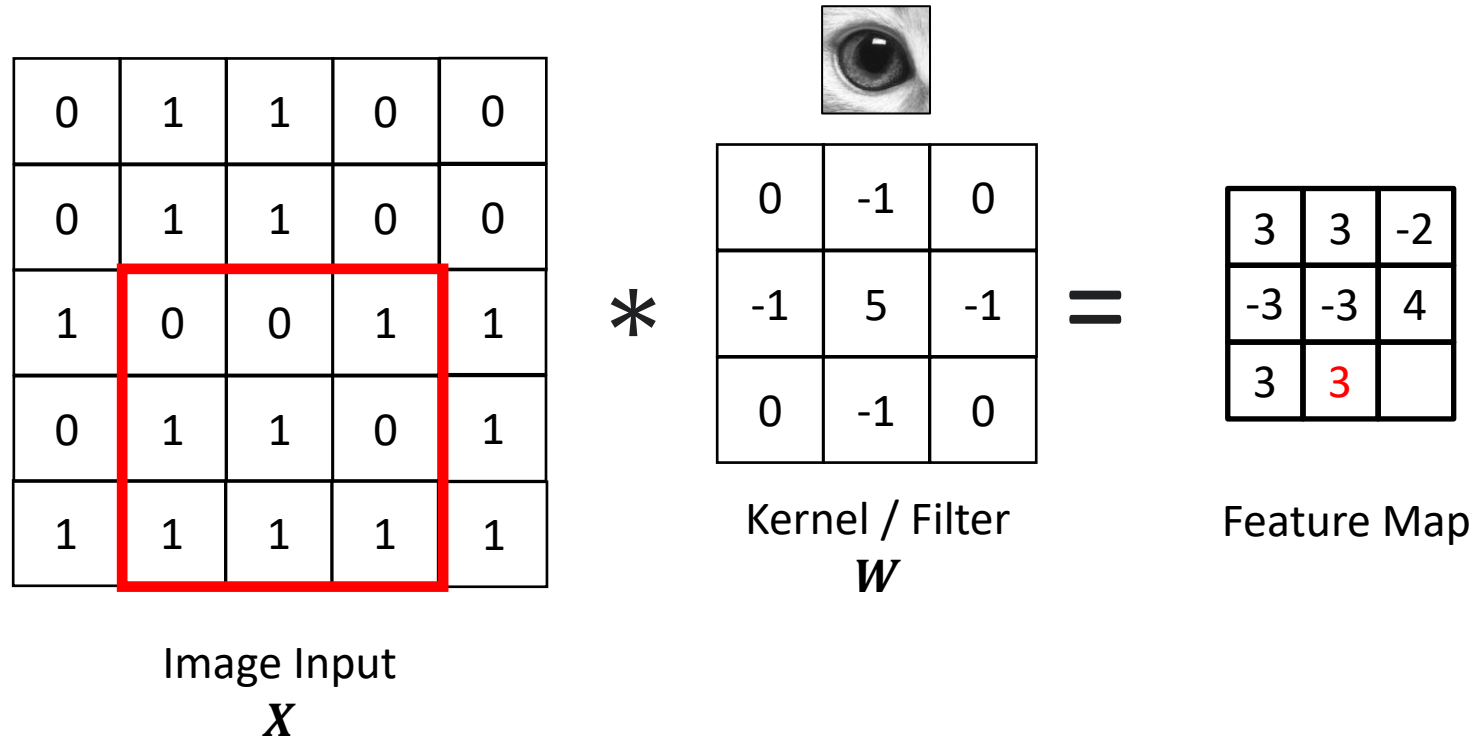
$$f_{conv}(X) = W * X$$



**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

$$f_{conv}(X) = W * X$$



**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

$$f_{conv}(X) = W * X$$

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

Image Input  
 $X$

\*

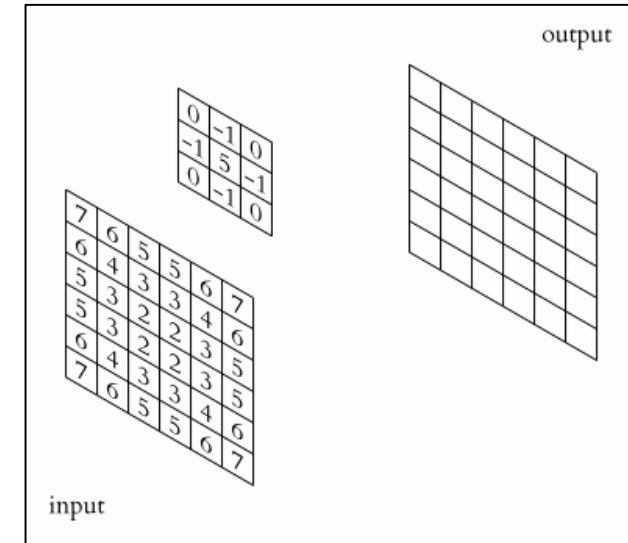
0	-1	0
-1	5	-1
0	-1	0

Kernel / Filter  
 $W$

=

3	3	-2
-3	-3	4
3	3	4

Feature Map



What if we want to detect other features (e.g., mouth)?

**Multiply** the sliding input window with kernel then **sum**

# Convolution: 2D

$$f_{conv}(X) = W * X$$

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

Image Input  
 $X$

\*

0	-1	0
-1	5	-1
0	-1	0

Kernel / Filter  
 $W^{[1]}$

=

3	3	-2
-3	-3	4
3	3	4

Feature Map 1

\*

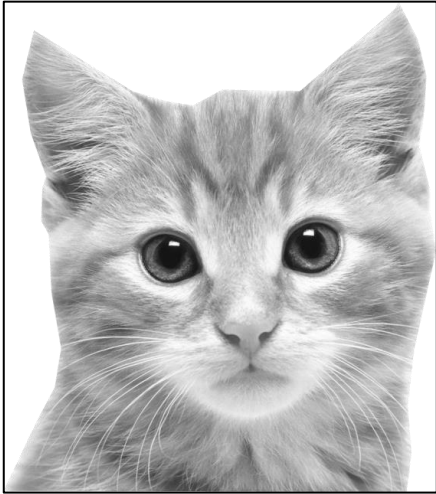
0	1	0
1	0	1
0	1	0

Kernel / Filter  
 $W^{[2]}$

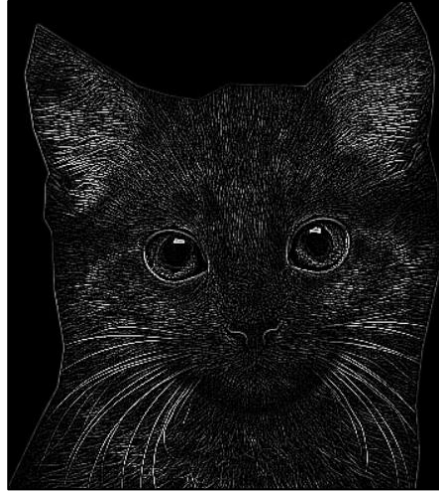
=


Feature Map 2

# Convolution: 2D Example

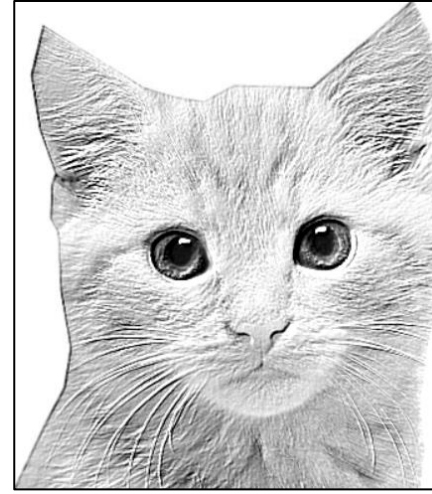


Original



Outline

-1	-1	-1
-1	8	-1
-1	-1	-1



Sobel

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



Blur

.06	.01	.06
.12	.25	.12
.06	.12	.06

**Try it yourself:** <https://setosa.io/ev/image-kernels/>

# Convolution or Cross-correlation?

$$f_{conv}(X) = W * X$$

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

Image Input  
 $X$

\*

1	2	3
4	5	6
7	8	9

Kernel / Filter  
 $W$

=


Feature Map

This is **cross-correlation**!  
Need to flip the kernel



# Convolution or Cross-correlation?

$$f_{conv}(X) = W * X$$

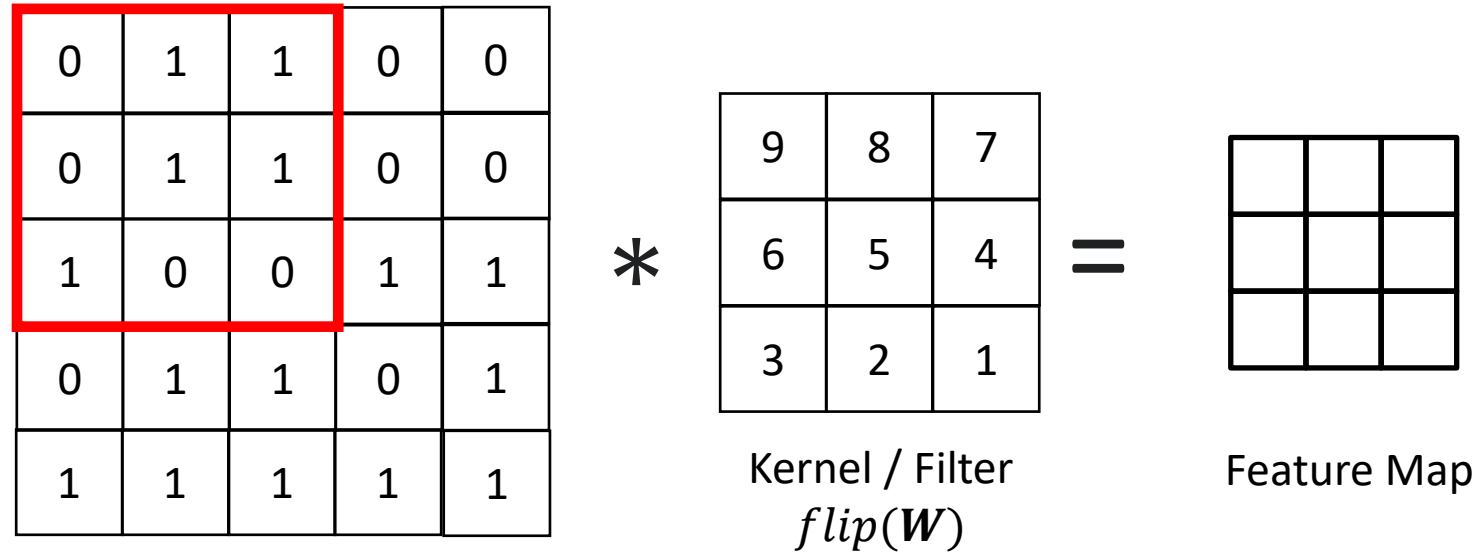
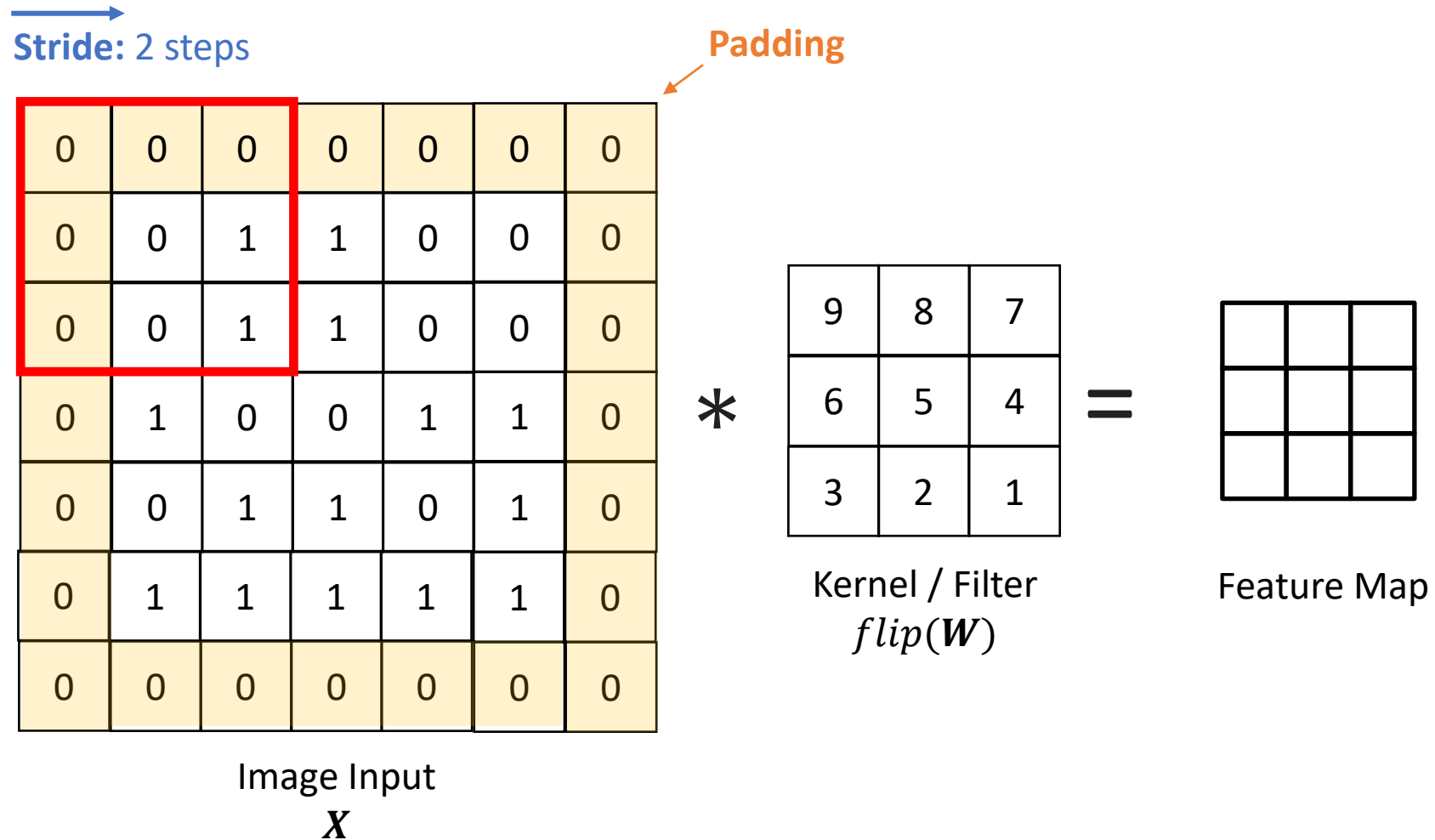


