

Deep Learning

10

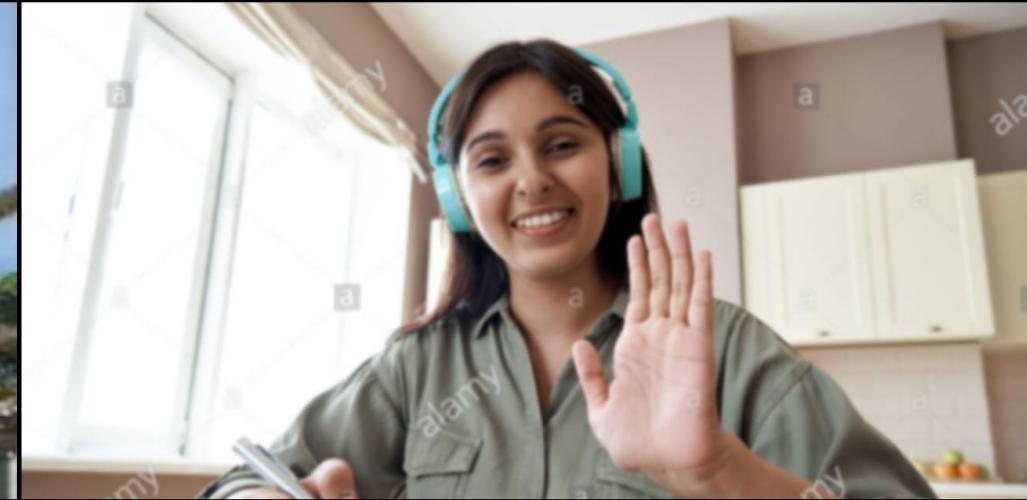
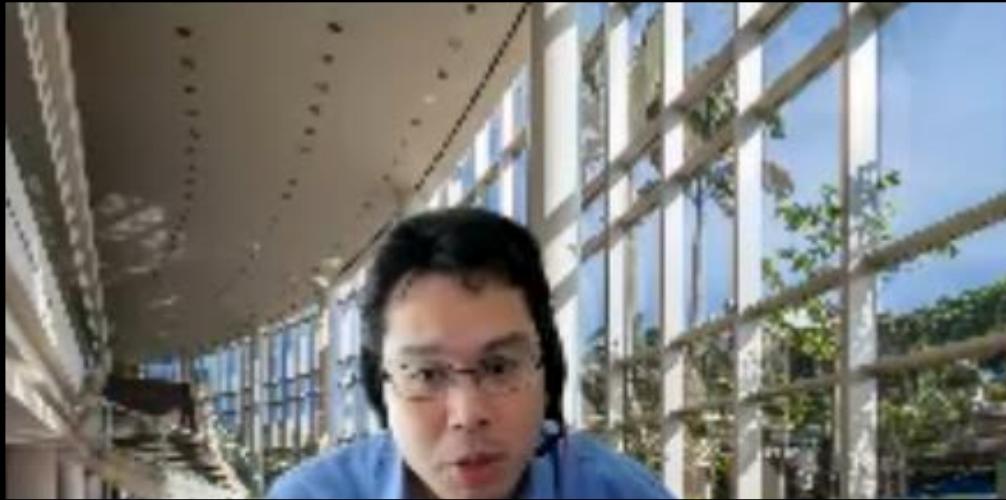
A

CS 3244
Machine Learning



NUS | Computing

National University
of Singapore



Please turn on your webcam



Mystery Student

SUSTAINABILITY IRL:

Love Data Science?
Concerned about
sustainability?

WE WANT YOU!



Partner Treatsure, Charlotte Mei, and Upcircle to unlock data-driven insights around topics like food waste and individual impact.

The challenge:

Creatively visualise an answer to their problem statements, using our open dataset of >100k conversations. Or, create your own dataset.

