

# Deep Learning (RNN+DNN)

# 10

B

**CS 3244**  
**Machine Learning**



**Computing**

Convolutional Layer:  
Feature Kernels & Feature Maps

$$k^{[l-1]} = c^{[l]}$$

# filters from previous layer  $k^{[l-1]}$  is equal to  
#channels into current layer  $c^{[l]}$

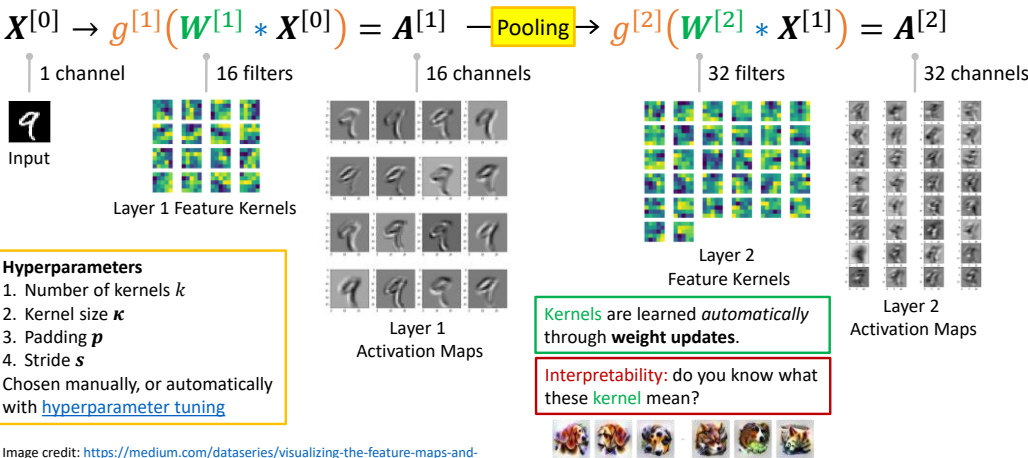


Image credit: <https://medium.com/dataseries/visualizing-the-feature-maps-and-filters-by-convolutional-neural-networks-e1462340518e>

Pooling Layer

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

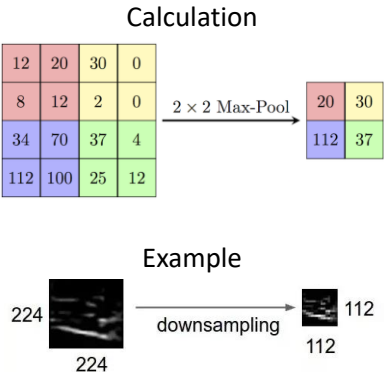


Image credit: <https://computersciencewiki.org/index.php/Max-pooling> / Pooling

Fully-Connected Neuron vs. Convolutional Activation Map

Weights  $\rightarrow$  Kernels

Sum across channels

2D Convolution

Activation function can be the same

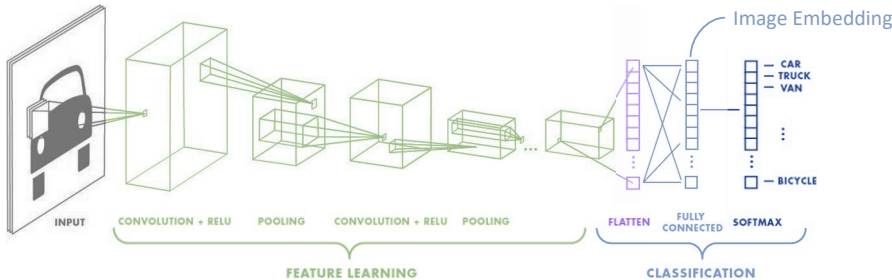
Activation is a 2D Matrix for each channel

Layer stores activation of **each neuron**  $p$  separately

Layer stores activation map of **each kernel**  $k$  separately

$$a_p^{[l]} = g^{[l]} \left( \sum_r^{n^{[l-1]}} w_{rp}^{[l]} a_r^{[l-1]} \right) \text{ vs. } A_k^{[l]} = g^{[l]} \left( \sum_c^{c^{[l-1]}} w_{ck}^{[l]} * A_c^{[l-1]} \right)$$

Convolutional Neural Network



Key concepts

- ❶ **Learn Spatial Feature**
  - Series of multiple convolution + pooling layers
  - Progressively learn more diverse and higher-level features
- ❷ **Flattening**
  - Convert to fixed-length 1D vector
- ❸ **Learn Nonlinear Features**
  - With fully connected layers (regular neurons)
  - Learns nonlinear relations with multiple layers
- ❹ **Classification**
  - Softmax  $\equiv$  Multiclass Logistic Regression
  - Feature input = image embedding vector (typically large vector)

Image credit: <https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53>

# Week 10A&B: Learning Outcomes

1. Understand how deep learning enables better model performance than shallow machine learning
2. Explain how CNNs and RNNs are different from feedforward neural networks
3. Appropriately choose and justify when to use each architecture
4. Explain how to mitigate training issues in deep learning

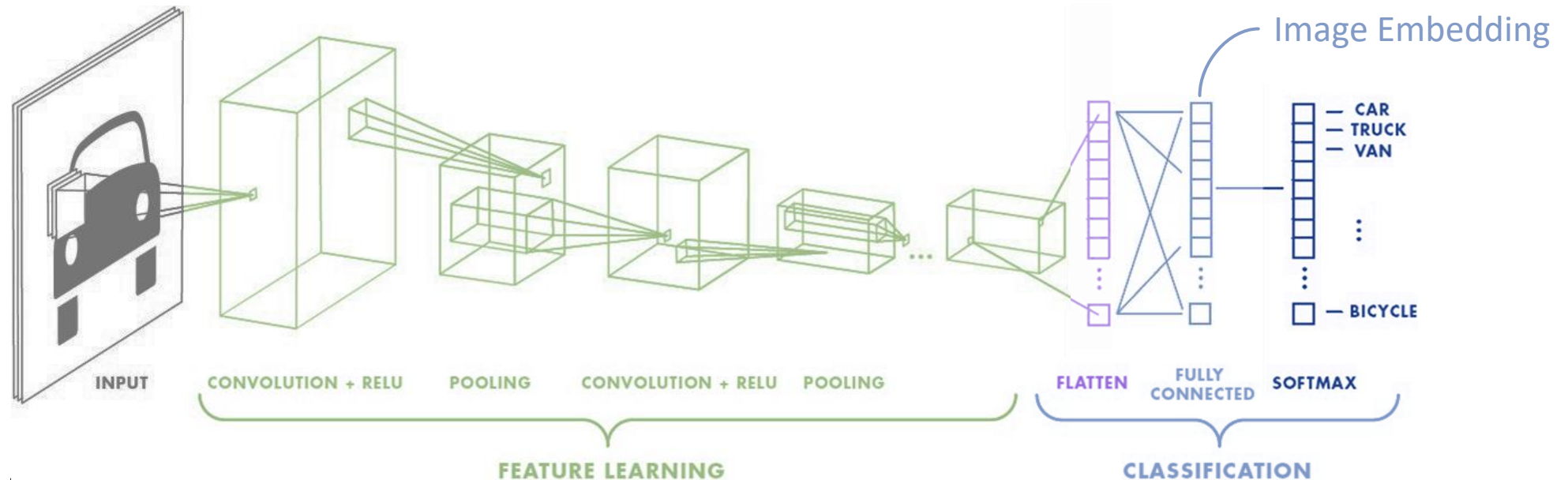
# Week 10B: Lecture Outline

1. Deep learning motivation
2. Popular Architectures
  1. Convolutional Neural Networks
  2. Recurrent Neural Networks
3. Deep learning training issues



# Convolutional Neural Networks (CNN)

# Convolutional Neural Network



## Key concepts

### ① Learn Spatial Feature

- Series of multiple convolution + pooling layers
- Progressively learn more diverse and higher-level features
- Analogy: human visual cortex

### ② Flattening

- Convert to fixed-length 1D vector

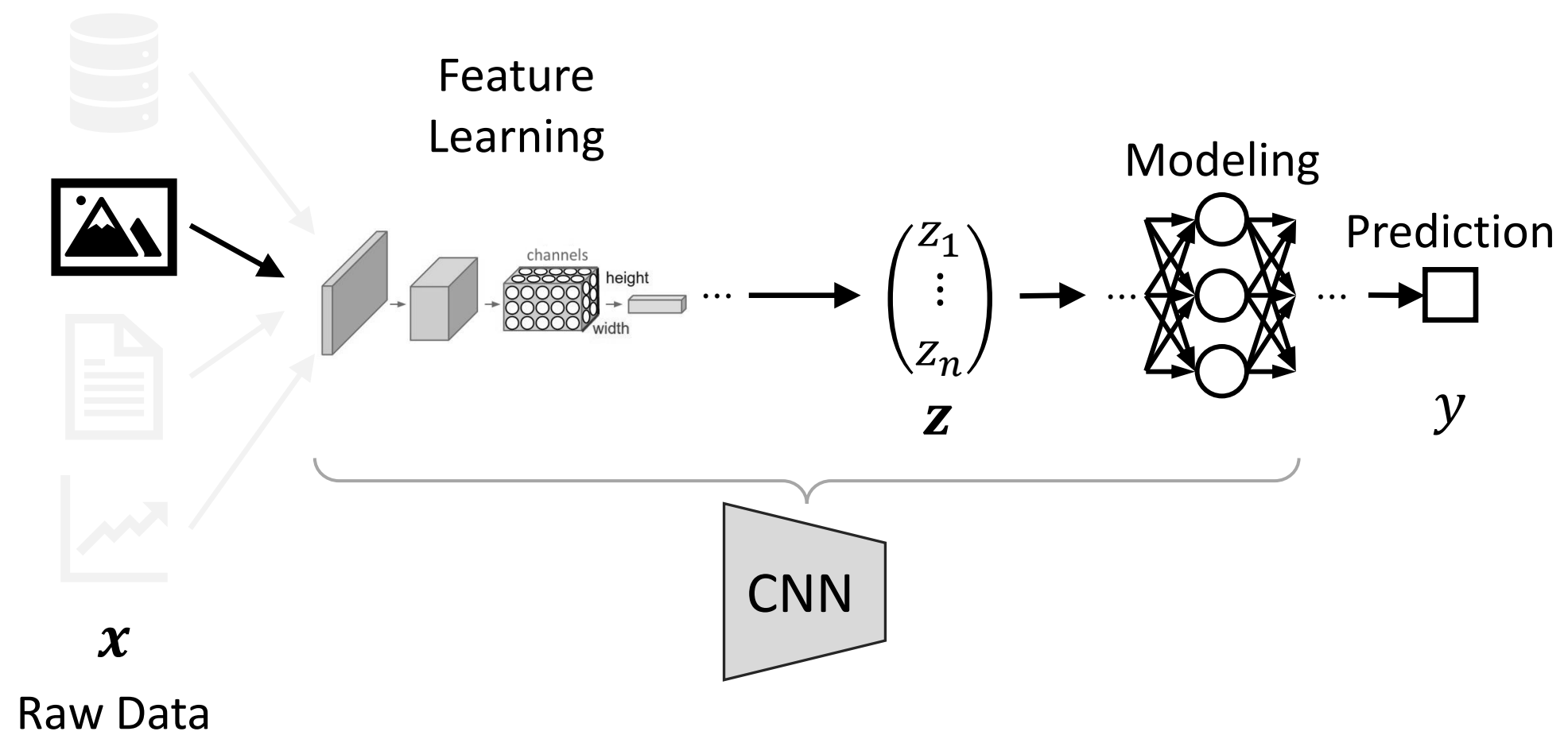
### ③ Learn Nonlinear Features

- With fully connected layers (regular neurons)
- Learns nonlinear relations with multiple layers
- Analogy: semantic reasoning