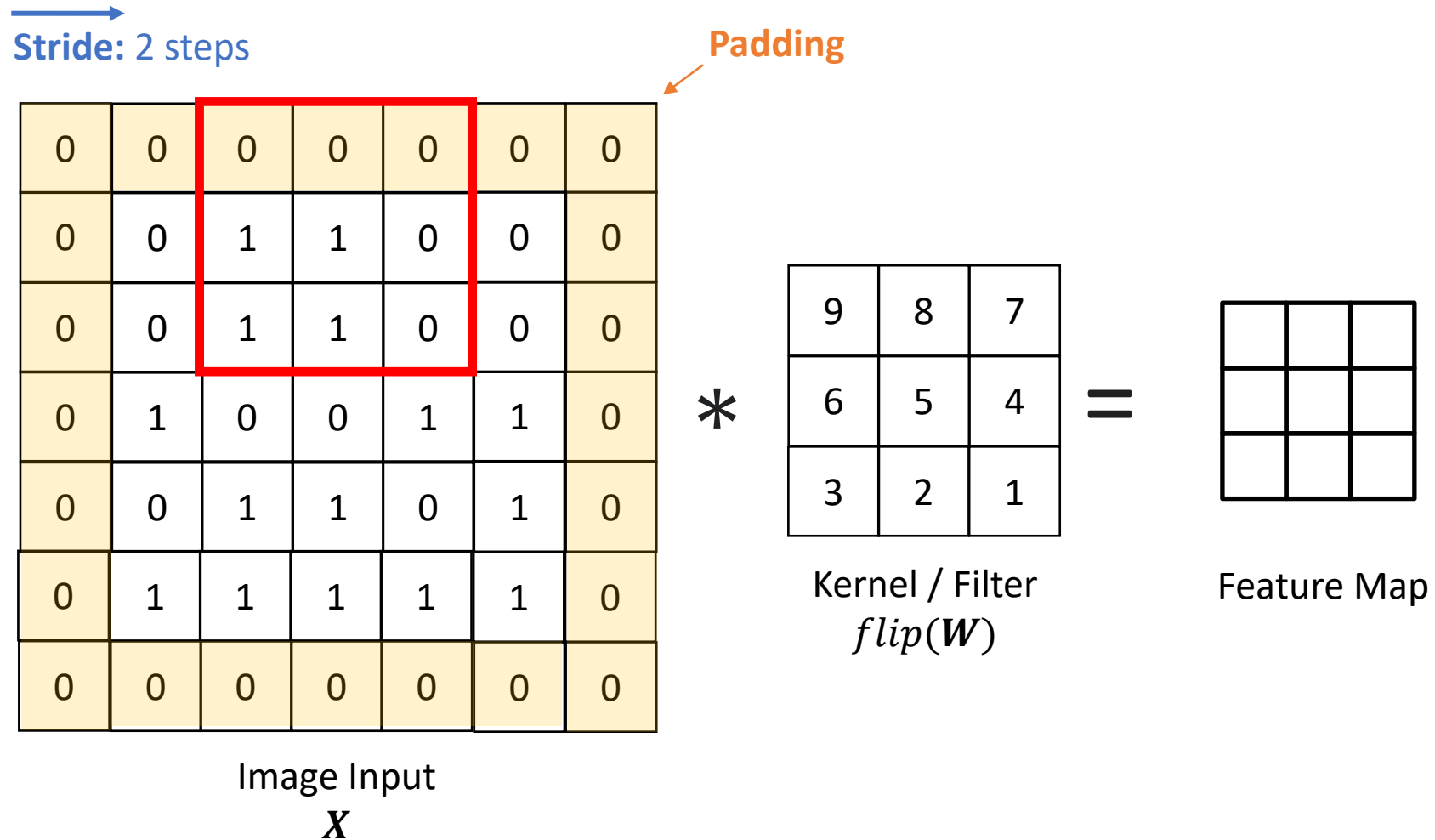
Image Input  
 $X$

This is now **convolution**!  
Does this matter?

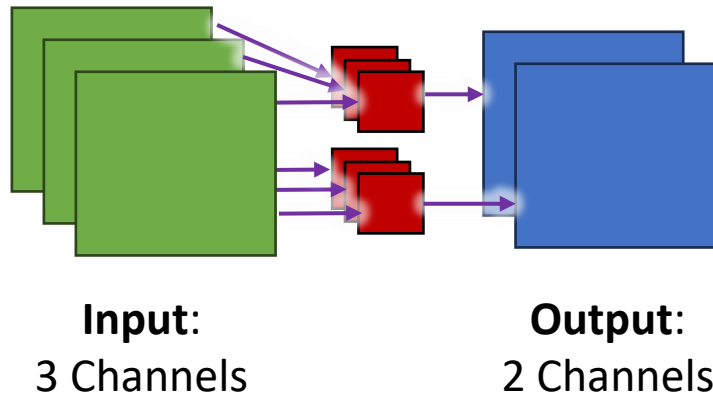
# Convolution: Common Practice



# Convolution: Common Practice



# Convolution Layer vs Feedforward Layer

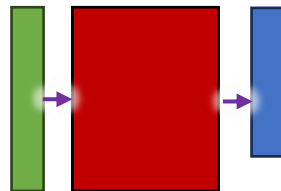


## Convolution Layer

$$\bullet A^{[l]} = W^{[l]} * A^{[l-1]}$$

Weights = **Kernels**  
(4D tensor)

Input and Output (3D Matrix)  
Concatenation of feature maps



## Feedforward (Linear) Layer

$$\bullet a^{[l]} = (W^{[l]})^T a^{[l-1]}$$

Weights  
(2D matrix)

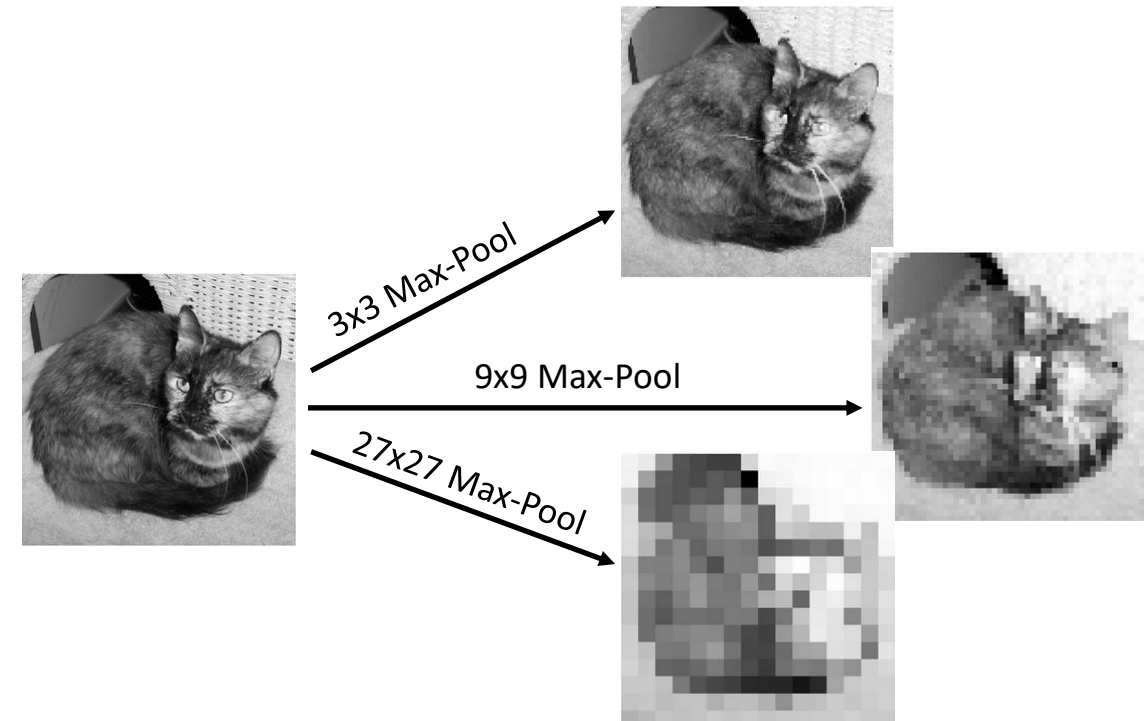
# Pooling Layer

- **Downsamples** Feature Maps
- Helps to train later kernels to detect **higher-level** features
- Reduces **dimensionality**
- Aggregation methods
  - Max-Pool (most common)
  - Average-Pool
  - Sum-Pool