Prizes:

Up to \$1,000 to be won

JOIN US AT THE
KICK-OFF EVENT

28th Oct (Thursday),
6 - 7.30pm Online

Register here

sustainabilityirl.
synthesis.partners



Participate:

<https://www.sustainabilityirl.synthesis.partners/>

Hosted by

synthesis

Our Partners

XDS Experimental Data Science

treatsure



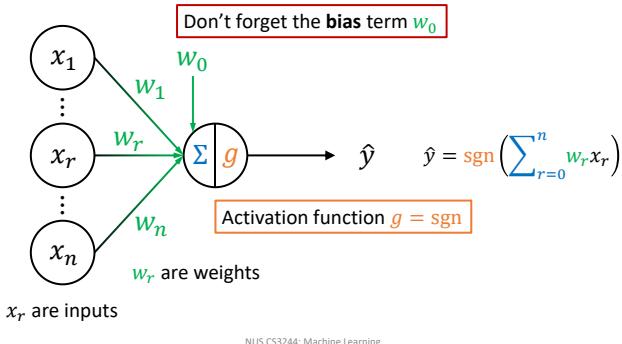
TheCharlotteMei

Mid-Semester Anonymous Survey



- What worked
 - Slack is good for comms & awareness
 - Exercises (pre-lecture, in-class) are engaging
 - Easier midterm (less stressful)!
- What to improve
 - Want more coding teaching: Bonus programming lectures to be conducted by Prof Min and TAs.
 - Want more examples: Included in lecture.
 - Somewhat out-of-scope tutorial questions: We will align tutorial questions more closely with lectures.

Perceptron



NUS CS3244: Machine Learning

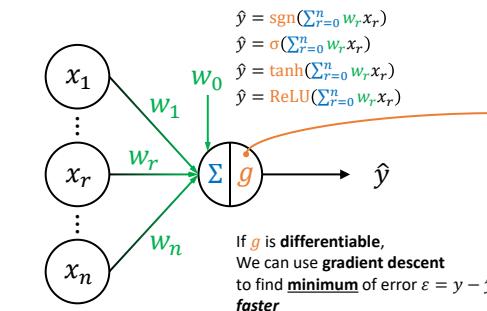
Perceptron Learning Algorithm

1. Initialize weights \mathbf{w}
 - Could be all zero, or random small values
2. For each instance i with features $\mathbf{x}^{(i)}$
 - Classify $\hat{y}^{(i)} = \text{sgn}(\mathbf{w}^\top \mathbf{x}^{(i)})$
3. Select one misclassified instance
 - Update weights: $\mathbf{w} \leftarrow \mathbf{w} + \eta(y - \hat{y})\mathbf{x}$
4. Iterate steps 2 to 3 until
 - Convergence (classification error < threshold), or
 - Maximum number of iterations

$$\begin{aligned} \mathbf{w}_0 &\leftarrow \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{pmatrix} + \eta(y - \hat{y}) \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{pmatrix} \\ w_r &\leftarrow w_r + \eta(y - \hat{y})x_r \end{aligned}$$

NUS CS3244: Machine Learning

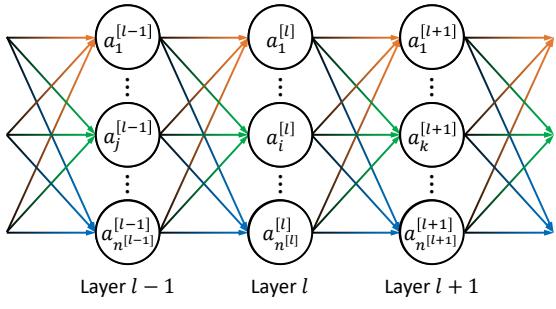
Differentiable Activation Functions



https://miro.medium.com/max/1400/0*sl-gbj0rrzlb.png

30

Multi-Layer Perceptron (Neural Network)



NUS CS3244: Machine Learning

Chain Rule

Consider composite function

$$g(x) = g(f(x))$$

$$g = g(f), f = f(x)$$

$$g'(x) = \frac{dg}{dx} = \frac{dg}{df} \frac{df}{dx}$$

Intuition

Rate of change of g relative to x is the product of

- rates of change of g relative to f and
- rates of change of f relative to x

"if

- a car travels 2x fast as a bicycle and
- the bicycle is 4x as fast as a walking man,