### ④ Classification

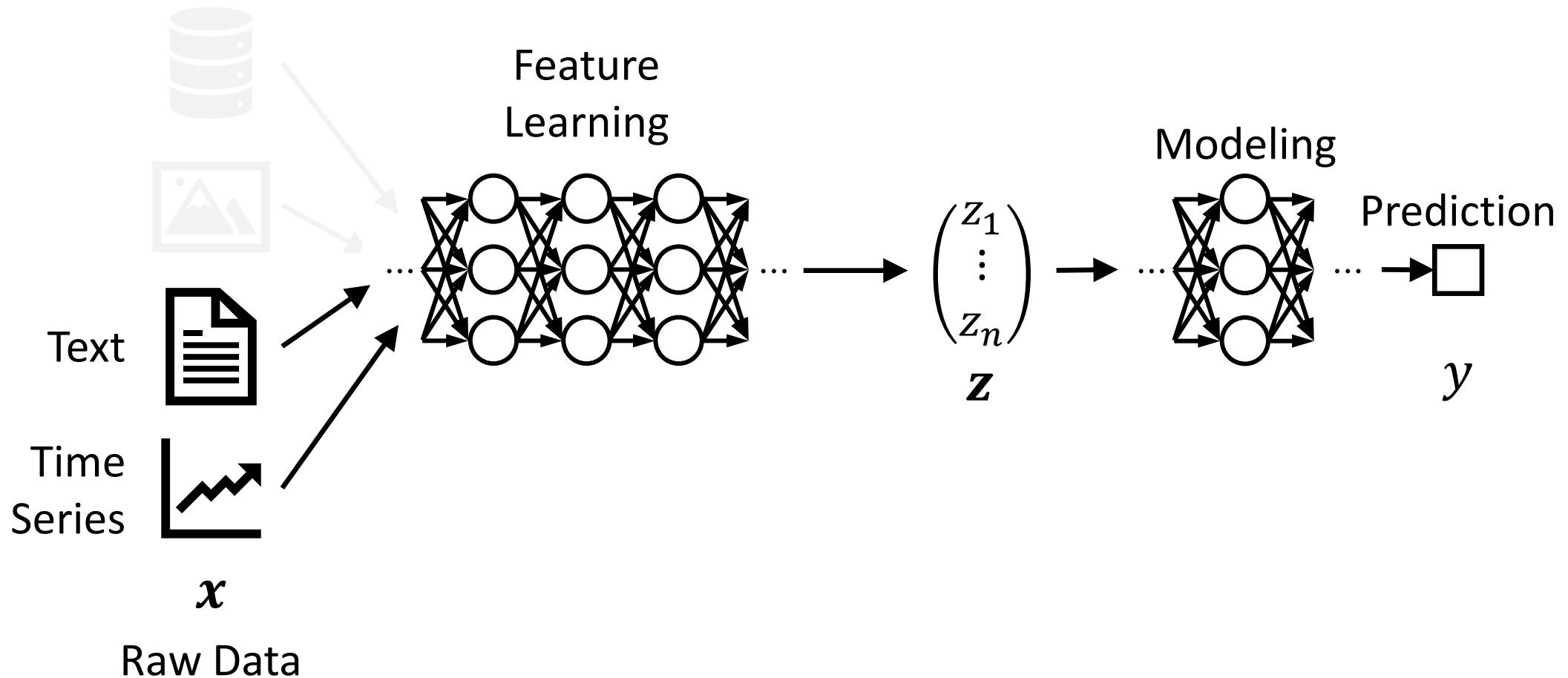
- Softmax := Multiclass Logistic Regression
- Feature input = image embedding vector (typically large vector)
- Analogy: decision making

# From Manual Feature Engineering To Automatic Feature Learning

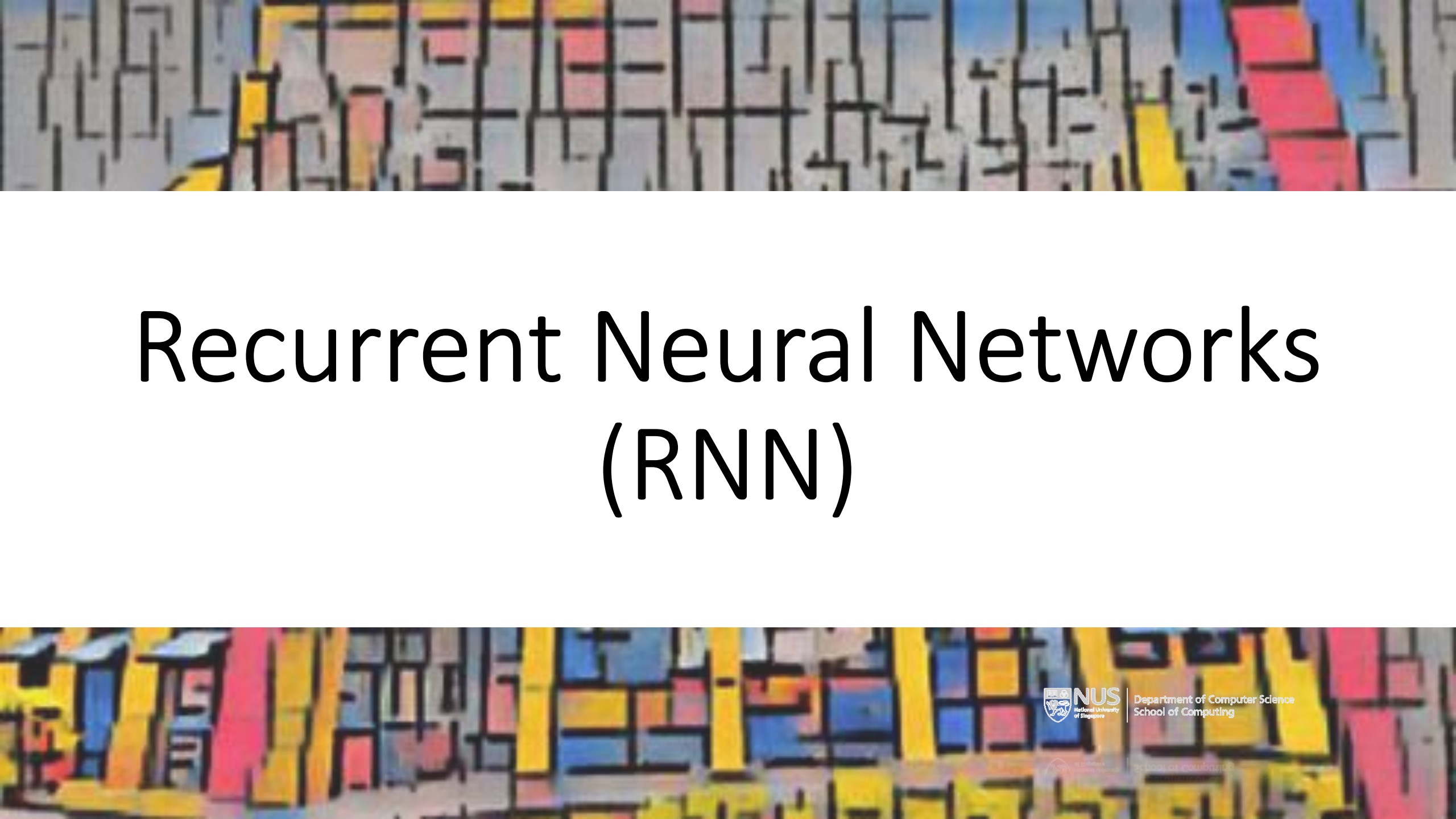




# From Manual Feature Engineering To Automatic Feature Learning







# Recurrent Neural Networks (RNN)

# Applications of RNN

Speech recognition



"The quick brown fox jumped  
over the lazy dog."

Music generation

∅



Sentiment classification

"There is nothing to like  
in this movie."



DNA sequence analysis

AGCCCCTGTGAGGAACTAG



AG**CCCCTGTGAGGAACT**AG

Machine translation

Voulez-vous chanter avec  
moi?



Do you want to sing with  
me?

Image credit: <https://laptrinhx.com/understanding-of-recurrent-neural-networks-lstm-gru-3720007533/>

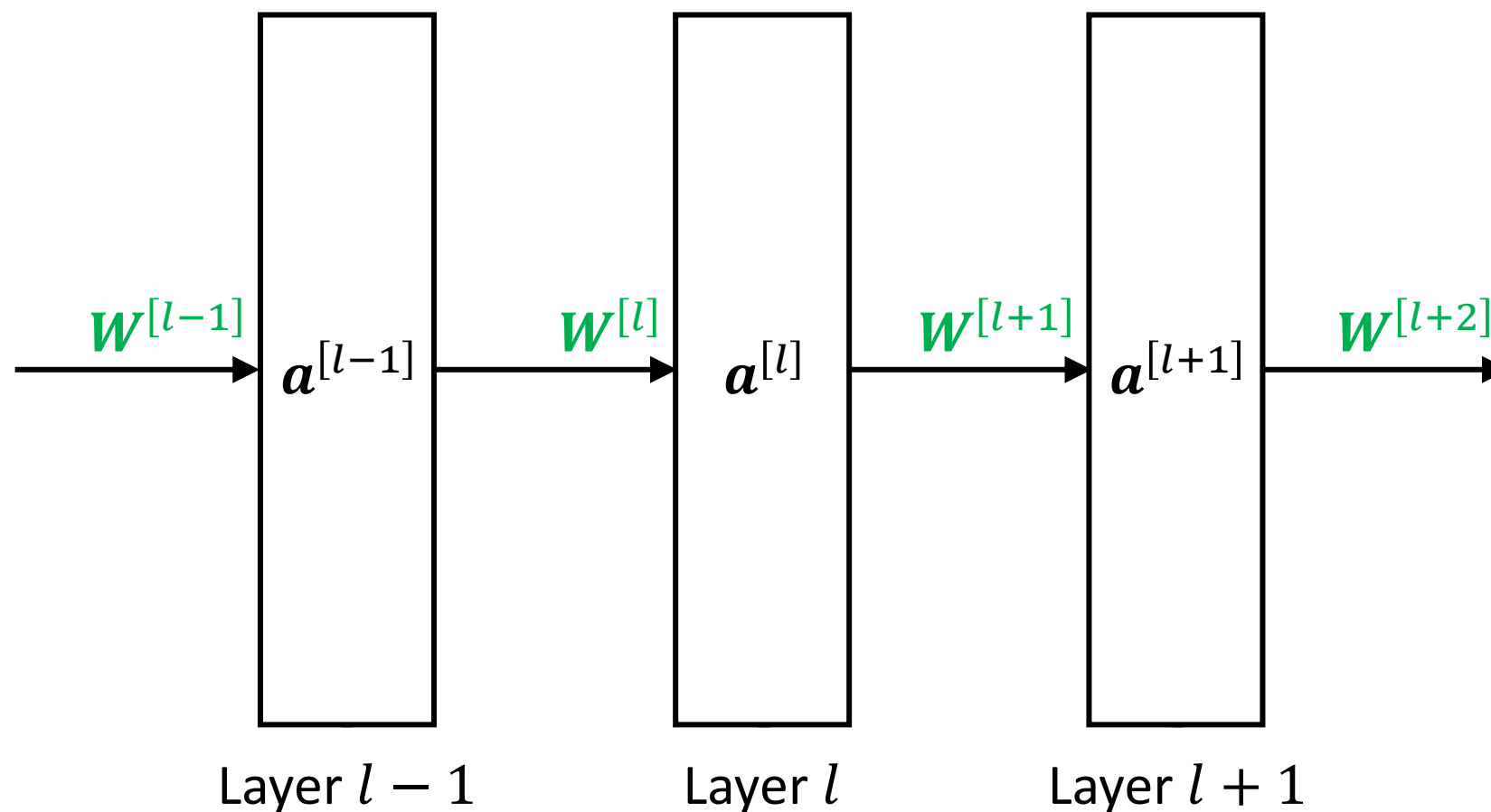
draw a shape, any shape.