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

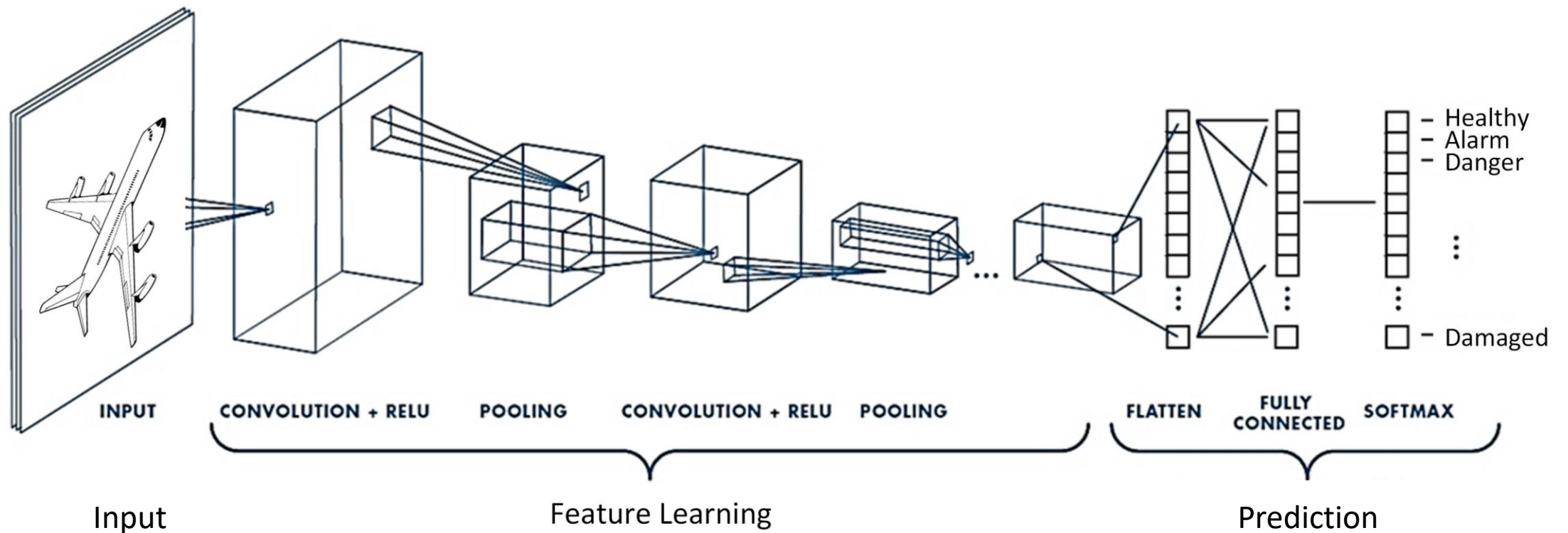
2x2 Max-Pool  
Stride = 2

20	30
112	37

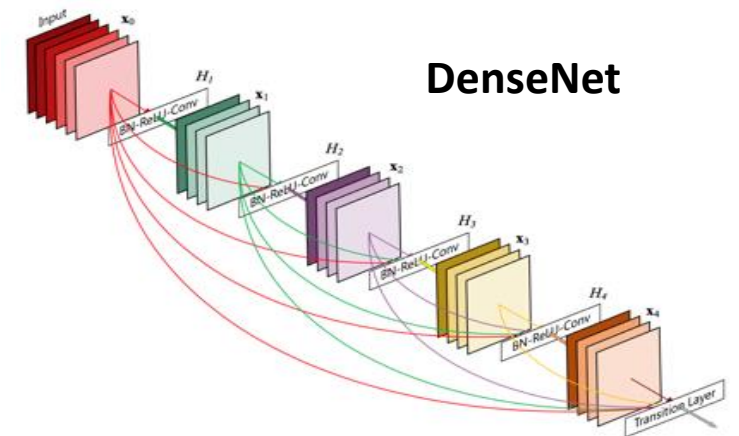
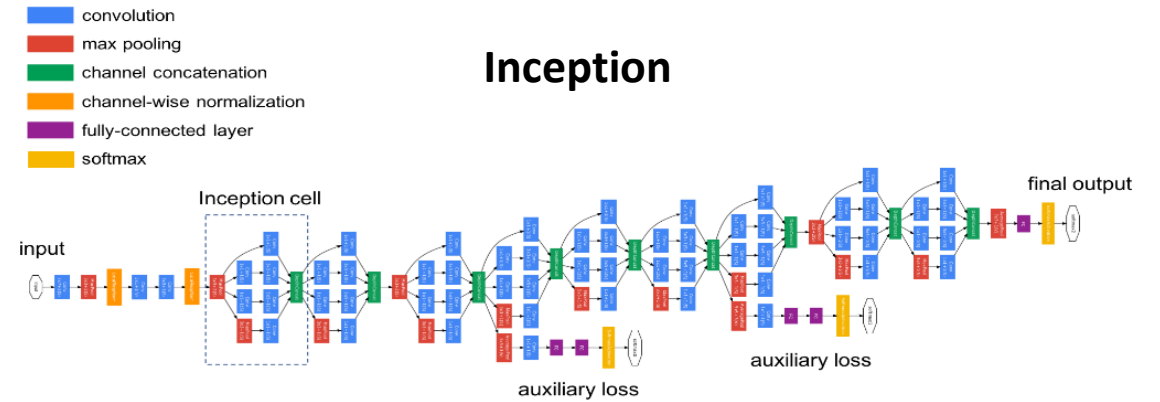
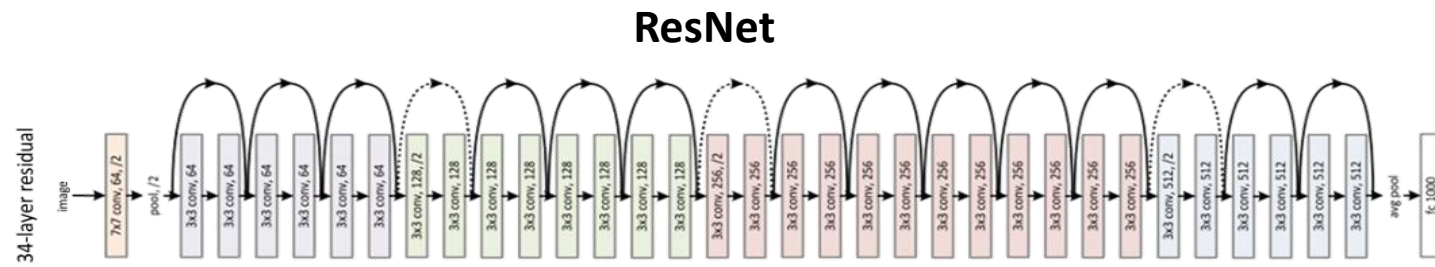
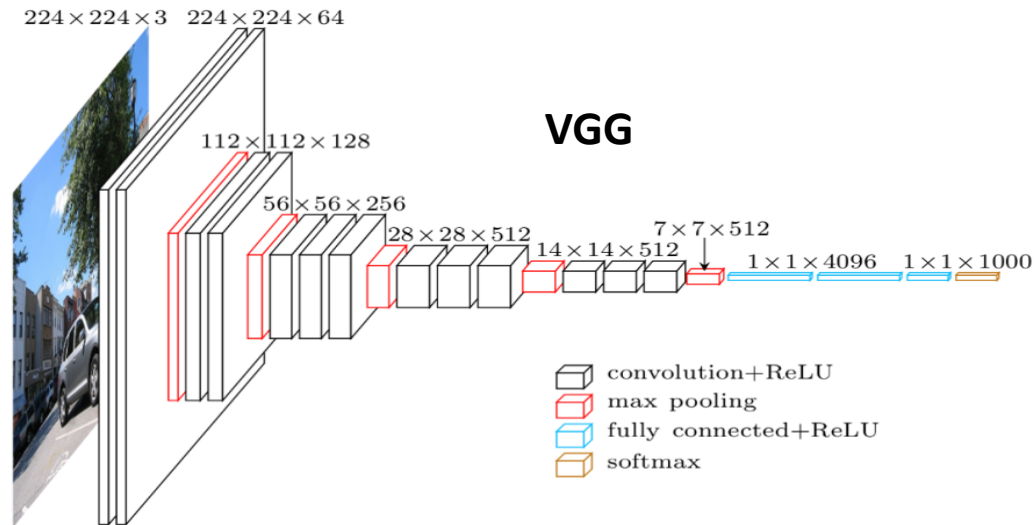


Typical

# Convolutional Neural Networks (CNN)



# Popular CNN Architectures



Further reading: <https://www.jeremyjordan.me/convnet-architectures/>

# Applications of CNN



Image Classification  
e.g., face emotions

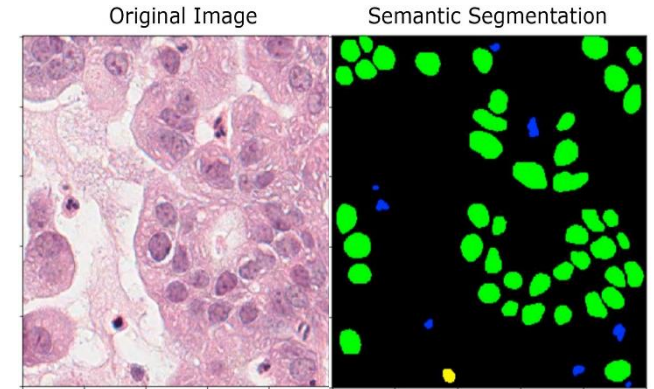
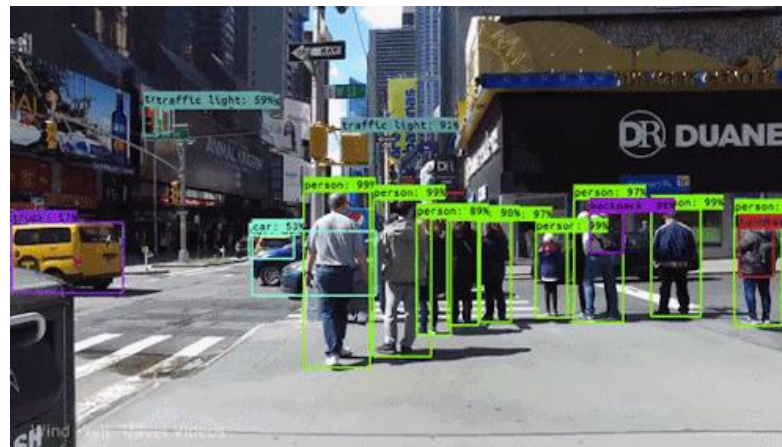


Image Segmentation  
e.g., cancer cell detection



Object Detection  
e.g., self-driving cars



# Outline

- Deep Neural Networks
- Convolution Neural Networks
  - Motivation: handling spatial structure
  - Convolution, Pooling Layer, and Common Architectures
  - Applications
- **Recurrent Neural Networks**
  - Motivation: handling sequential data
  - Recurrent Neural Networks and Variants
  - Applications
- Attention, Transformers, GPT, and ChatGPT (if time permits)
- Issues with Deep Learning

“Indomie is the best noodles ever”

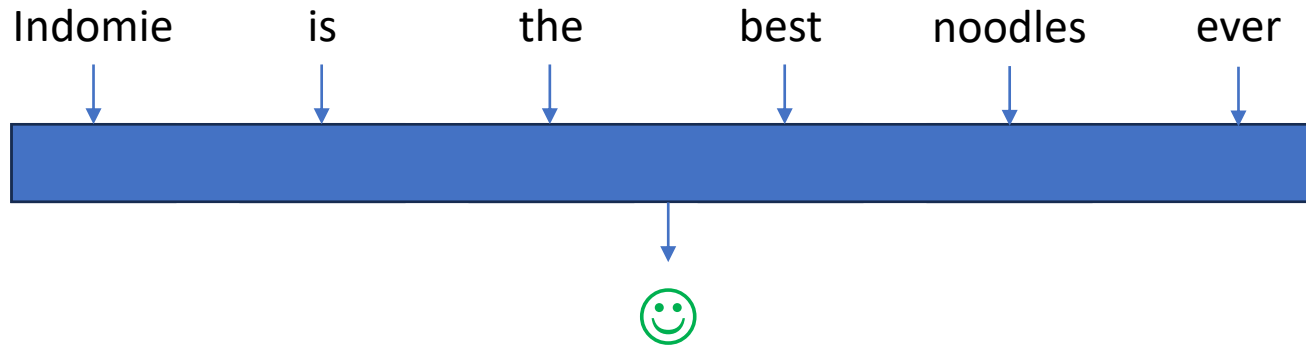
Credit: Indomie



# Sequential Data: 1<sup>st</sup> Attempt

## Sentiment Analysis

“Indomie is the best noodles ever”



Credit: Indomie



# Sequential Data: 1<sup>st</sup> Attempt

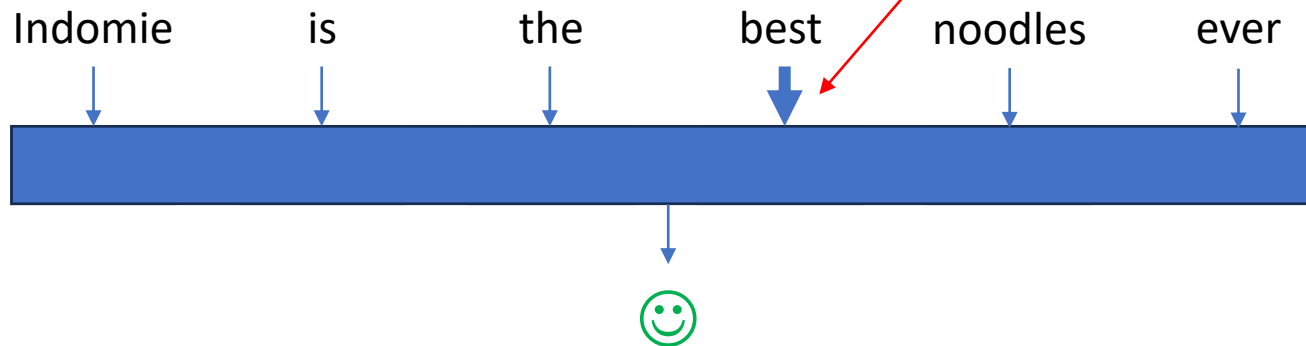
## Sentiment Analysis

“Indomie is the best noodles ever”

Credit: Indomie



Learn that this input/feature is important



# Sequential Data: 1<sup>st</sup> Attempt

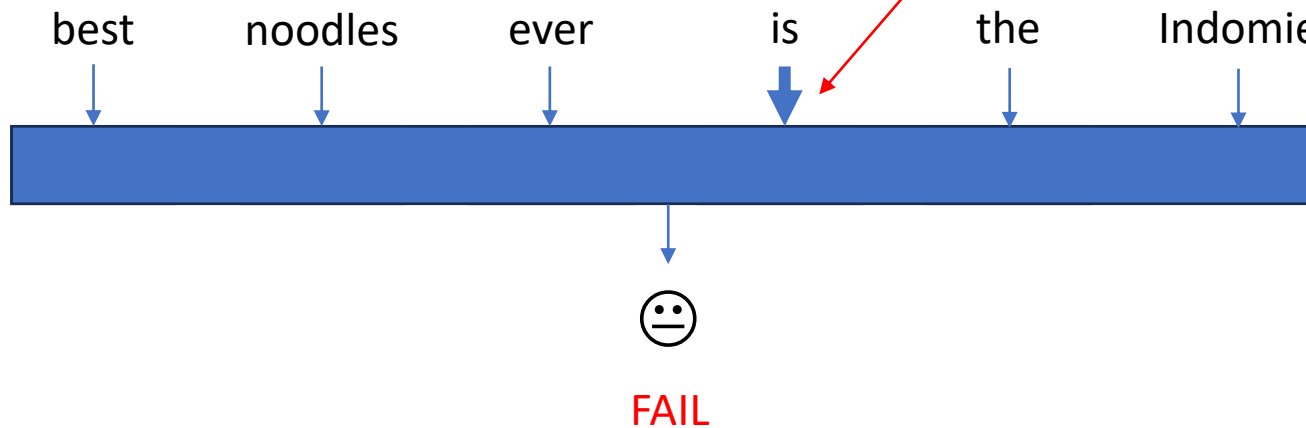
## Sentiment Analysis

“Indomie is the best noodles ever”