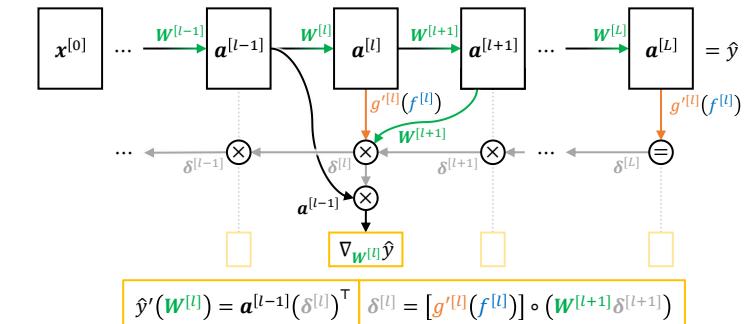
then the car travels $2 \times 4 = 8$ times as fast as the man."

– George F. Simmons, Calculus with Analytic Geometry (1985)

NUS CS3244: Machine Learning

22

Backward Propagation



NUS CS3244: Machine Learning

23

Notation

n = Number of features in x
 m = Number of instances in dataset

- **Scalar:** not bolded, lower case

x

- **Vector:** bolded, lower case

$$\boldsymbol{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

- **Matrix:** bolded, upper case

$$\boldsymbol{X} = \begin{pmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{pmatrix}$$

Functions with Vectors and Matrices

- Scalar-by-scalar:
 - $y(x) = wx$ for scaling input
- Scalar-by-vector:
 - $y(x) = \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^\top \mathbf{x} = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = w_1 x_1 - \cancel{w_1 x_2} - \cancel{w_2 x_1} + w_2 x_2$ for weighted sum
- Vector-by-vector:
 - $\mathbf{y}(x) = w\mathbf{x} = w \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} wx_1 \\ wx_2 \end{pmatrix}$ for scaled outputs (same weight)

Even more Chain Rule:

Gradient of Neural Network

$$\hat{y}(\mathbf{x}) = g^{[L]}(f^{[L]}(g^{[L-1]}(\cdots(g^{[l]}(f^{[l]}(g^{[l-1]}(\cdots(g^{[1]}(f^{[1]}(\mathbf{x}^{[0]})\cdots)\cdots)\cdots)\cdots)))\cdots)))$$

Gradient relative to \mathbf{W}

$$\hat{y}'(\mathbf{W}^{[L-1]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[L-1]}} = \frac{\partial f^{[L]}}{\partial \mathbf{W}^{[L]}} \frac{\partial g^{[L]}}{\partial f^{[L]}} \delta^{[L-1]}$$

Reference

$$\mathbf{a}^{[l]} = g^{[l]}(f^{[l]})$$

$$\hat{y}'(\mathbf{W}^{[l+1]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l+1]}} = \frac{\partial f^{[l+1]}}{\partial \mathbf{W}^{[l+1]}} \underbrace{\frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \cdots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}}} \delta^{[l+1]}$$

$$f^{[l]} = (\mathbf{W}^{[l]})^\top \mathbf{a}^{[l-1]}$$

$$\hat{y}'(\mathbf{W}^{[l]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l]}} = \frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} \frac{\partial g^{[l]}}{\partial f^{[l]}} \underbrace{\frac{\partial f^{[l+1]}}{\partial g^{[l]}} \frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \cdots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}}} \delta^{[l+1]}$$

Recursive

$$\hat{y}'(\mathbf{W}^{[l]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l]}} = \frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} \frac{\partial g^{[l]}}{\partial f^{[l]}} \frac{\partial f^{[l+1]}}{\partial g^{[l]}} \delta^{[l+1]}$$

$$\hat{y}'(\mathbf{W}^{[1]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[1]}} = \frac{\partial f^{[1]}}{\partial \mathbf{W}^{[1]}} \frac{\partial g^{[1]}}{\partial f^{[1]}} \cdots \frac{\partial g^{[l]}}{\partial f^{[l]}} \frac{\partial f^{[l+1]}}{\partial g^{[l]}} \frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \cdots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}}$$