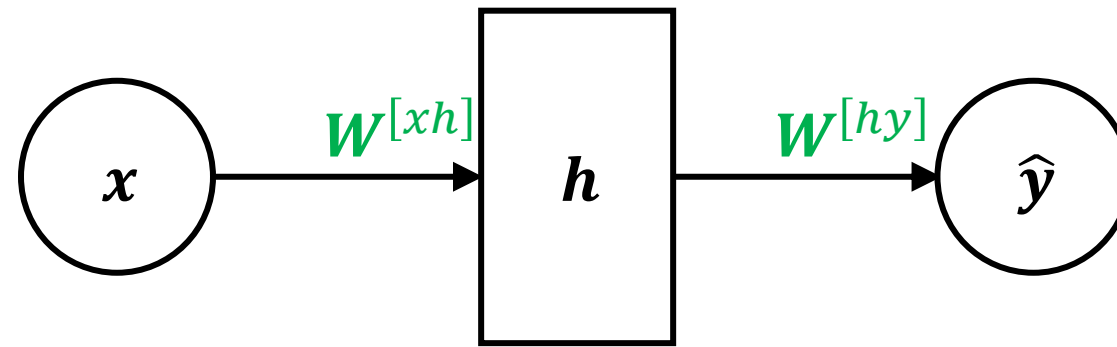
- **Draw something** and see what the AI will continue drawing
- **Take screenshot** and post your results
  - No obscenities please
- **Emote**
  - Up vote those you like
  - Down vote those with mistakes
- **Discuss** how you think the model predicted what to draw next

Try yourself: [https://magenta.tensorflow.org/assets/sketch\\_rnn\\_demo/index.html](https://magenta.tensorflow.org/assets/sketch_rnn_demo/index.html)

# Feedforward Neural Network



# Neural Network (simplified 1-hidden layer)

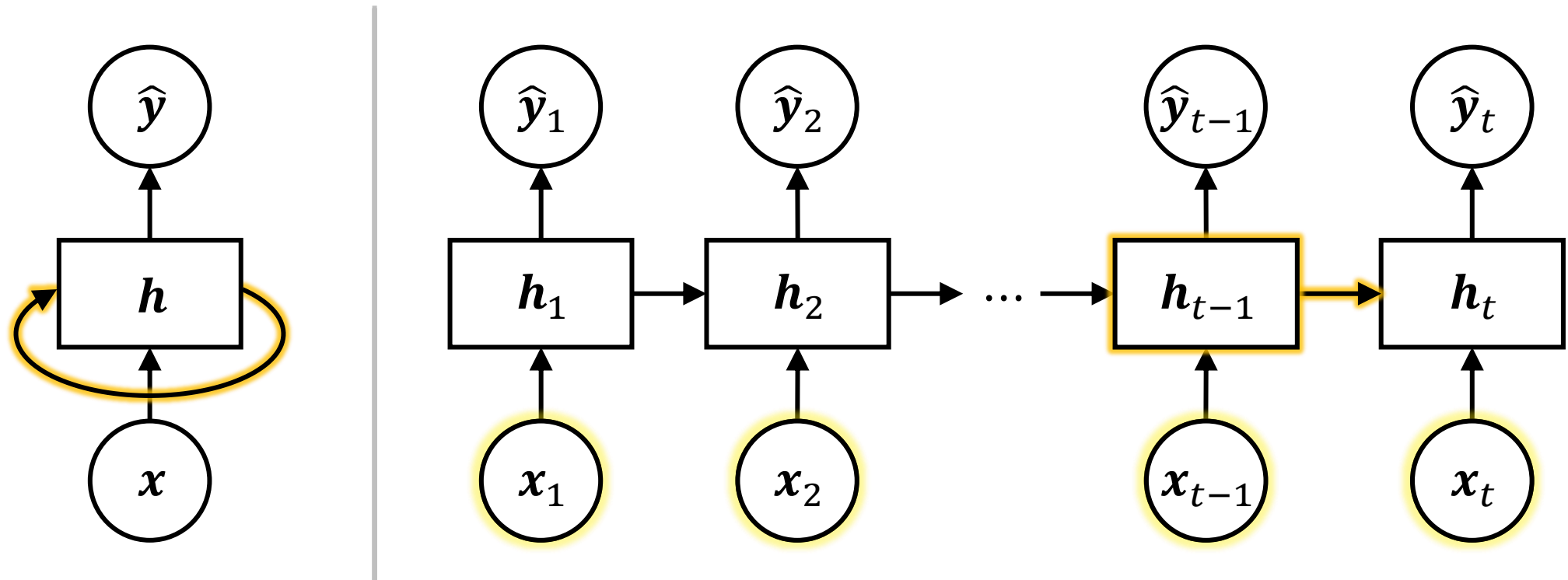


Input Layer

Hidden Layer

Output Layer

# Neurons with Recurrence

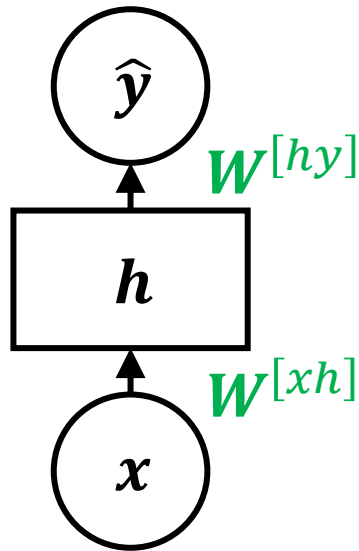


$$\hat{y} = g^{[y]} \left( g^{[h]}(x_t, \mathbf{h}_{t-1}) \right)$$

Recurrent Neural Network (RNN)

$$\hat{y}_t = g^{[y]}(\mathbf{h}_t)$$
$$\mathbf{h}_t = g^{[h]}(x_t, \mathbf{h}_{t-1})$$

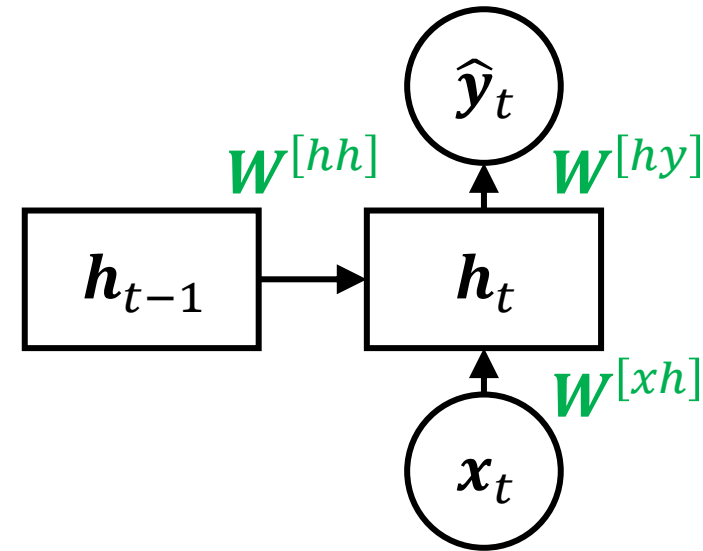
# RNN Weights



Feedforward Neural Network

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

**Question:** Do these **weights** change for different time  $t$ ?



Recurrent Neural Network

$$\hat{y}_t = g^{[y]} \left( (W^{[hy]})^\top h_t \right)$$
$$h_t = g^{[h]} \left( (W^{[xh]})^\top x_t + (W^{[hh]})^\top h_{t-1} \right)$$
$$= g^{[h]} \left( (W^{[xh]} \oplus W^{[hh]})^\top (x_t \oplus h_{t-1}) \right)$$