Credit: Indomie



Learn that this input/feature is important

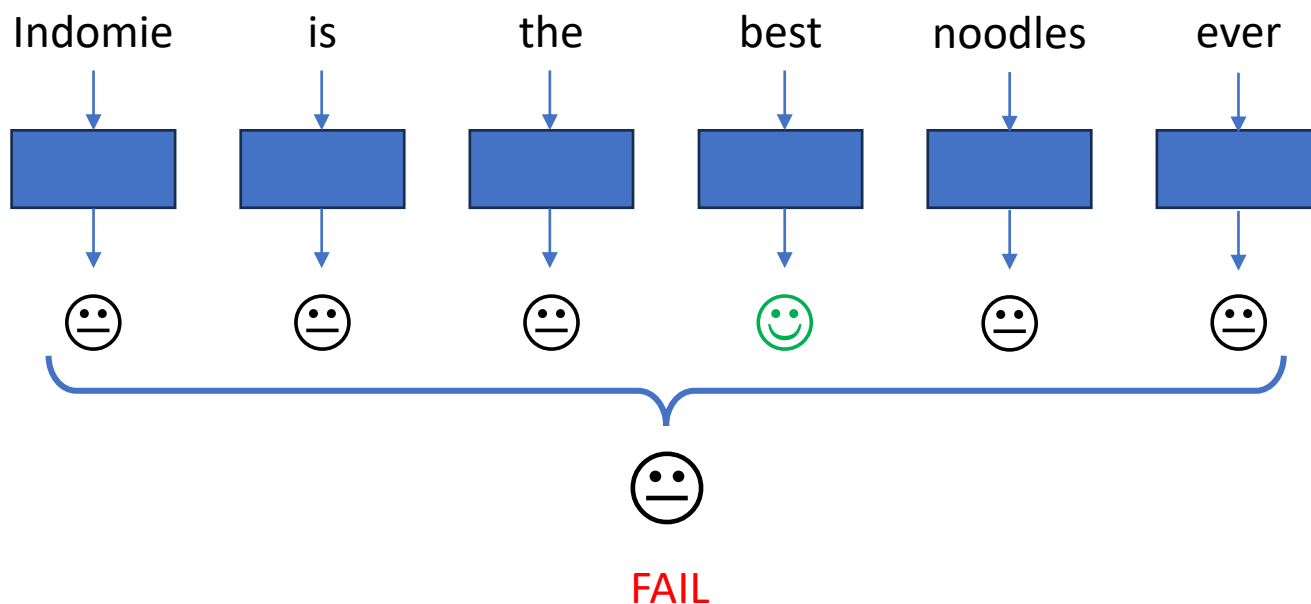


# Sequential Data: 2<sup>nd</sup> Attempt

## Sentiment Analysis

“Indomie is the best noodles ever”

Credit: Indomie

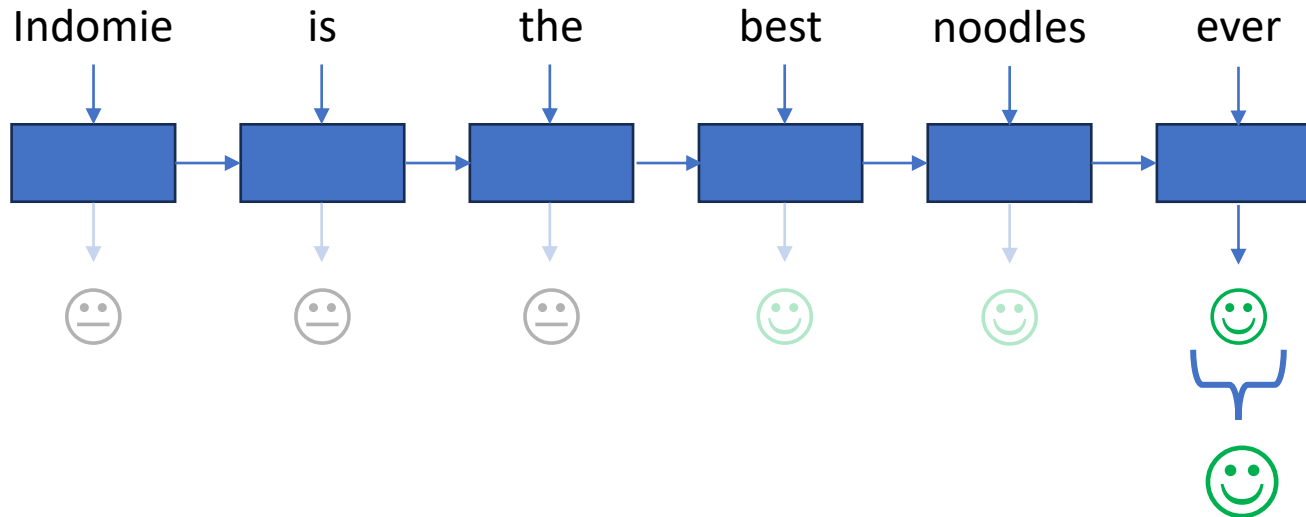


# Sequential Data: 3<sup>rd</sup> Attempt

## Sentiment Analysis

“Indomie is the best noodles ever”

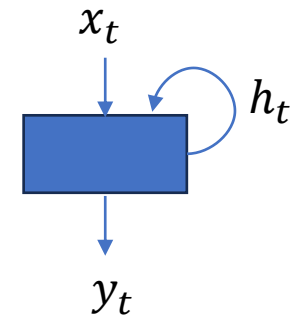
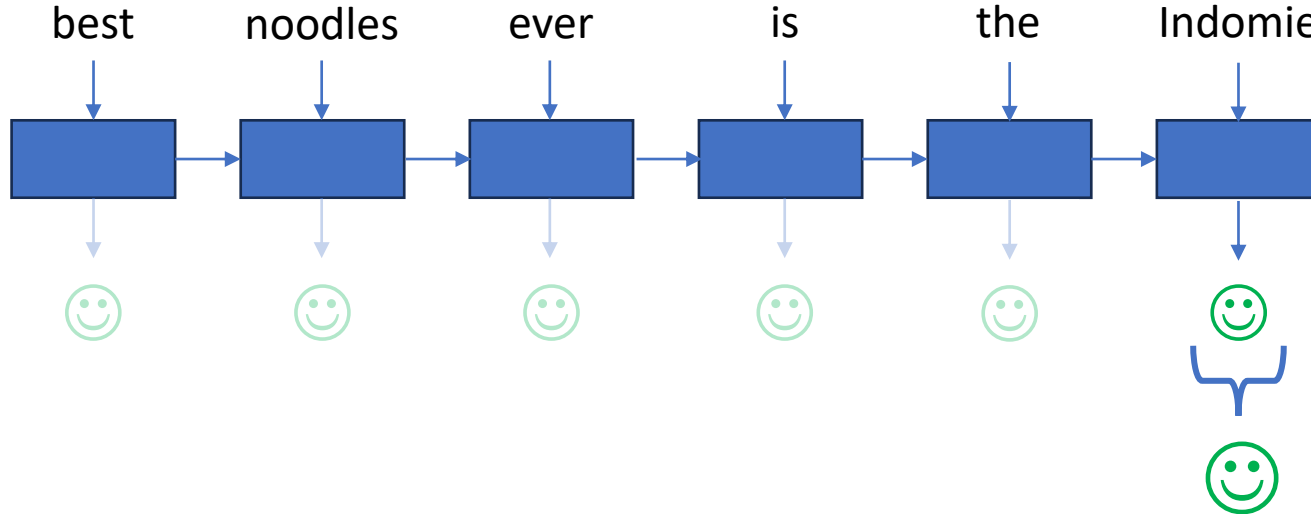
Credit: Indomie



# Sequential Data: 3<sup>rd</sup> Attempt

## Sentiment Analysis

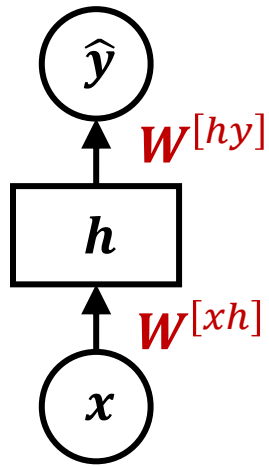
“Indomie is the best noodles ever”



Recurrent Neural Networks (RNN)



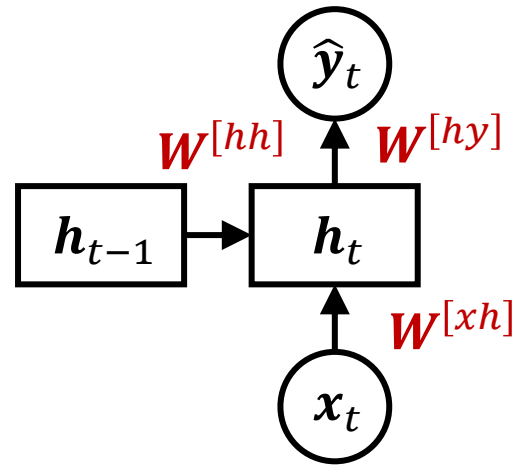
# Recurrent Neural Networks (RNN)



Feed-forward networks

$$\mathbf{y} = g^{[y]} \left( (\mathbf{W}^{[hy]})^\top \mathbf{h} \right)$$

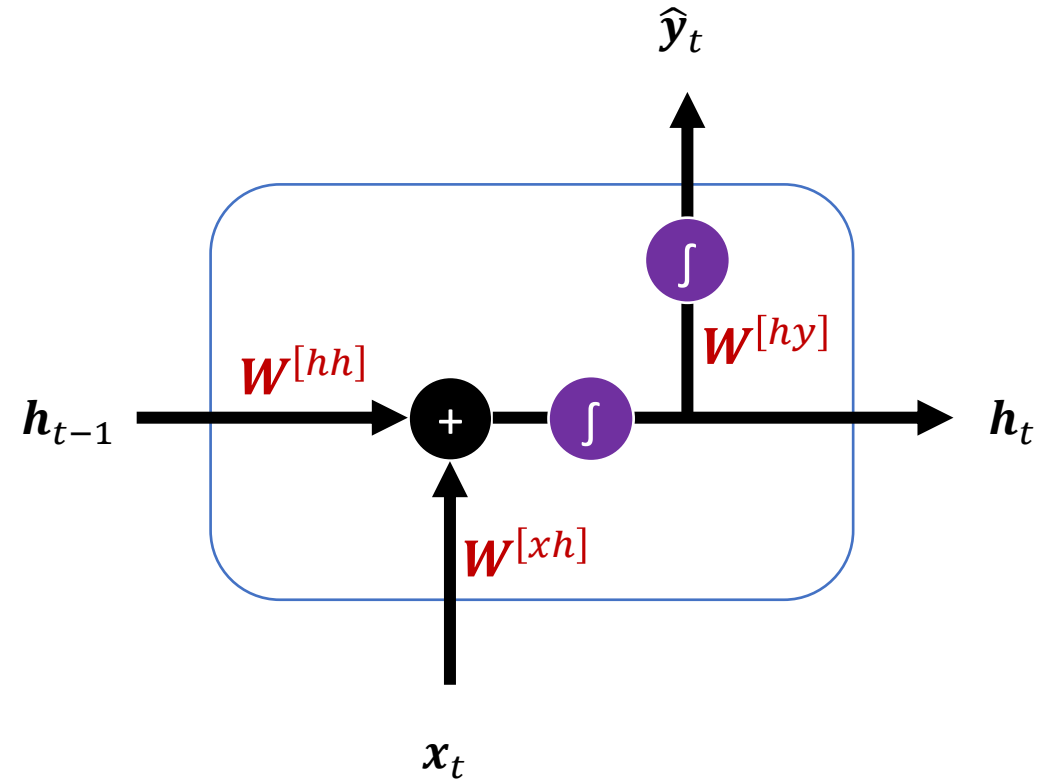
$$\mathbf{h} = g^{[h]} \left( (\mathbf{W}^{[xh]})^\top \mathbf{x} \right)$$