Week 10A: 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 10A: Lecture Outline

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

Deep Neural Network



NUS
National University
of Singapore

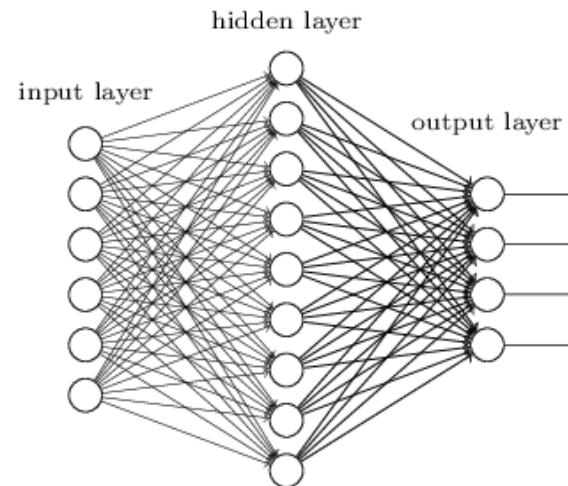
Department of Computer Science
School of Computing



SCHOOL OF COMPUTING

Deep Neural Network = many hidden layers (≥ 3)

Shallow Network



Deep Network

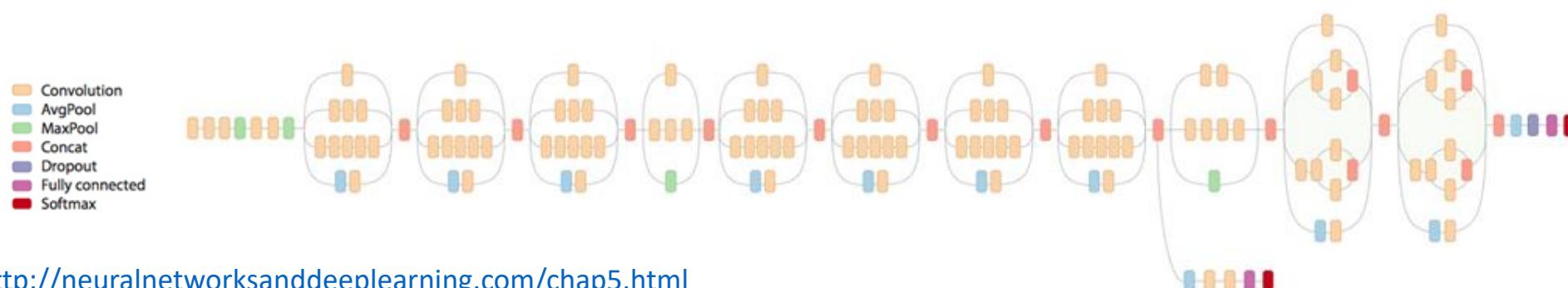
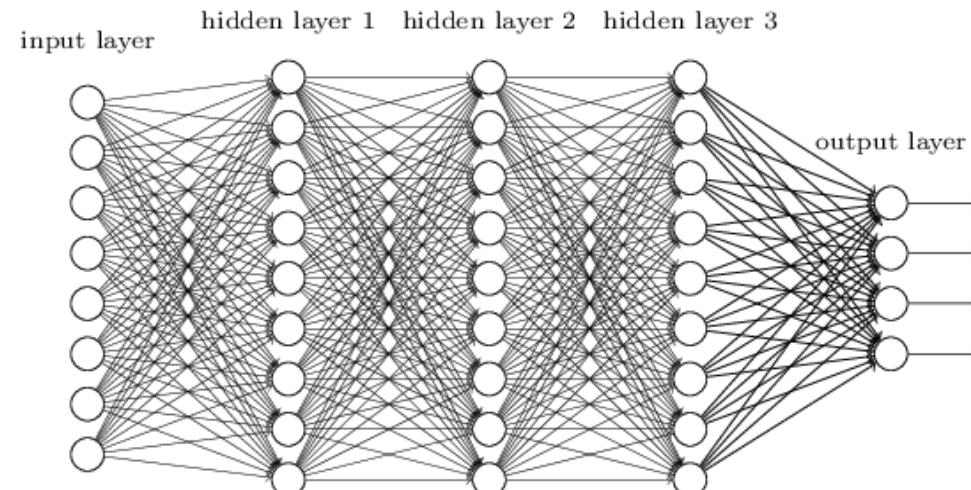
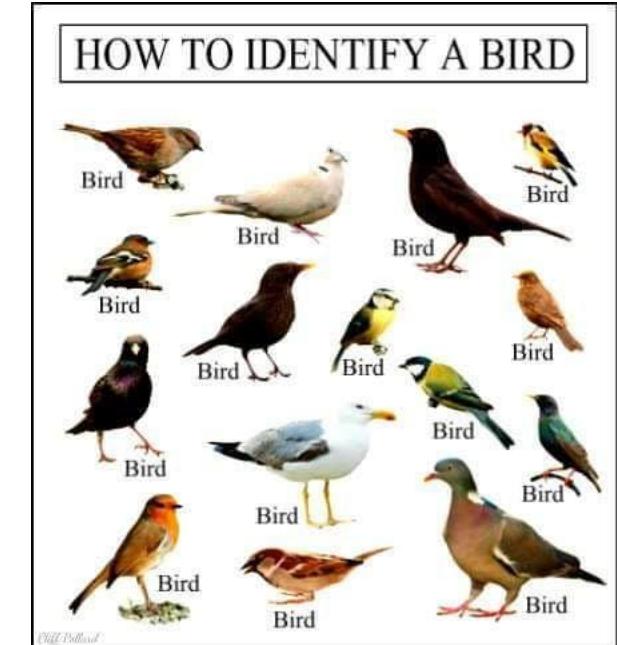