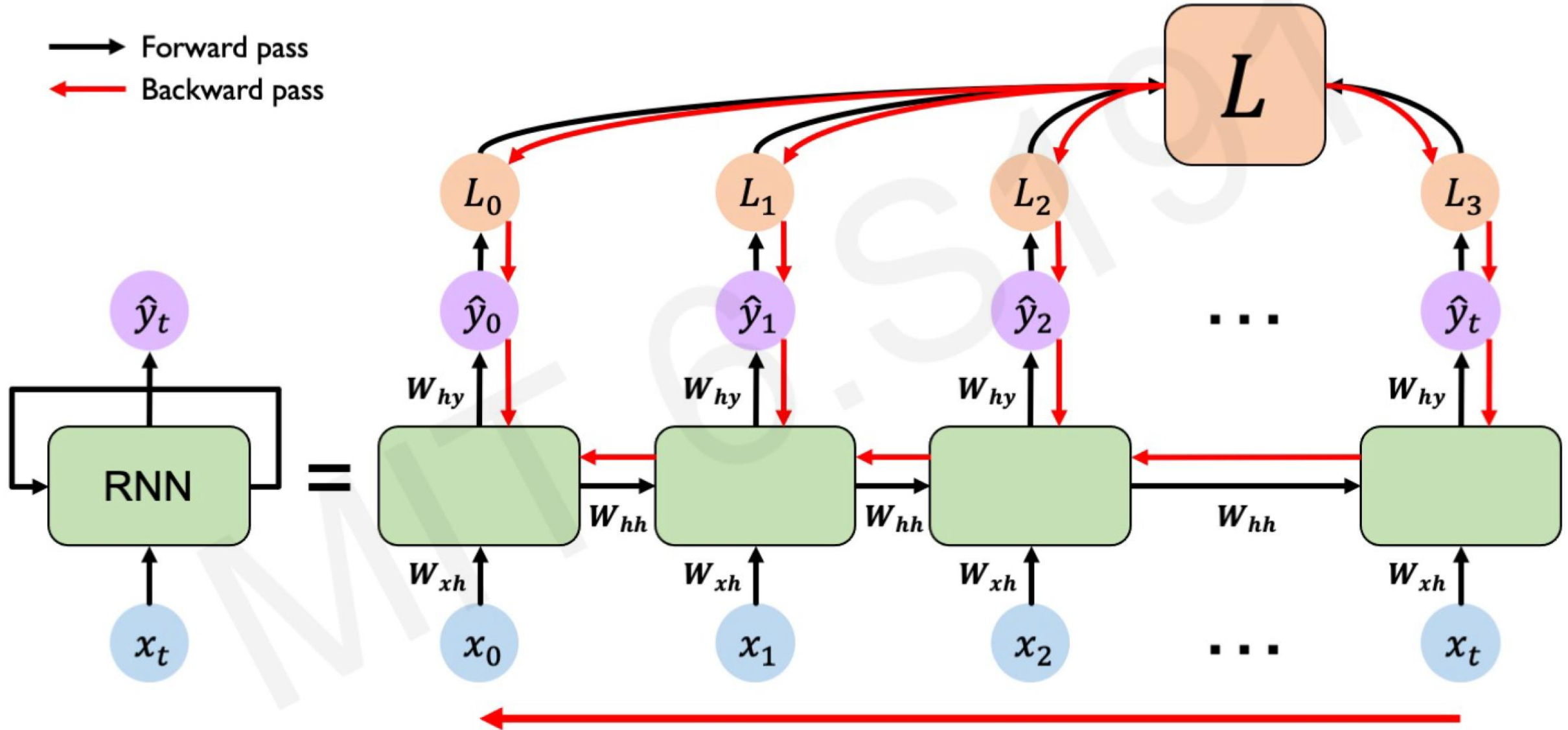




2021



# RNNs: Backpropagation Through Time



## Example RNN

# Training text prediction

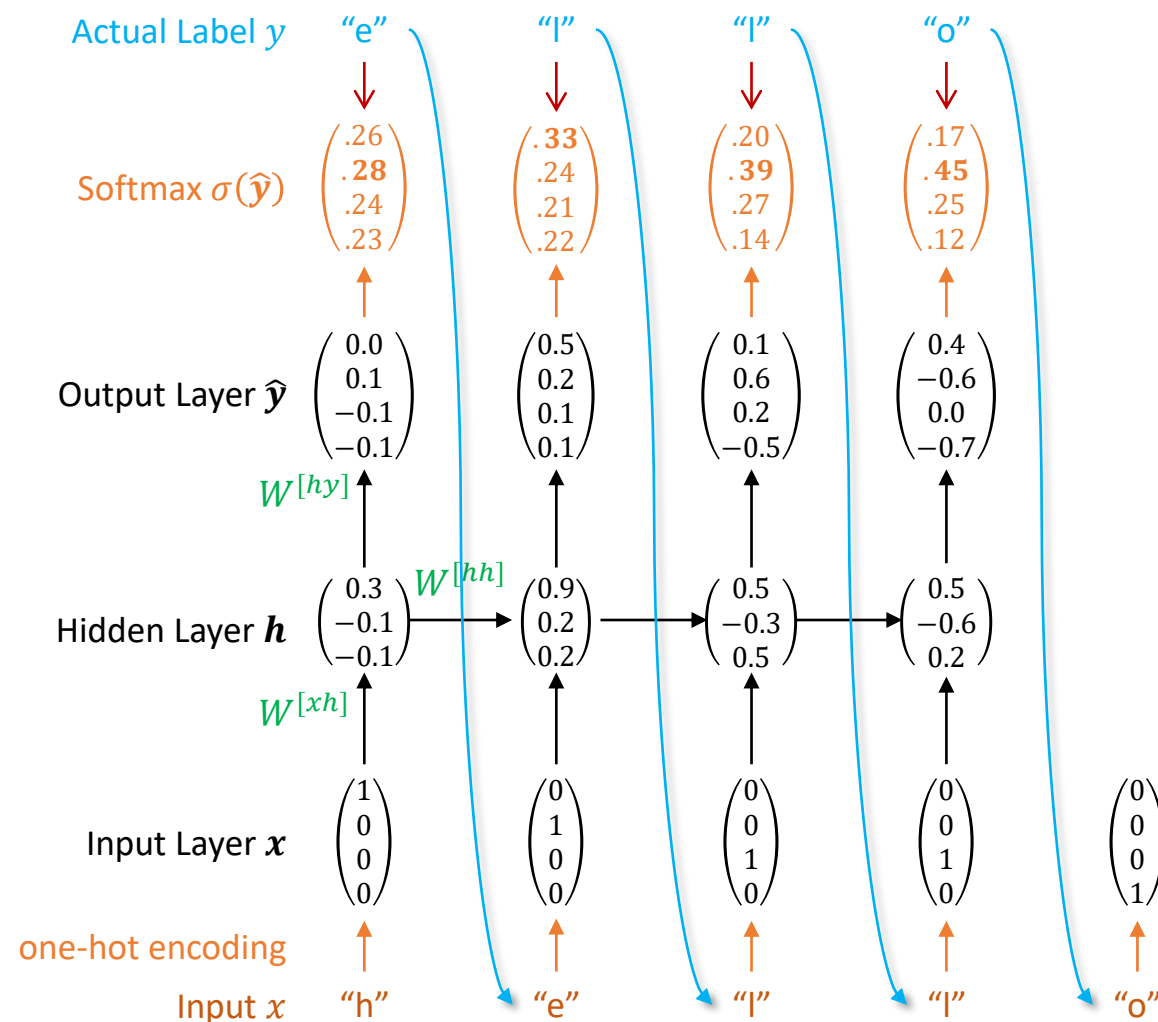
- Dictionary
  - [h, e, l, o]
- Encoding and Decoding chars
  - One-hot encoding (e.g., BOW)
  - Softmax classification
- At training time,
  - $x_t = y_{t-1}$
  - **Loss** is calculated as Cross-Entropy Error between  $\hat{y}_t$  and  $y_t$

$$W^{[xh]} = \begin{pmatrix} 0.3 & 1.0 & 0.1 & 0.5 \\ -0.1 & 0.3 & -0.3 & 0.1 \\ -0.1 & 0.1 & -0.5 & -0.5 \end{pmatrix}$$

$$W^{[hy]} = \begin{pmatrix} 0.3 & 0.9 & 0.1 \\ 0.2 & -0.6 & 0.5 \\ -0.1 & 0.1 & 0.5 \\ -0.1 & 1.0 & -0.2 \end{pmatrix}$$

$$W^{[hh]} = \begin{pmatrix} 0.1 & 0.4 & 0.8 \\ -0.1 & 0.5 & -0.2 \\ 0.9 & 0.2 & 0.6 \end{pmatrix}$$