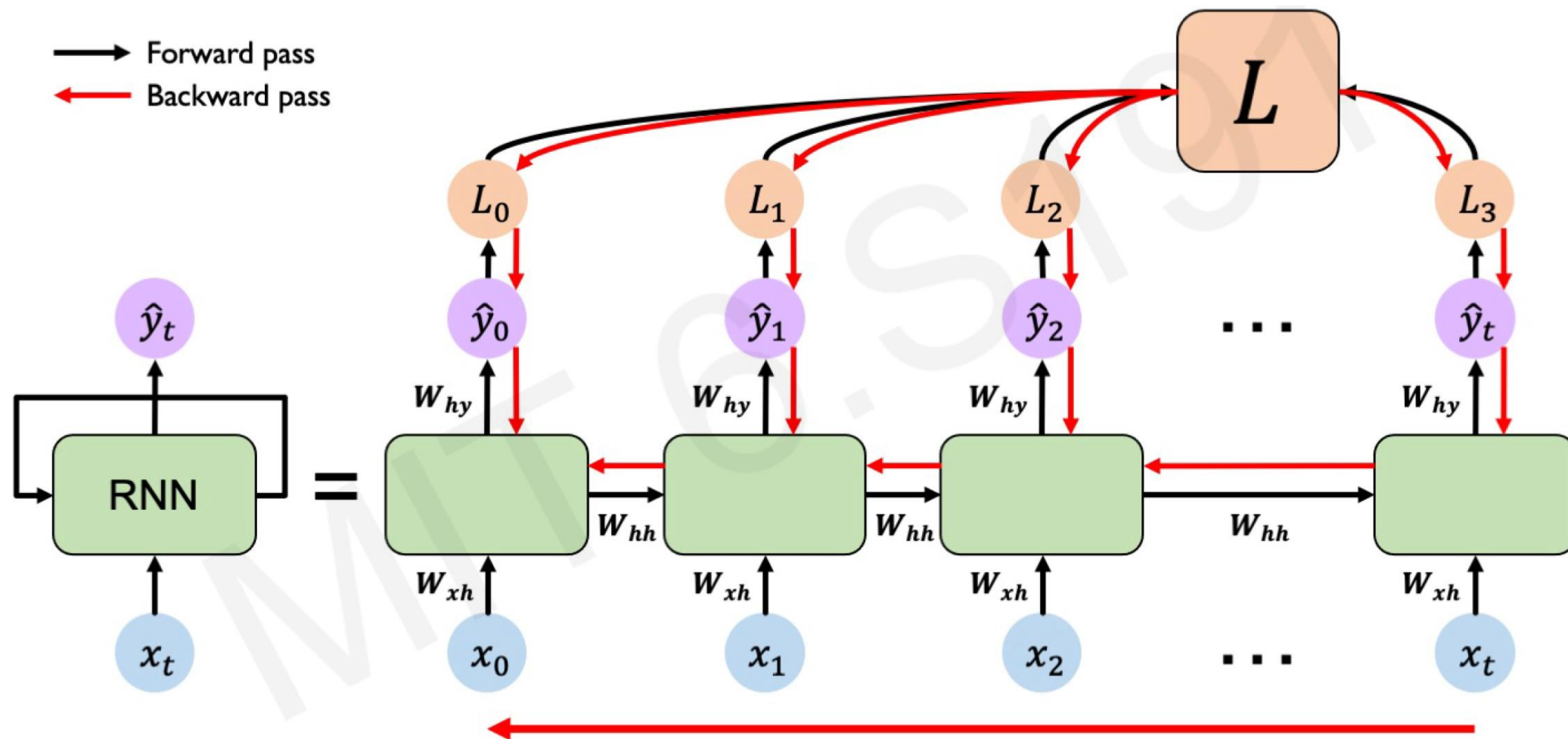
Recurrent Neural Networks

$$\mathbf{y}_t = g^{[y]} \left( (\mathbf{W}^{[hy]})^\top \mathbf{h}_t \right)$$

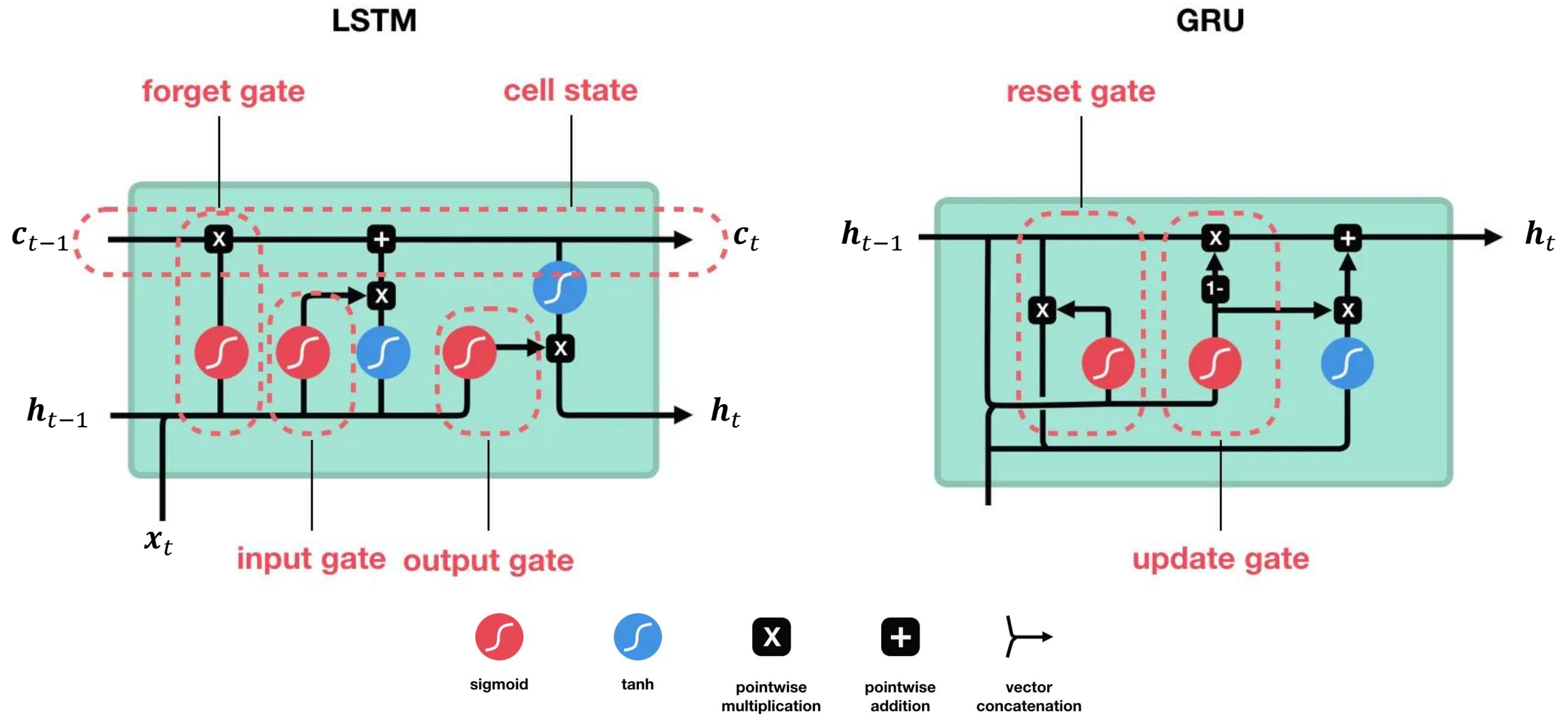
$$\mathbf{h}_t = g^{[h]} \left( (\mathbf{W}^{[xh]})^\top \mathbf{x}_t + (\mathbf{W}^{[hh]})^\top \mathbf{h}_{t-1} \right)$$



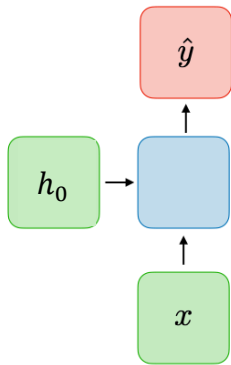
# Backpropagation Through Time (BPTT)



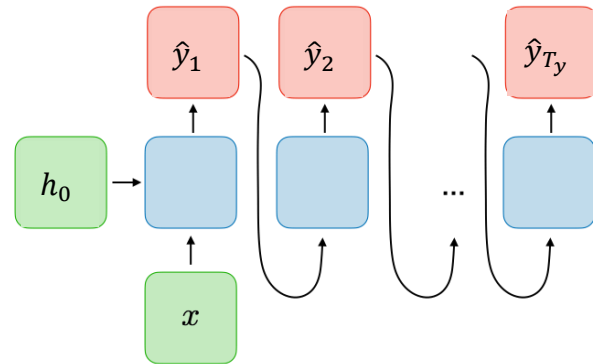
# Long Short-Term Memory (LSTM) & Gated Recurrent Unit (GRU)



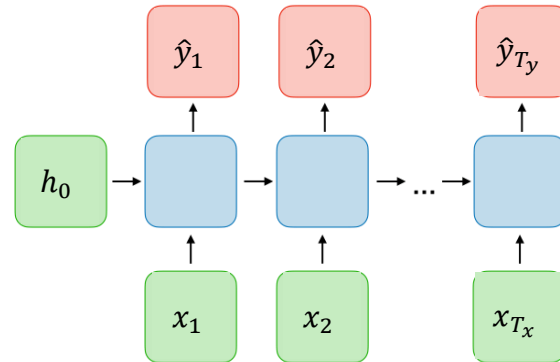
# Sequence Modelling



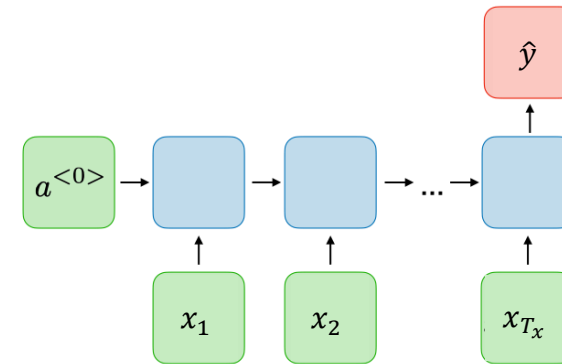
**One-to-one**  
 $T_x = T_y = 1$   
 Feedforward Network



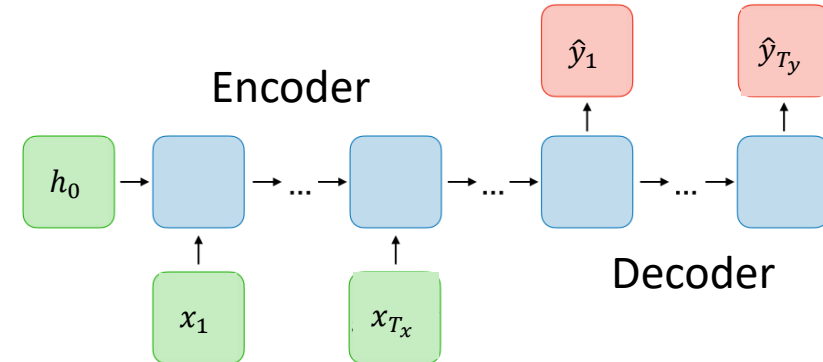
**One-to-many**  
 $T_x = 1, T_y > 1$   
 Image captioning, text generation



**Many-to-many**  
 $T_x > 1, T_y > 1$   
 Name entity recognition



**Many-to-one**  
 $T_x > 1, T_y = 1$   
 Sentiment classification



**Many-to-many**  
 $T_x \neq T_y$   
 Machine translation

# Applications of RNN

	$x$ - input		$y$ - output
Speech recognition		→	"Fuzzy Wuzzy was a bear. Fuzzy Wuzzy had no hair."
Music generation	$\emptyset$	→	
Sentiment classification	"Decent effort. The plot could have been better."	→	
DNA sequence analysis	ACTGTACCCATGTGACTGCCC	→	ACTGTACCCATGTGACTGCCC
Machine translation	"El que no arriesga, no gana."	→	"If you don't take risks, you cannot win."
Video activity recognition		→	Running
Name entity recognition	"Ygritte says Jon Snow knows nothing."	→	"Ygritte says Jon Snow knows nothing."

# Outline

- Deep Neural Networks
- Convolution Neural Networks
  - Motivation: handling spatial structure
  - Convolution, Pooling Layer, and Common Architectures
  - Applications
- Recurrent Neural Networks
  - Motivation: handling sequential data
  - Recurrent Neural Networks and Variants
  - Applications
- **Attention, Transformers, GPT, and ChatGPT** (if time permits)
- Issues with Deep Learning

# Long term dependencies problem

## Machine Translation (English → Japanese)

“Indomie is the best noodles ever created by humankind. ... <100 words later> ... This concludes my argument why I like the noodles.”



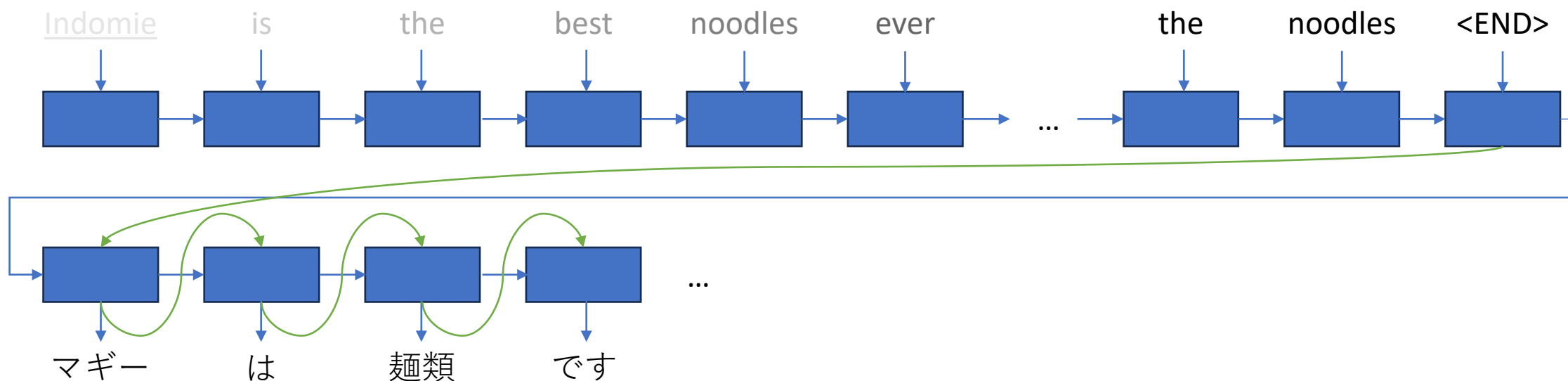


# Long term dependencies problem

## Machine Translation (English → Japanese)

“Indomie is the best noodles ever created by humankind. ... **<100 words later>** ... This concludes my argument why I like the noodles.”

Forgot about what it has read before!



“Maggi is noodles...”