$$\sigma(\hat{y}) = \frac{e^{\hat{y}}}{1 \cdot (1 + e^{\hat{y}})}$$

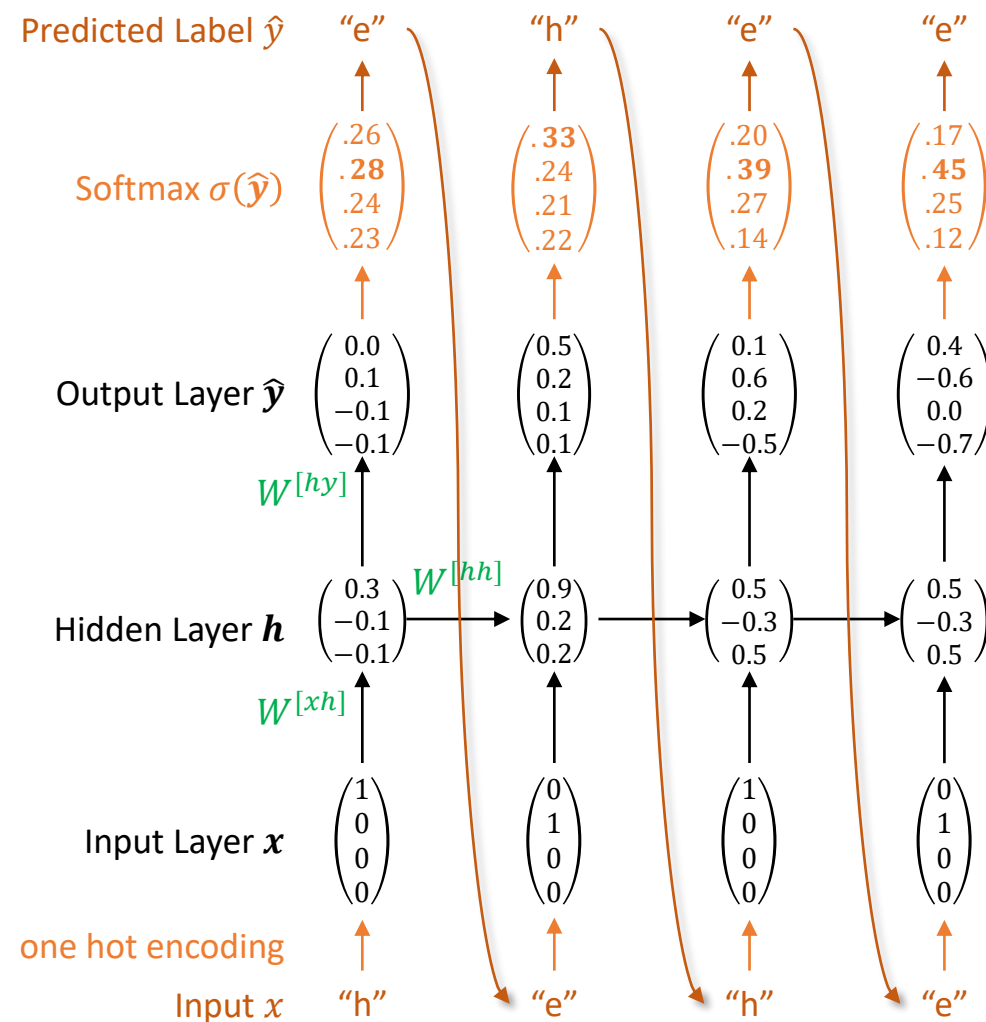


# Example RNN

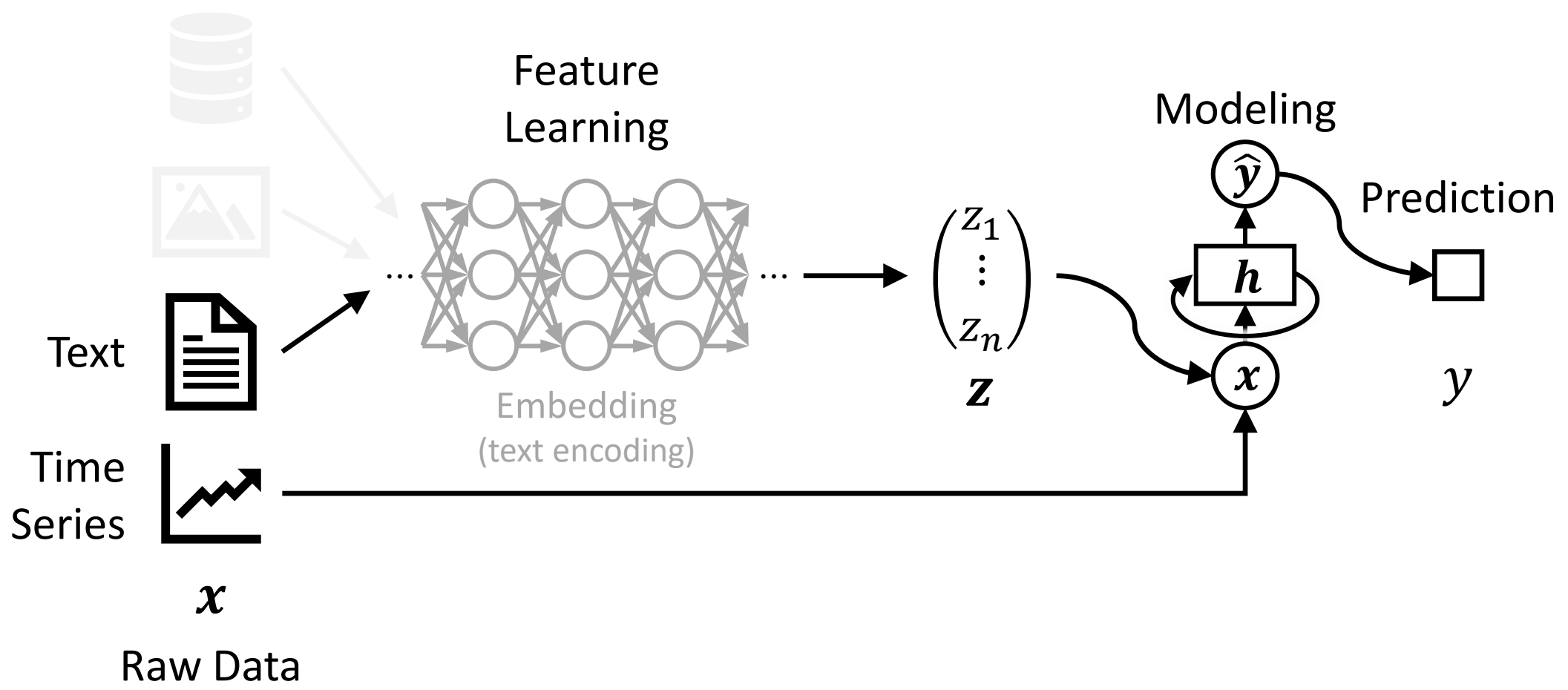
## Predicting Text

- At prediction time,
  - Forward propagate calculating activations to generate sequence of characters
  - $x_t = \hat{y}_{t-1}$

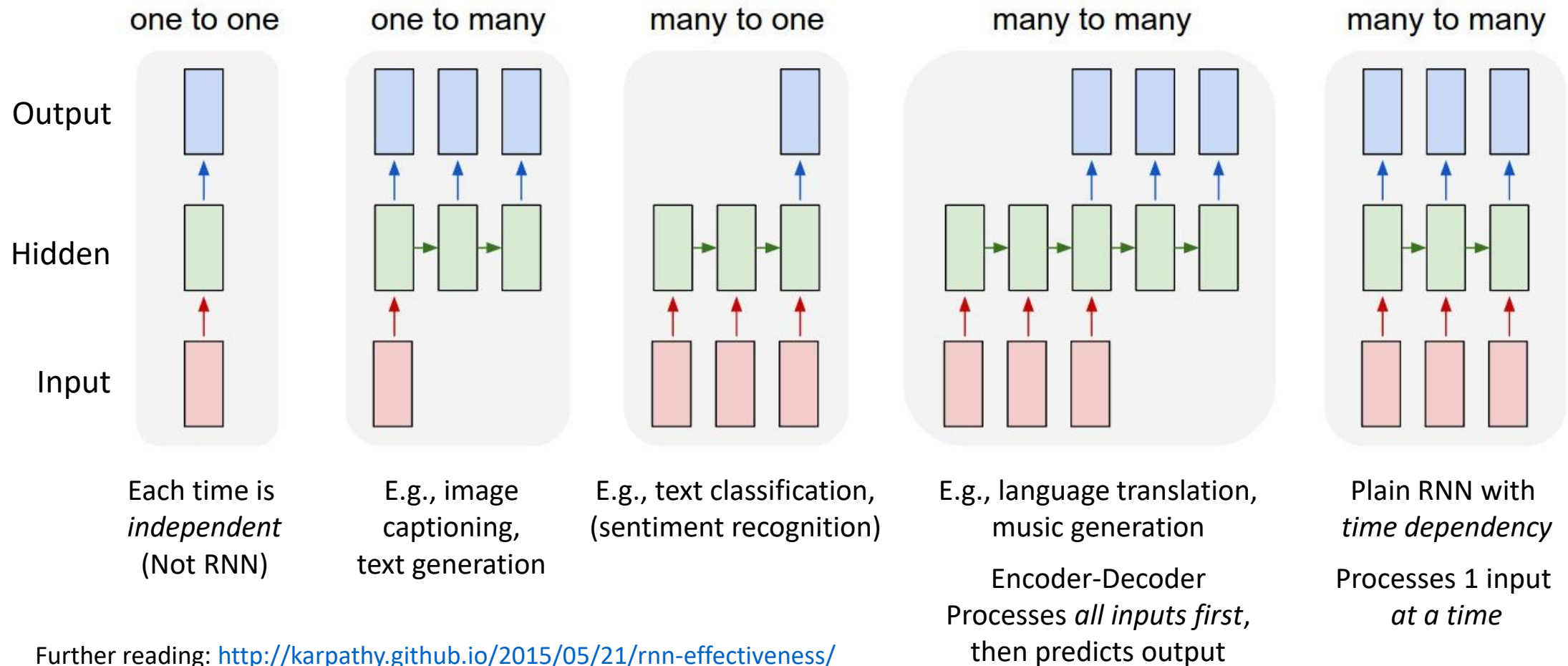
$$W^{[xh]} = \begin{pmatrix} 0.3 & 1.0 & 0.1 & 0.5 \\ -0.1 & 0.3 & -0.3 & 0.1 \\ -0.1 & 0.1 & -0.5 & -0.5 \end{pmatrix} \quad W^{[hy]} = \begin{pmatrix} 0.3 & 0.9 & 0.1 \\ 0.2 & -0.6 & 0.5 \\ -0.1 & 0.1 & 0.5 \\ -0.1 & 1.0 & -0.2 \end{pmatrix} \quad W^{[hh]} = \begin{pmatrix} 0.1 & 0.4 & 0.8 \\ -0.1 & 0.5 & -0.2 \\ 0.9 & 0.2 & 0.6 \end{pmatrix} \quad \sigma(\hat{y}) = \frac{\exp(\hat{y})}{\sum_{c=1}^C \exp(\hat{y}_c)}$$



# From Manual Feature Engineering To Automatic Feature Learning



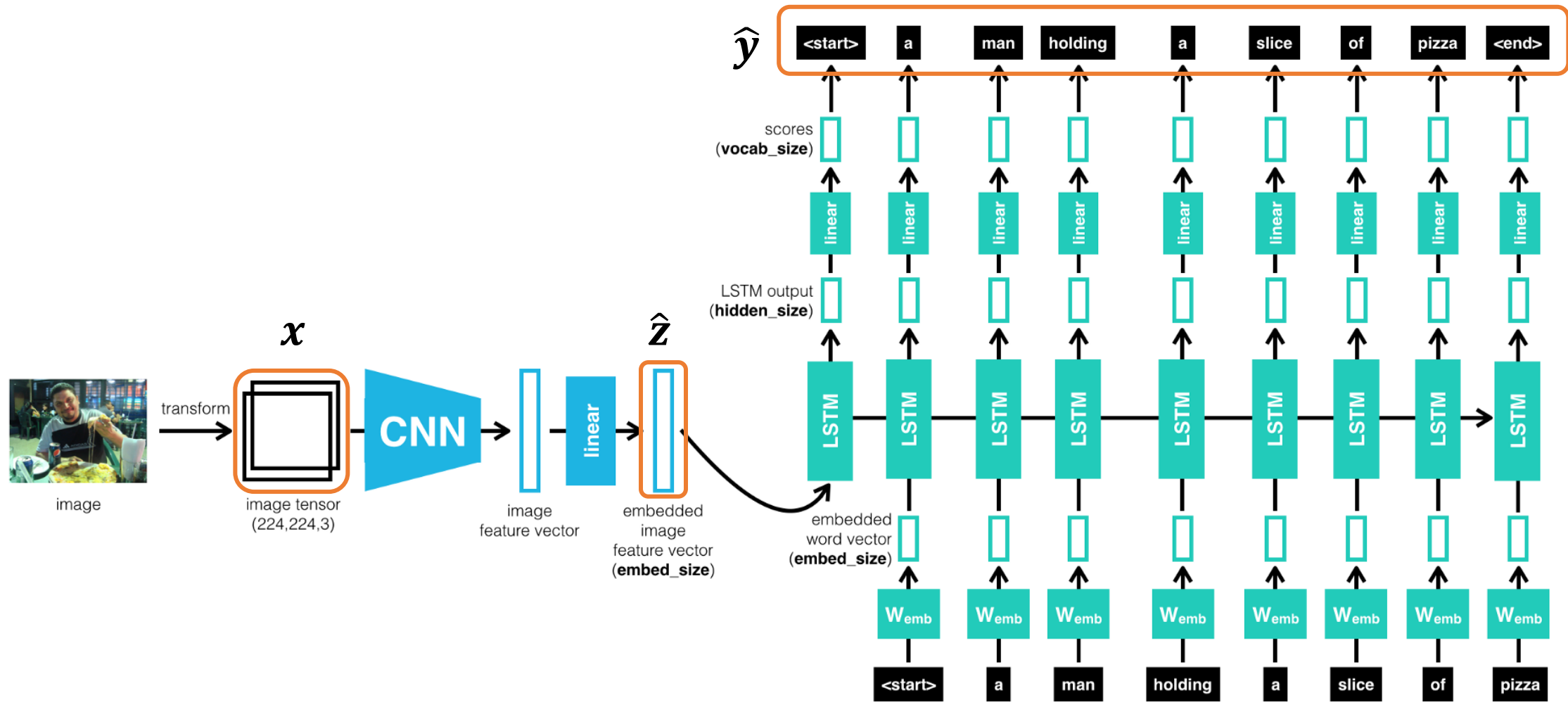
# Sequence Modeling Applications



Further reading: <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>



# Image Captioning: CNN + RNN (LSTM) – not in exam








*Questions!*





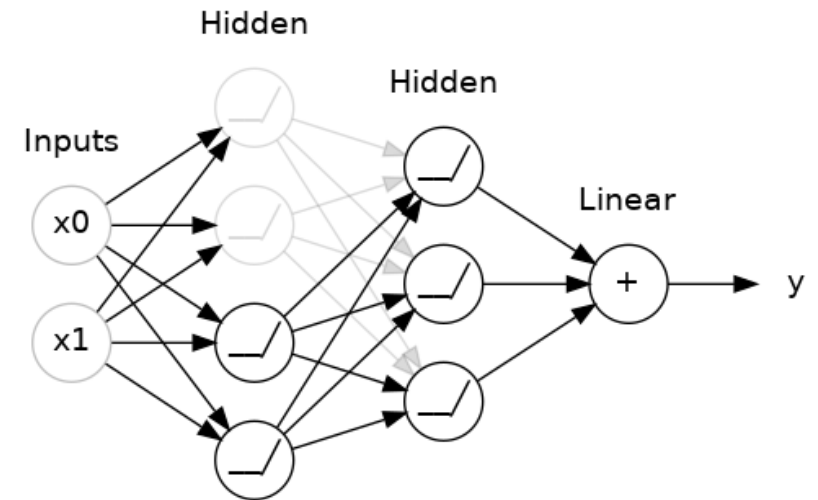
# Deep Learning Training Issues

# Deep Learning Training Issues

- Overfitting
- Saturating Gradient Problem
- Vanishing Gradient Problem

# Overfitting in deep neural networks

- Recall: what is overfitting?
  - Performance: Validation  $<$  Training
- Why can deep learning overfit?
  - Too **many** parameters!
  - Model is more expressive than needed
- Mitigation?
  - Dropout
    - **Randomly “drop out”** some neurons during batch **training**
    - **Cannot propagate** through those neurons during training
    - Note: **all nodes** are still used for prediction



Further reading: <https://towardsdatascience.com/12-main-dropout-methods-mathematical-and-visual-explanation-58cdc2112293>

# Deep Learning Training Issues

- Overfitting
- Saturating Gradient Problem
- Vanishing Gradient Problem

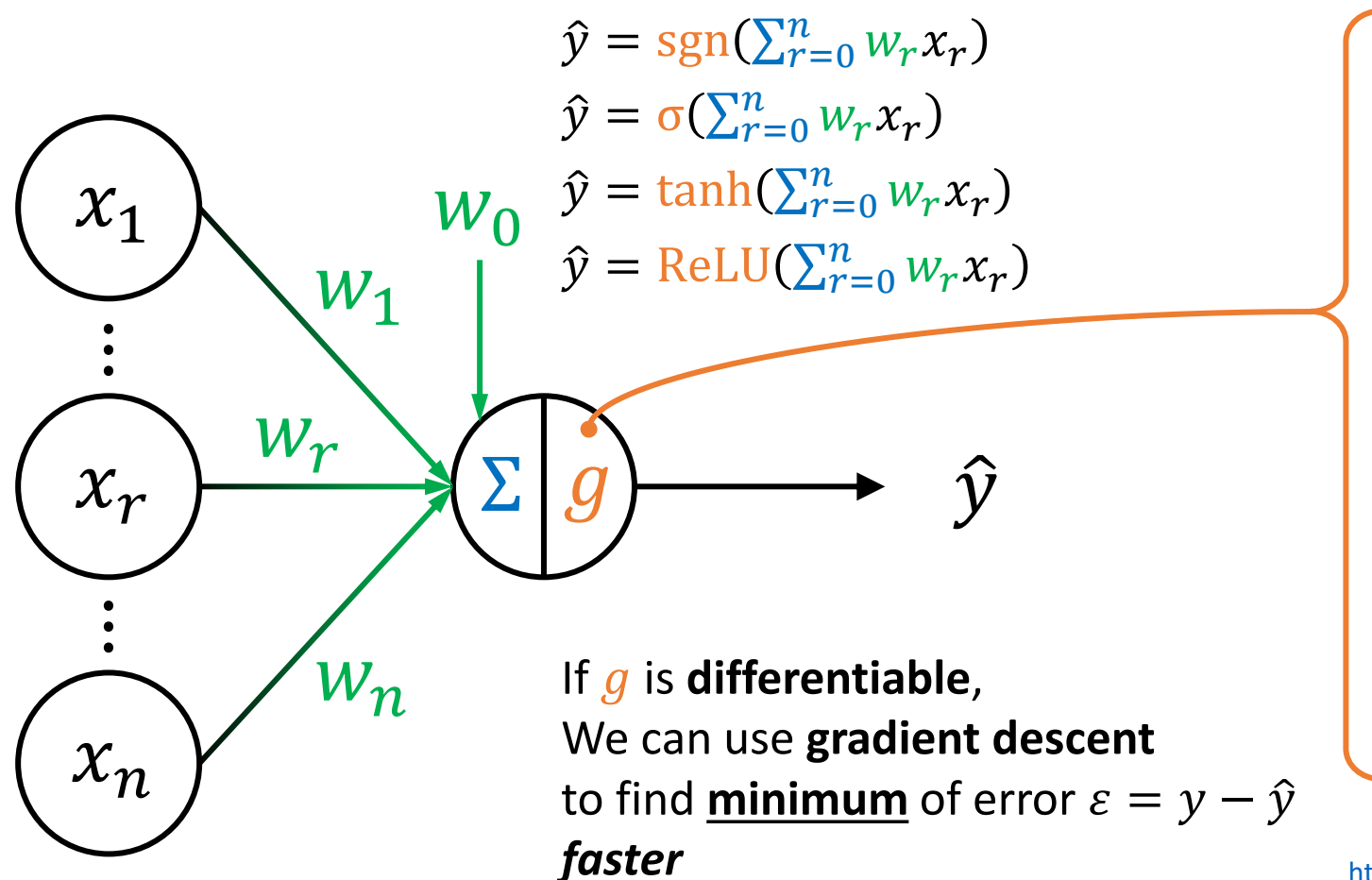
# Gradient Descent Weight Update (Neural Network)

$$\underset{\substack{\text{New} \\ \text{weight}}}{\mathbf{W}} \leftarrow \underset{\substack{\text{Old} \\ \text{weight}}}{\mathbf{W}} - \underset{\substack{\text{Learning} \\ \text{Rate}}}{\eta} \underset{\substack{\text{Direction of} \\ \text{fastest error} \\ \text{increase}}}{\nabla \epsilon}$$

Gradient of error

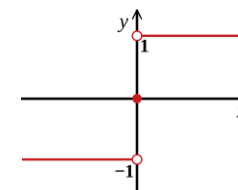
$$\nabla \epsilon = \frac{d\epsilon}{d\mathbf{W}} = \frac{d\epsilon}{d\hat{y}} \underbrace{\frac{d\hat{y}}{d\mathbf{W}}}_{\frac{df}{d\mathbf{W}} \frac{dg}{df}}$$

# Differentiable Activation Functions



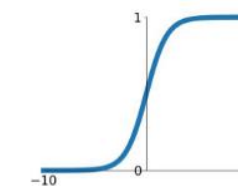
**Step**

$$\text{sgn}(x) = \begin{cases} +1 & z > 0 \\ -1 & z \leq 0 \end{cases}$$



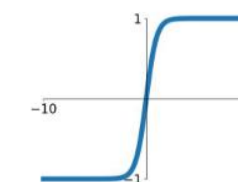
**Sigmoid**

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



**tanh**

$$\tanh(x)$$



**ReLU**

$$\max(0, x)$$

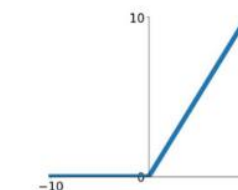


Image credit:

[https://miro.medium.com/max/1400/0\\*sIJ-gbjlz0zrz8lb.png](https://miro.medium.com/max/1400/0*sIJ-gbjlz0zrz8lb.png)

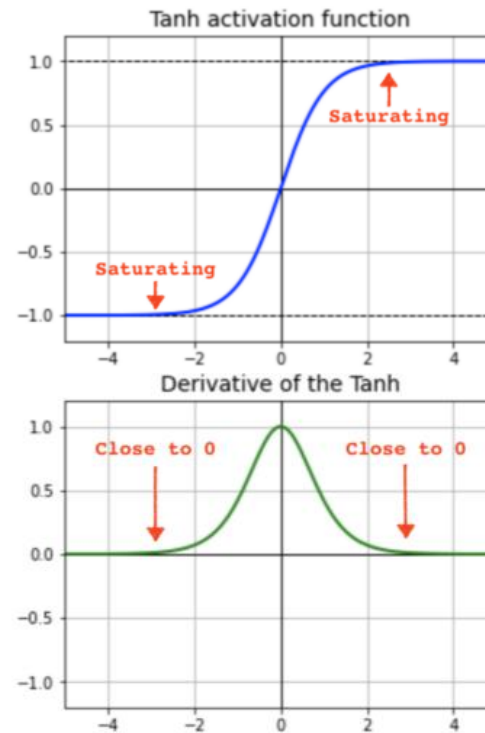
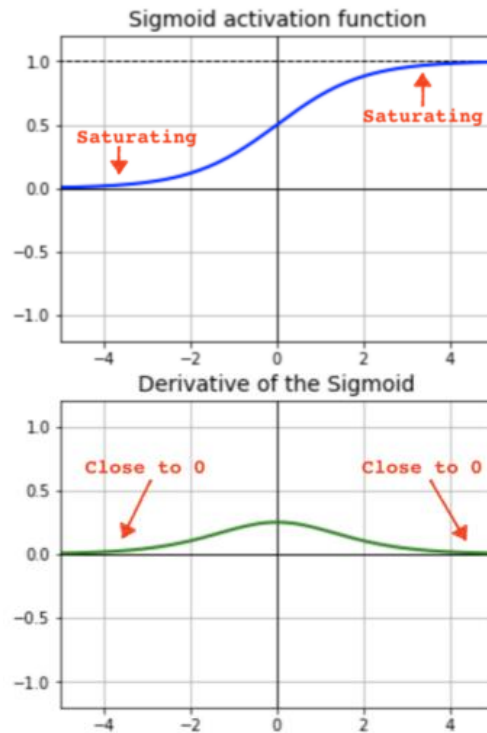


# Saturating Gradient Problem due to activation functions

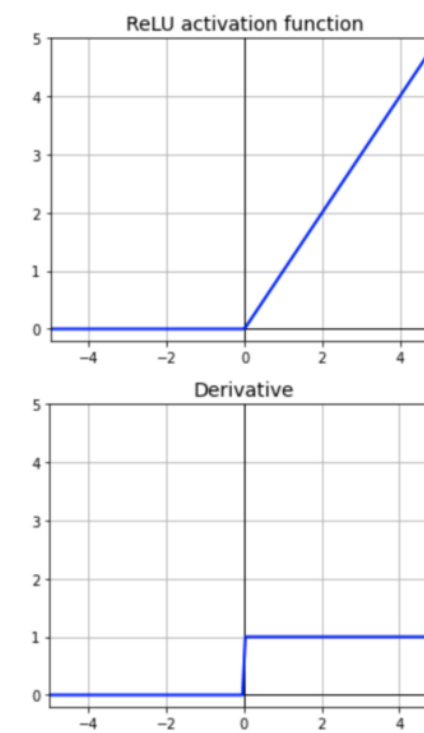
## Mitigate with ReLU activation function