Credit: Indomie

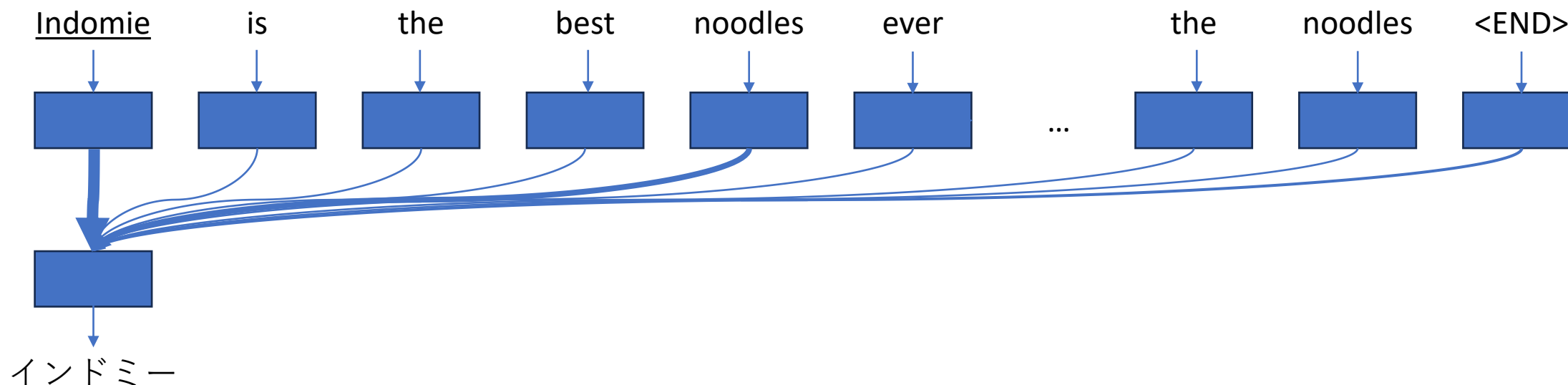


# Long term dependencies problem

## Machine Translation (English → Japanese)

“Indomie is the best noodles ever created by humankind. ... **<100 words later>** ... This concludes my argument why I like the noodles.”

Attends to (focus on) what's necessary!



Credit: Indomie

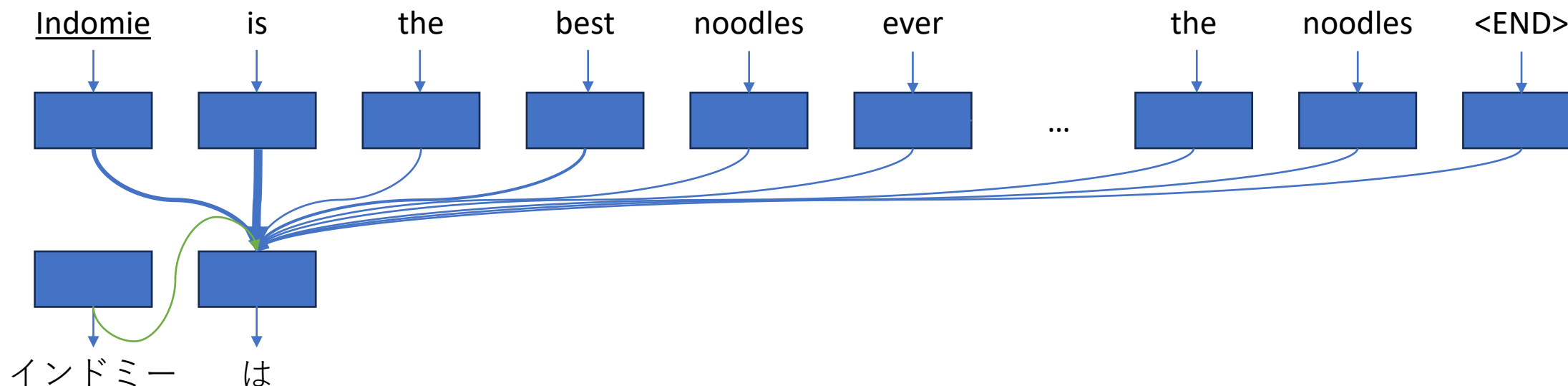


# Long term dependencies problem

## Machine Translation (English → Japanese)

“Indomie is the best noodles ever created by humankind. ... **<100 words later>** ... This concludes my argument why I like the noodles.”

Attends to (focus on) what's necessary!



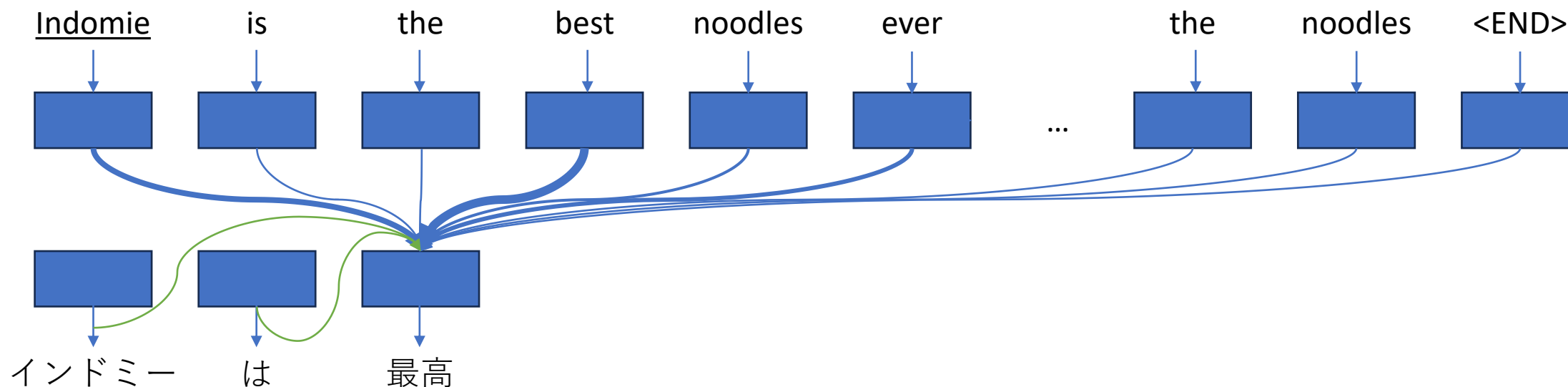


# Long term dependencies problem

## Machine Translation (English → Japanese)

“Indomie is the best noodles ever created by humankind. ... **<100 words later>** ... This concludes my argument why I like the noodles.”

Attends to (focus on) what's necessary!



Credit: Indomie

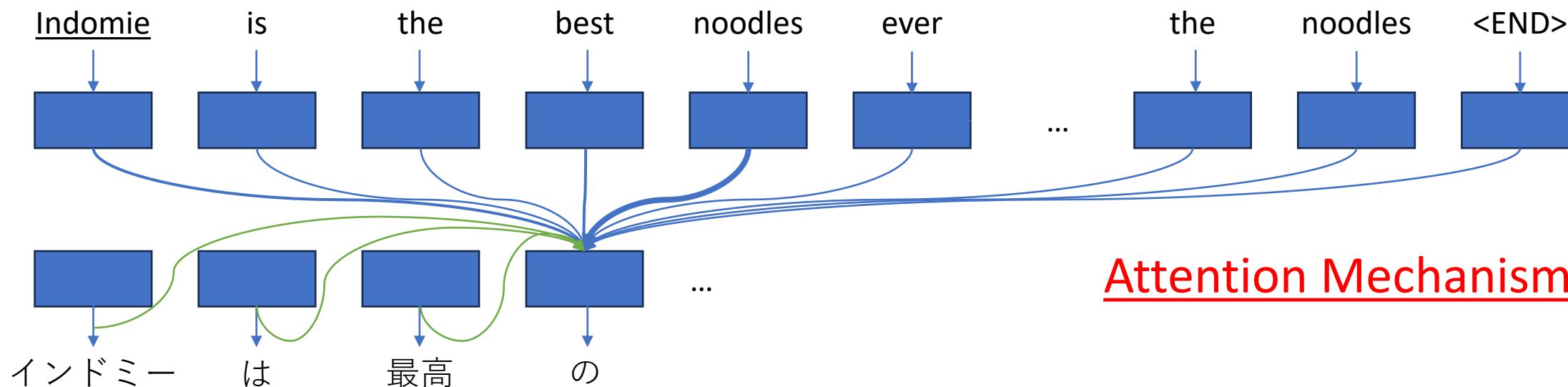


# Long term dependencies problem

## Machine Translation (English → Japanese)

“Indomie is the best noodles ever created by humankind. ... **<100 words later>** ... This concludes my argument why I like the noodles.”

Attends to (focus on) what's necessary!

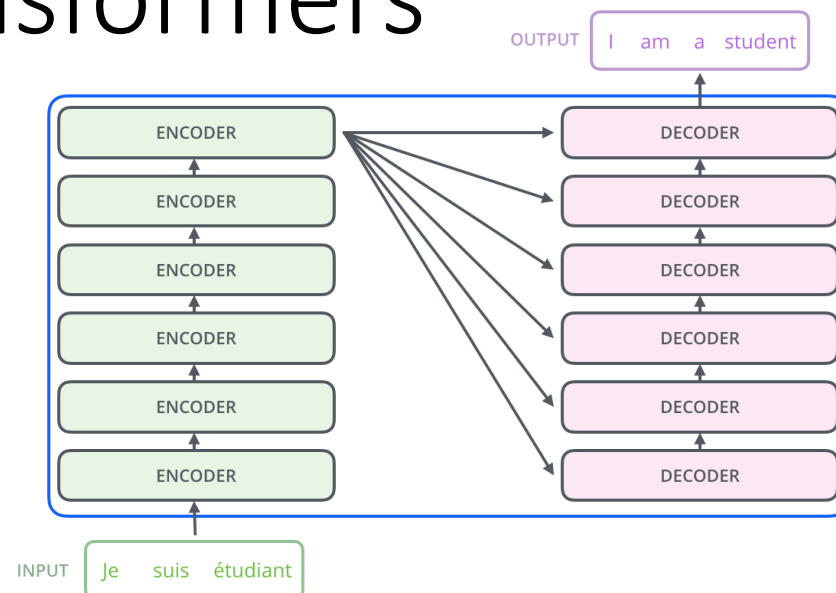


Attention Mechanism

“Indomie is the best ...”