$g(f)$

$\frac{\partial g}{\partial f}$



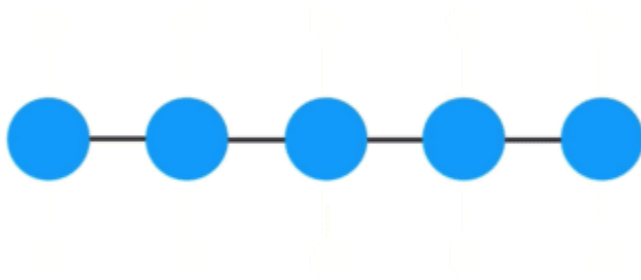
Mitigation



When  $x$  value far from 0, gradient  $\rightarrow 0$  (saturating)  
When gradient  $\approx 0$ , then  $\Delta W = \eta \nabla \varepsilon$  weights don't update much

With ReLU, gradient is always 1 (for  $x > 0$ )  
Can always update weights (for  $x > 0$ )

# Vanishing Gradient Problem



$$\frac{\partial \hat{y}}{\partial \mathbf{W}^{[1]}} = \frac{df^{[1]}}{d\mathbf{W}^{[1]}} \frac{dg^{[1]}}{df^{[1]}} \cdots \frac{dg^{[l]}}{df^{[l]}} \frac{df^{[l+1]}}{dg^{[l]}} \frac{dg^{[l+1]}}{df^{[l+1]}} \cdots \frac{df^{[L]}}{dg^{[L-1]}} \frac{dg^{[L]}}{df^{[L]}}$$

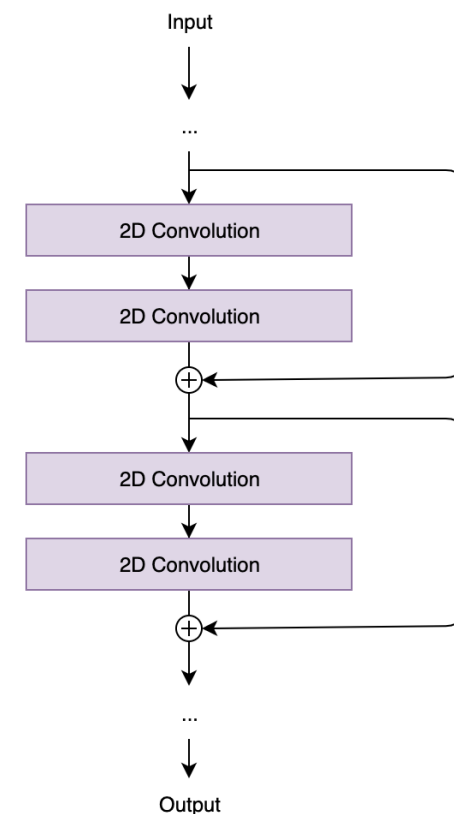
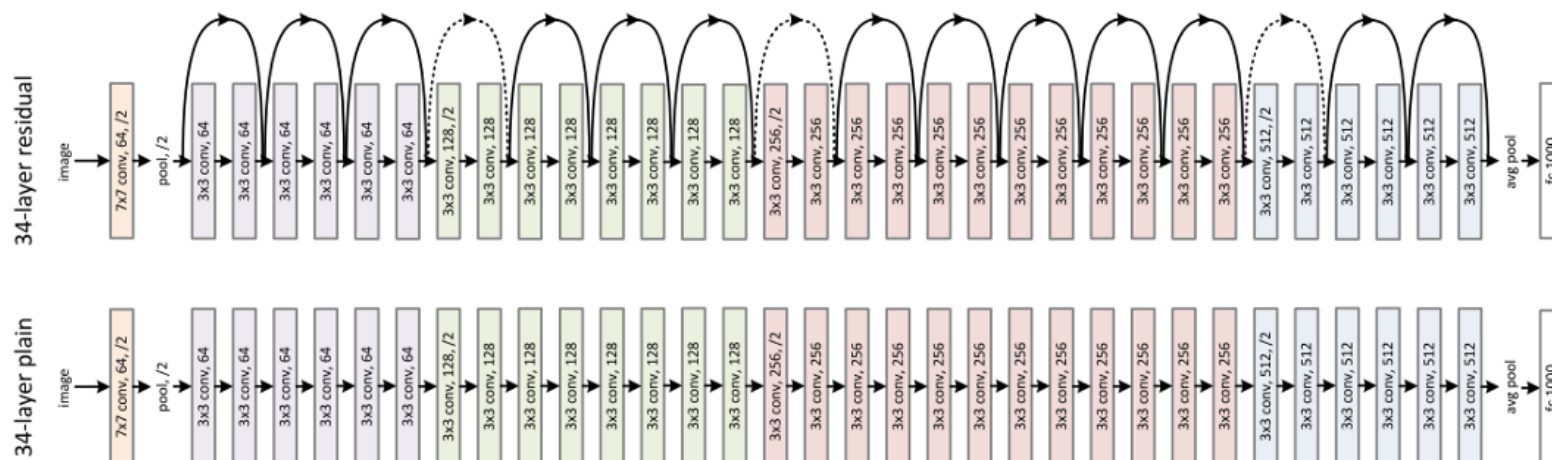
If some gradients are small ( $< 1$ ),  
multiplying **many small numbers** equals a **very small number**.

E.g.,  $0.5^{15} \approx 0.0003$

Image credit: <https://towardsdatascience.com/understanding-rnns-lstms-and-grus-ed62eb584d90>

# Mitigating Vanishing Gradients in CNN: Using architecture with “shortcut” connections

- ResNet (Residual Networks)
- Propagates residuals (forward) and gradients (backwards) through “**shortcut connections**”
- Gradients through shortcuts will not be as small



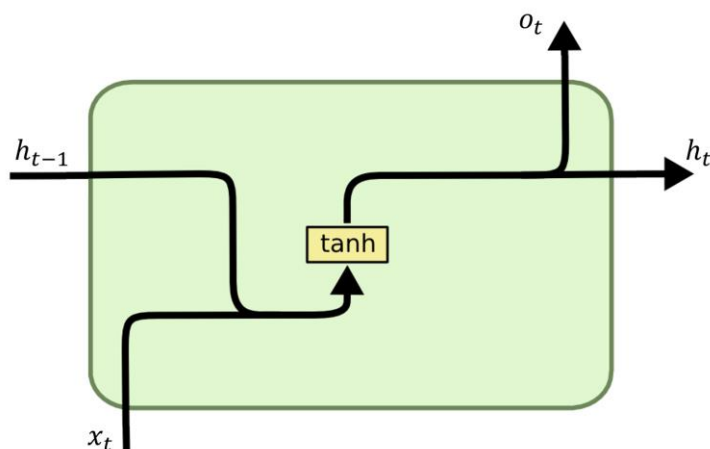
Further reading: <https://towardsdatascience.com/vggnet-vs-resnet-924e9573ca5c>

Image credit: <https://www.kaggle.com/keras/resnet50>

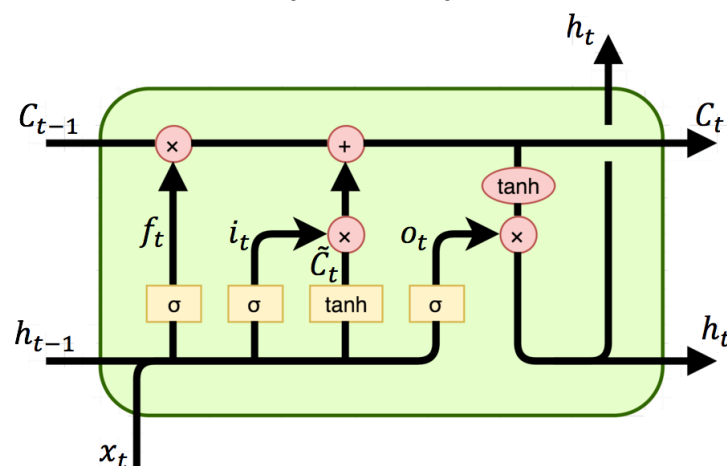
# Mitigating Vanishing Gradients in RNN

## Using architectures with “forget” gates

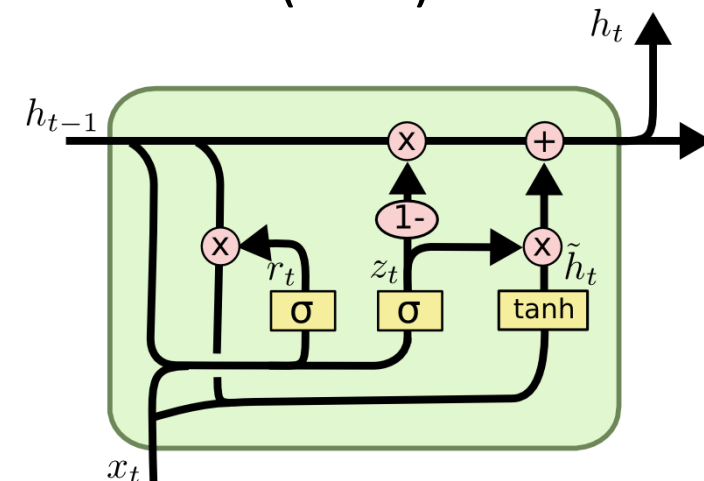
Plain RNN



Long-Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Includes “forget” gates

Image Credit: <http://dprogrammer.org/rnn-lstm-gru>

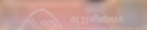
Further reading: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>



# Wrapping Up

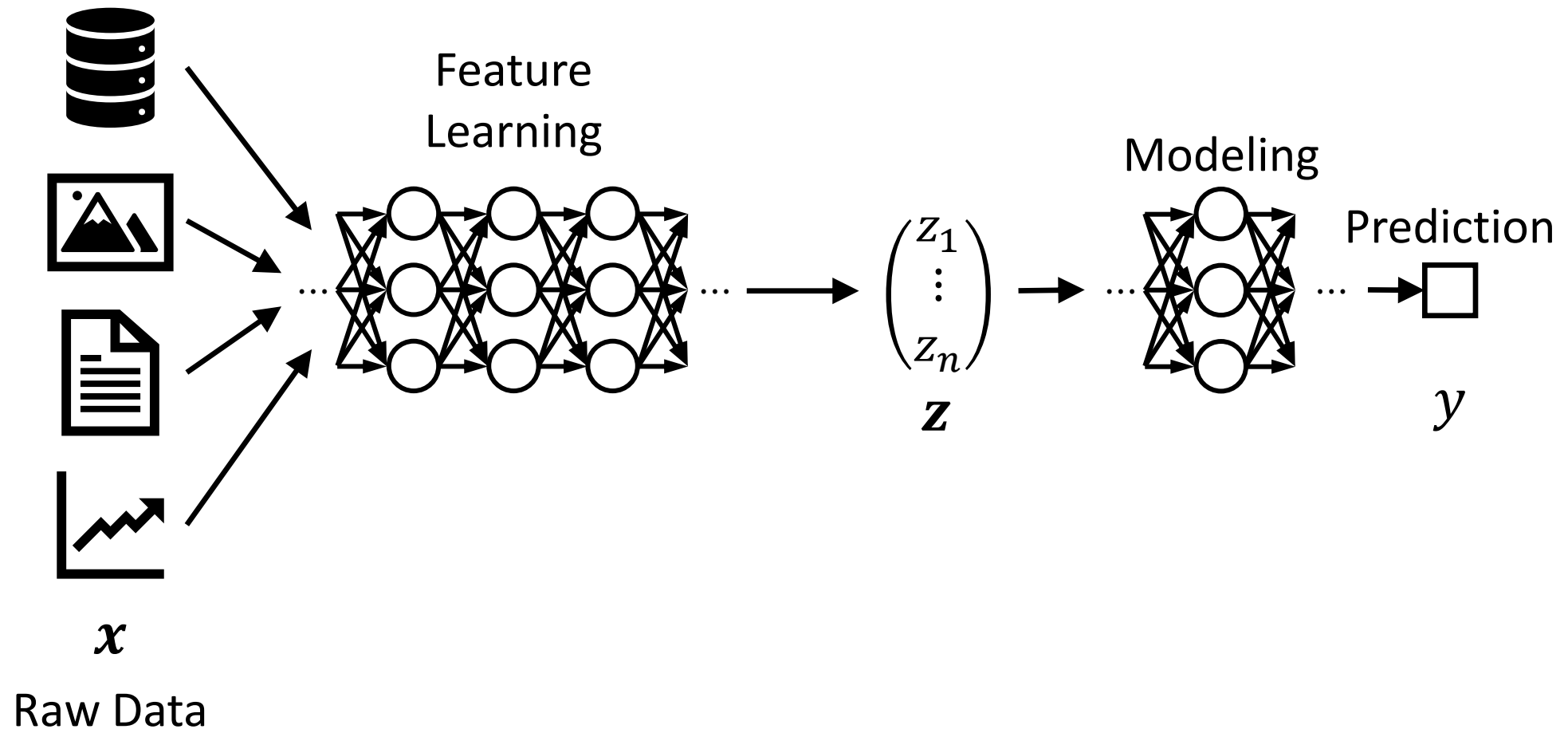


Department of Computer Science  
School of Computing



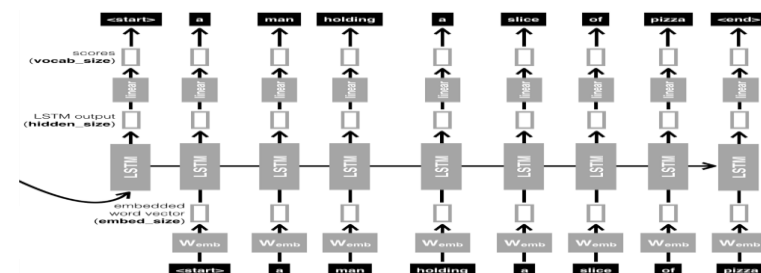
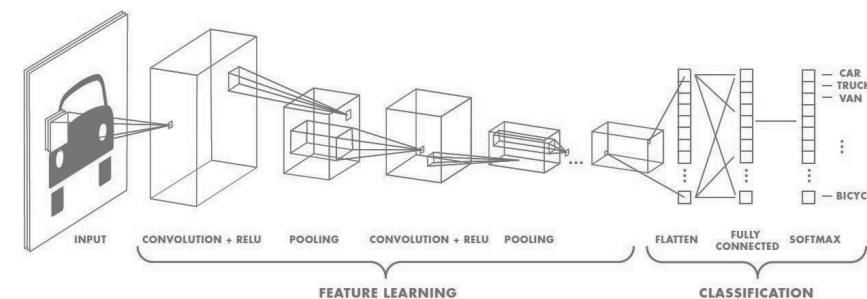
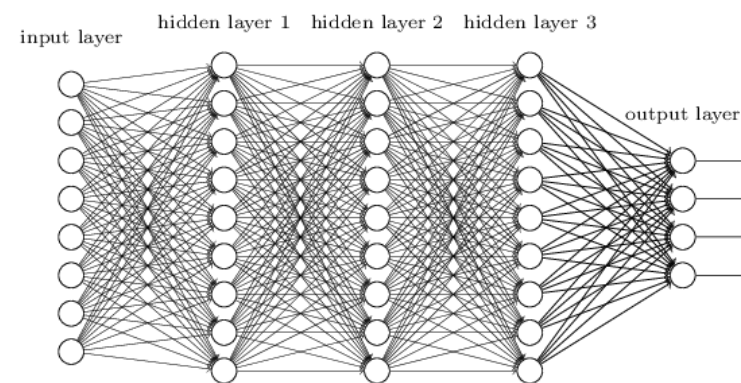
School of Computing

# From Manual Feature Engineering To Architecture Engineering



# What did we learn?

- Feature Engineering → Architecture Engineering
- **CNN**: exploits spatial information using **convolutions**
- **RNN**: exploits history information using **recurrence**





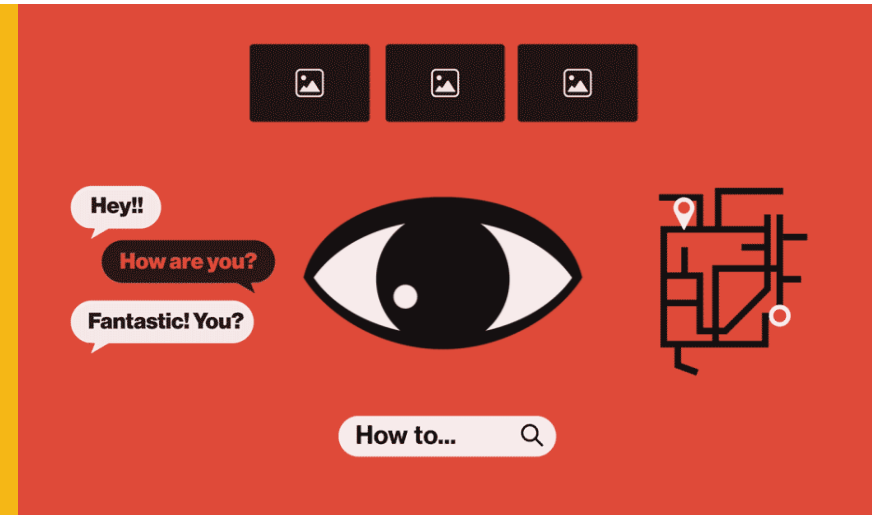
# Grand issues with AI (Deep Learning)



Lack of **Explainability**  
[W11b]



**Algorithmic Bias** (Societal)  
[W13a]

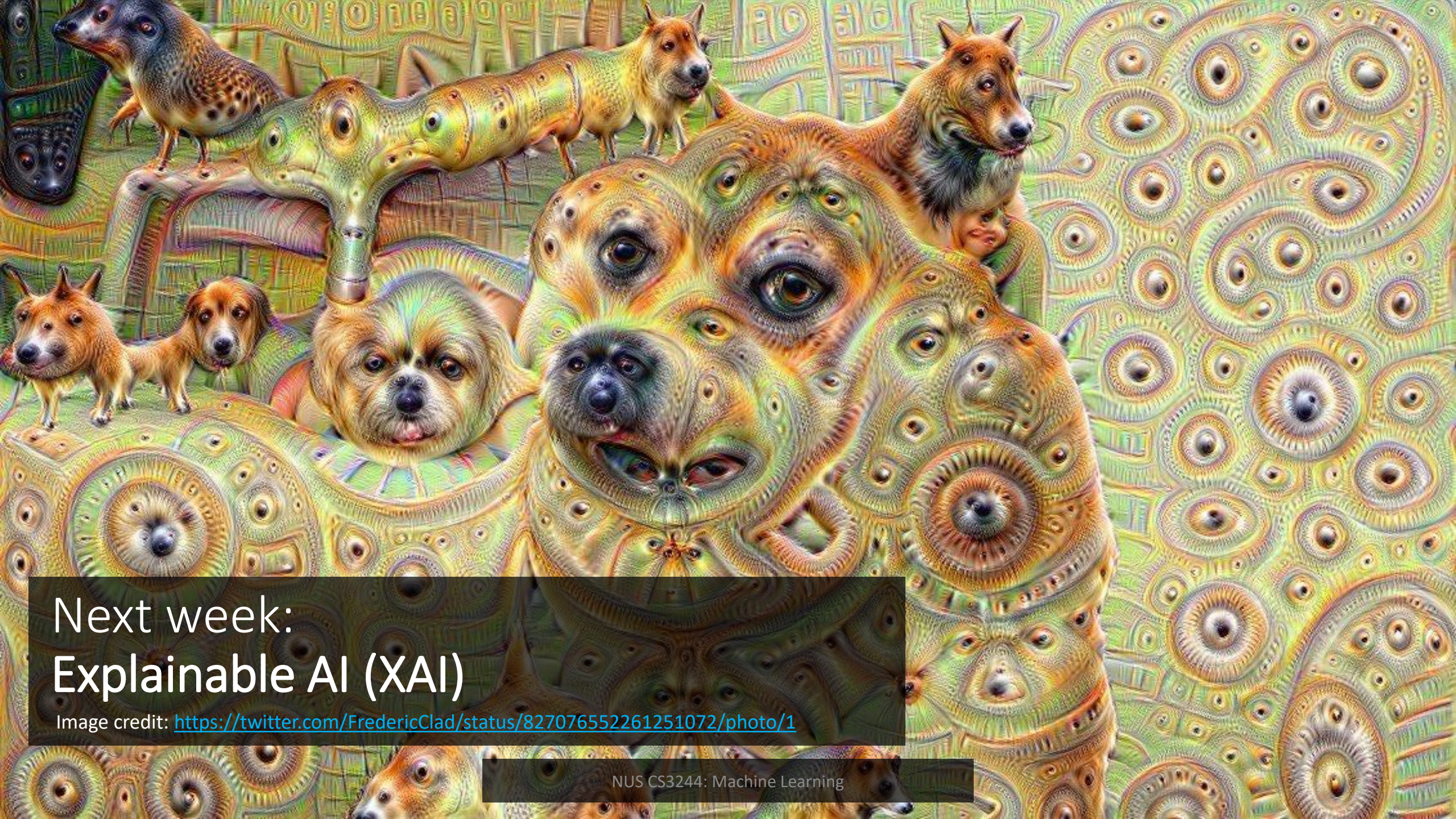


Data Privacy

Image credits:

[https://miro.medium.com/max/2000/1\\*H4cW-RCyHpu5FNtVaAPoQ.gif](https://miro.medium.com/max/2000/1*H4cW-RCyHpu5FNtVaAPoQ.gif)  
[https://www.insperity.com/wp-content/uploads/bias\\_1200x630.png](https://www.insperity.com/wp-content/uploads/bias_1200x630.png)  
<https://www.fightforprivacy.co/nuxt/img/512f421.gif>





# Next week: Explainable AI (XAI)

Image credit: <https://twitter.com/FredericClad/status/827076552261251072/photo/1>





# Following week: Unsupervised Learning

Image credit: <https://hip2save.com/2019/11/27/lego-classic-creative-fun-900-piece-set-only-20-at-walmart-regularly-40/>