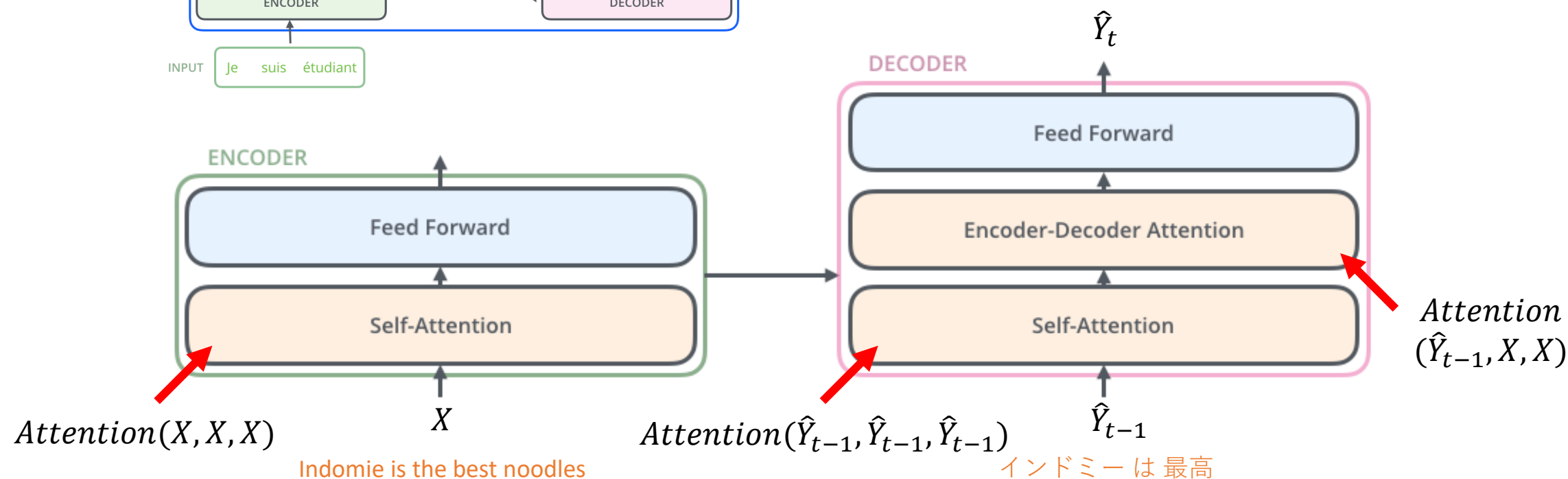


# Transformers

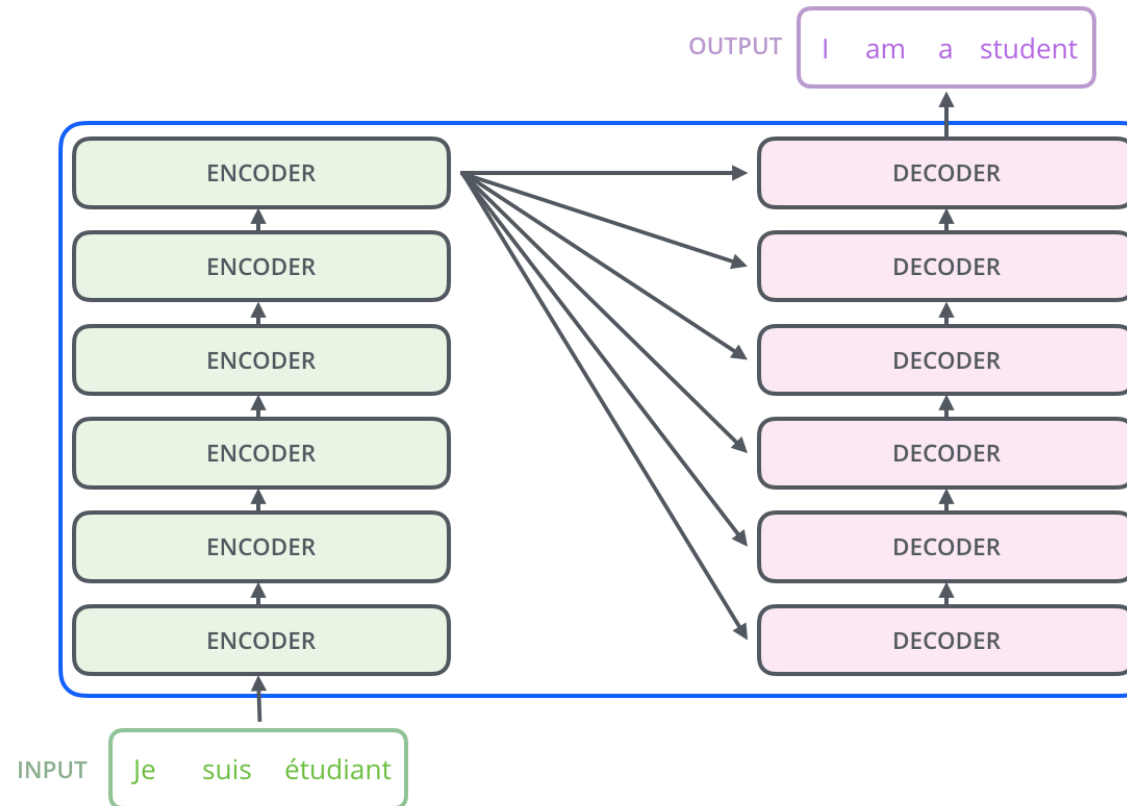


$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

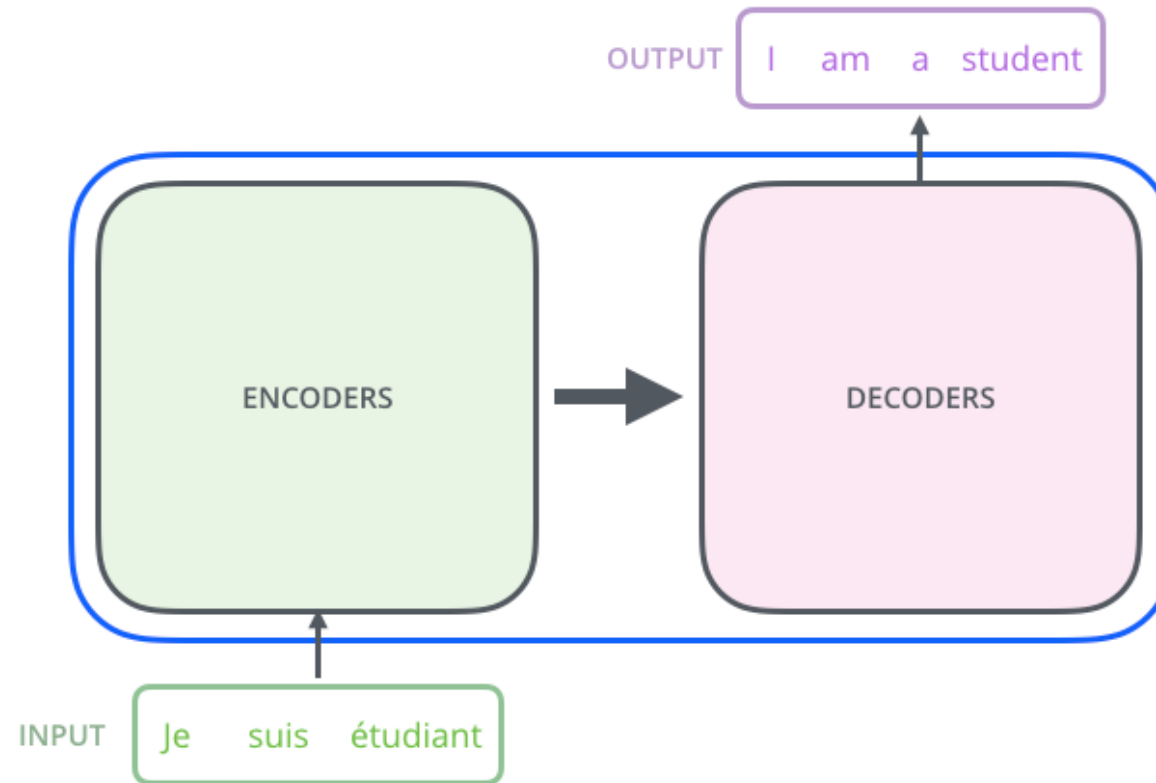
インドミーは最高の麺



# Transformers



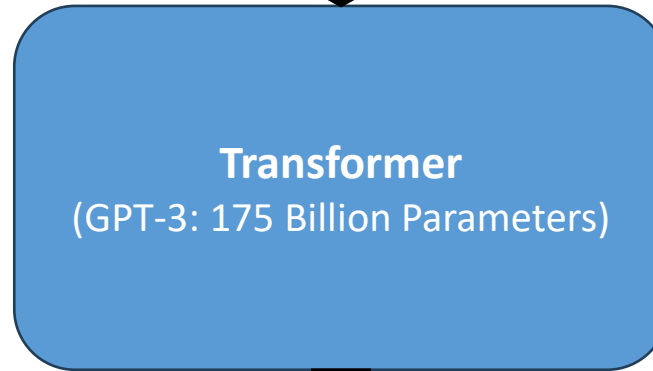
# Transformers





# Generative Pretrained Transformers (GPT)

Indomie is the best noodles ever created by \_



Indomie is the best noodles ever created by **humankind**

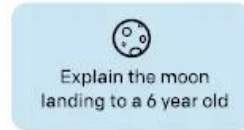
Trained to **predict the next word** on ~300 billion tokens (~words)

# ChatGPT

## Step 1

**Collect demonstration data, and train a supervised policy.**

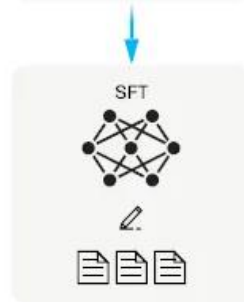
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



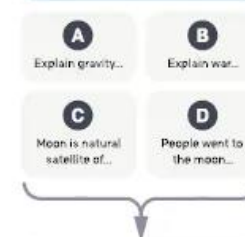
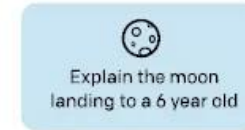
This data is used to fine-tune GPT-3 with supervised learning.



## Step 2

**Collect comparison data, and train a reward model.**

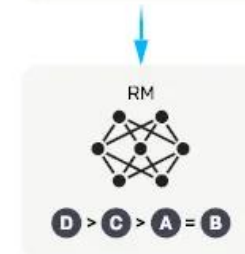
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



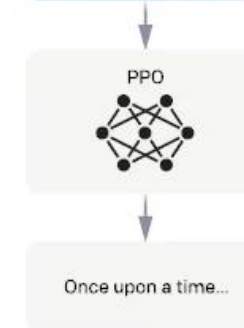
## Step 3

**Optimize a policy against the reward model using reinforcement learning.**

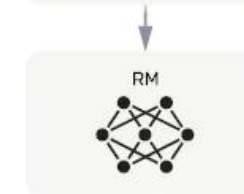
A new prompt is sampled from the dataset.



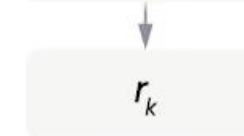
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



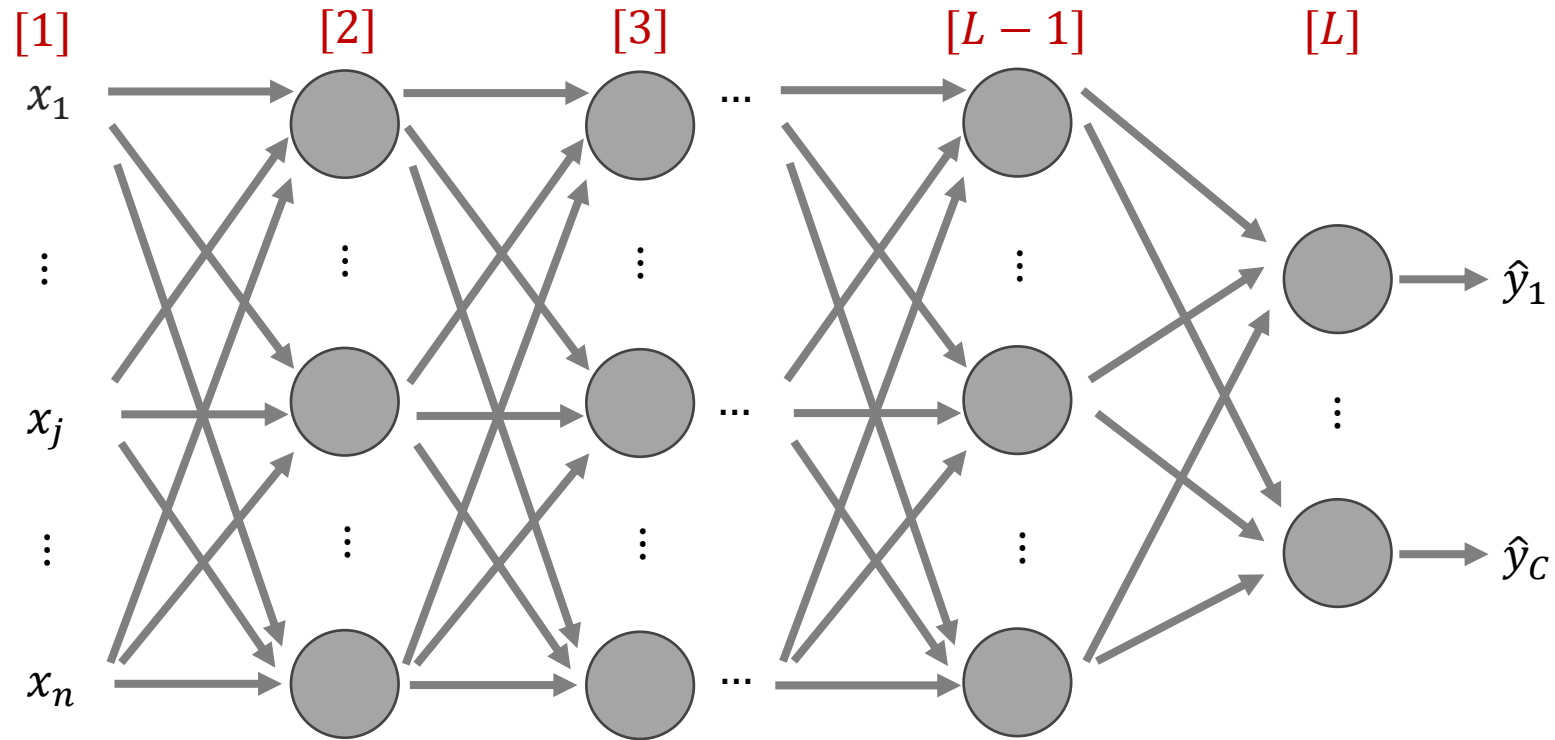
# Outline

- Deep Neural Networks
- Convolution Neural Networks
  - Motivation: handling spatial structure
  - Convolution, Pooling Layer, and Common Architectures
  - Applications
- Recurrent Neural Networks
  - Motivation: handling sequential data
  - Recurrent Neural Networks and Variants
  - Applications
- Attention, Transformers, GPT, and ChatGPT (if time permits)
- **Issues with Deep Learning**

# Issues with Deep Learning

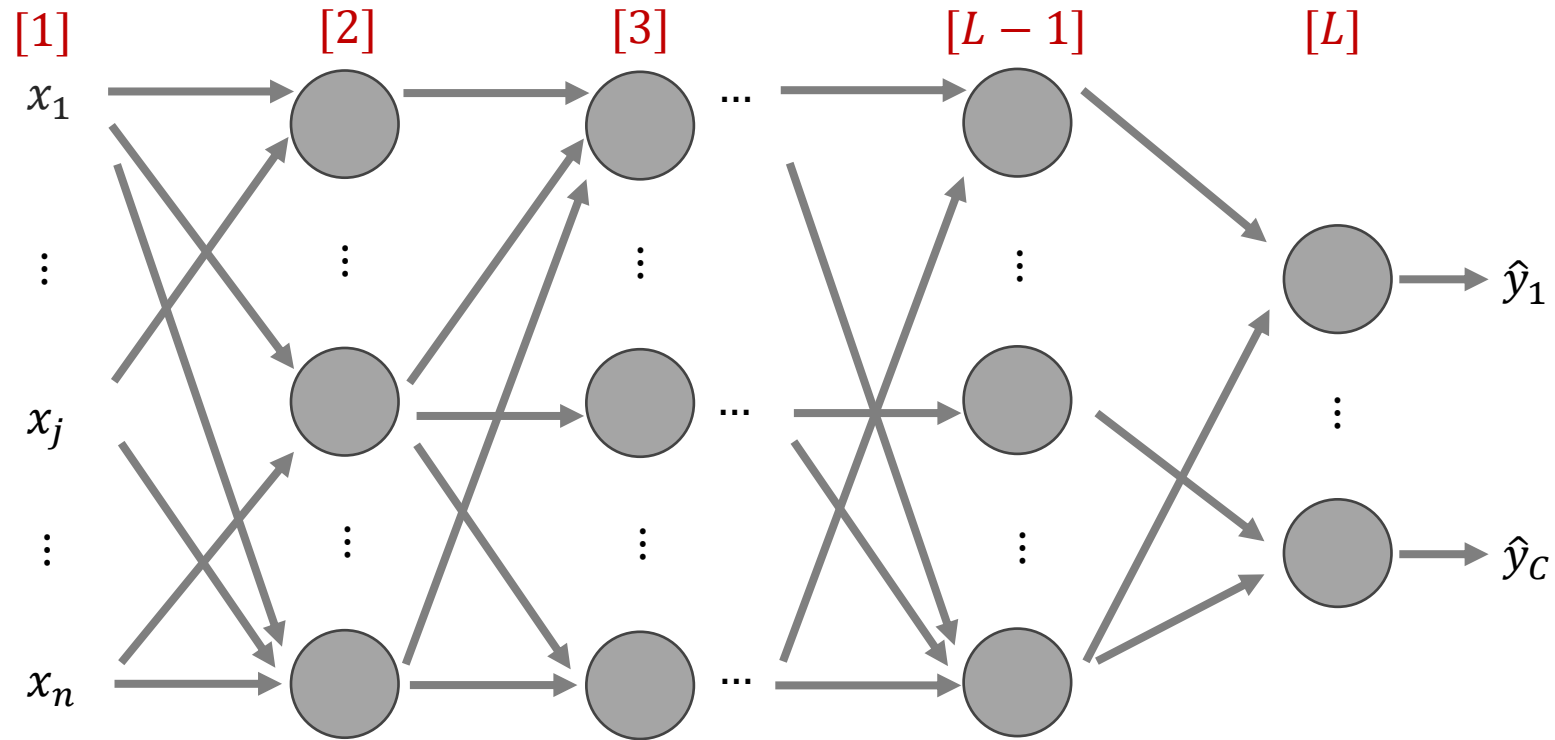
- Overfitting → Regularization
- Gradient Vanishing/Exploding

# Regularization: Dropout



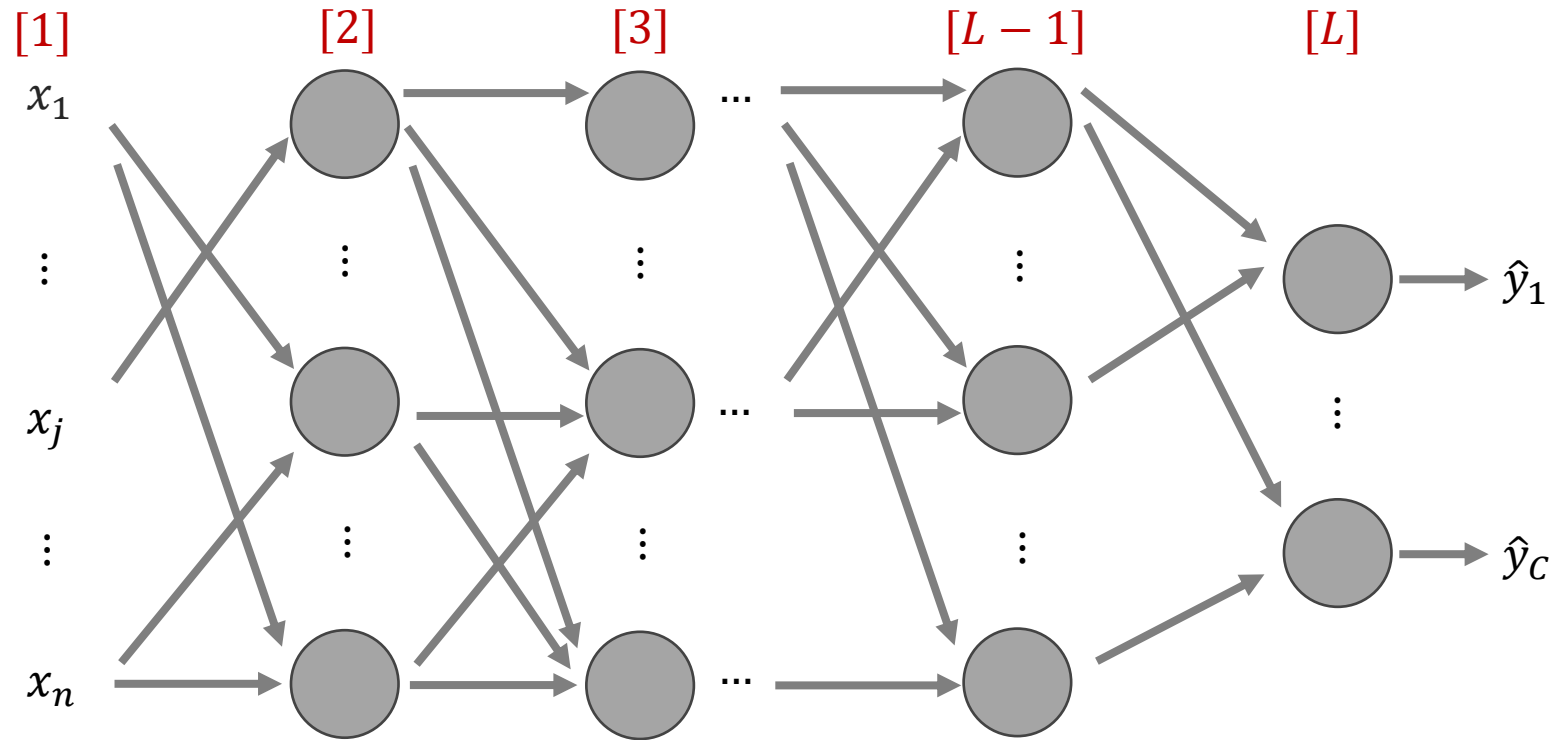
During training, **randomly** set some activations to 0

# Regularization: Dropout



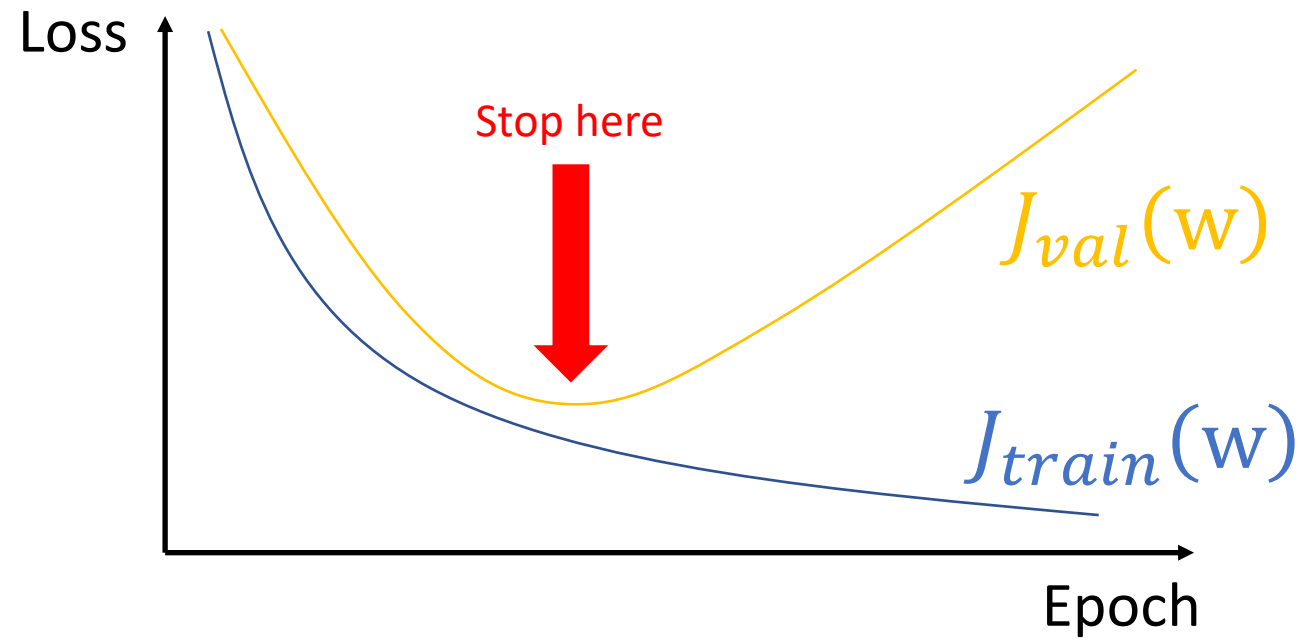
During training, **randomly** set some activations to 0

# Regularization: Dropout



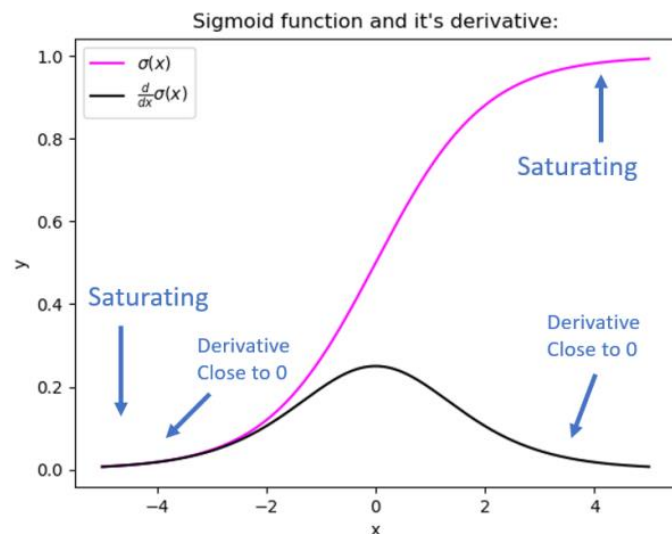
During training, **randomly** set some activations to 0

# Regularization: Early Stopping





# Vanishing/Exploding Gradient



- **Vanishing gradient:** **small** gradients got multiplied again and again until it reaches almost zero.
- **Exploding gradient:** **large** gradients got multiplied again and again until it overflows.
- **Mitigation:**
  - Proper Weight Initialization
  - Using Non-saturating Activation Functions, e.g. ReLU
  - Batch Normalization ~ **feature scaling at every layer**
  - Gradient Clipping ~ **clip gradient within range [min, max]**

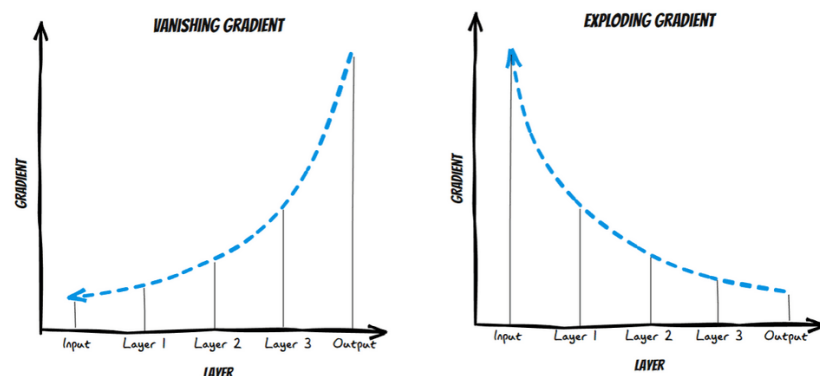


Image credit:

<https://neptune.ai/blog/vanishing-and-exploding-gradients-debugging-monitoring-fixing>

<https://medium.com/geekculture/how-to-deal-with-vanishing-and-exploding-gradients-in-neural-networks-24eb00c80e84>

# Summary

- Deep Neural Networks – neural networks with >3 layers
- Convolution Neural Networks
  - Motivation: handling **spatial structure**, translation **invariant**
  - Convolution (multiply-sum), Pooling (downsampling) Layer, and Common Architectures
  - Applications: image recognition, image segmentation, object detection
- Recurrent Neural Networks
  - Motivation: handling **sequential** data
  - Recurrent Neural Networks and Variants:
    - **neural networks with loop**  $y_t, h_t = RNN(x_t, h_{t-1})$
  - Applications: machine translation, summarization, etc
- Attention, Transformers, GPT, and ChatGPT
  - Attention: **focus on things that matters**. **Massive** neural networks; train with **billions of data**.
- Issues with Deep Learning: overfitting, gradient vanishing/exploding

# Coming Up Next Week

- **Unsupervised Learning**
- **Clustering**
  - K-means clustering
  - Hierarchical clustering
- **Dimensionality Reduction**
  - Principal component analysis (PCA) – Math!

# To Do

- **Lecture Training 10**
  - +100 Free EXP
  - +50 Early bird bonus
- **Problem Set 6 is due tomorrow!**