Image credit: <http://neuralnetworksanddeeplearning.com/chap5.html>

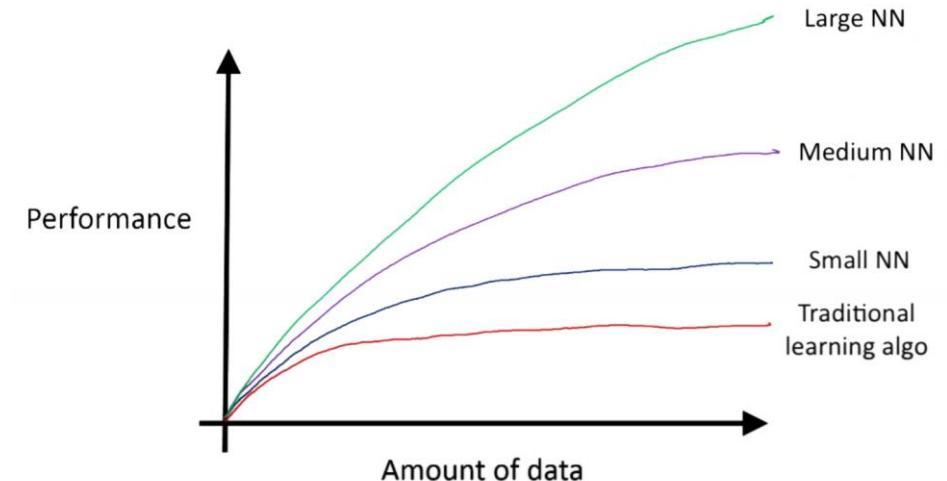
<https://adeshpande3.github.io/The-9-Deep-Learning-Papers-You-Need-To-Know-About.html>

Why Deep?

- Why need **so many layers**?
 - Need **many parameters**
 - Target functions of real-world tasks are **complex**
 - E.g., what is the function for recognizing birds or language?
- Why need **so much training data**?
 - Many parameters → Need more data
 - More data → **Better performance**

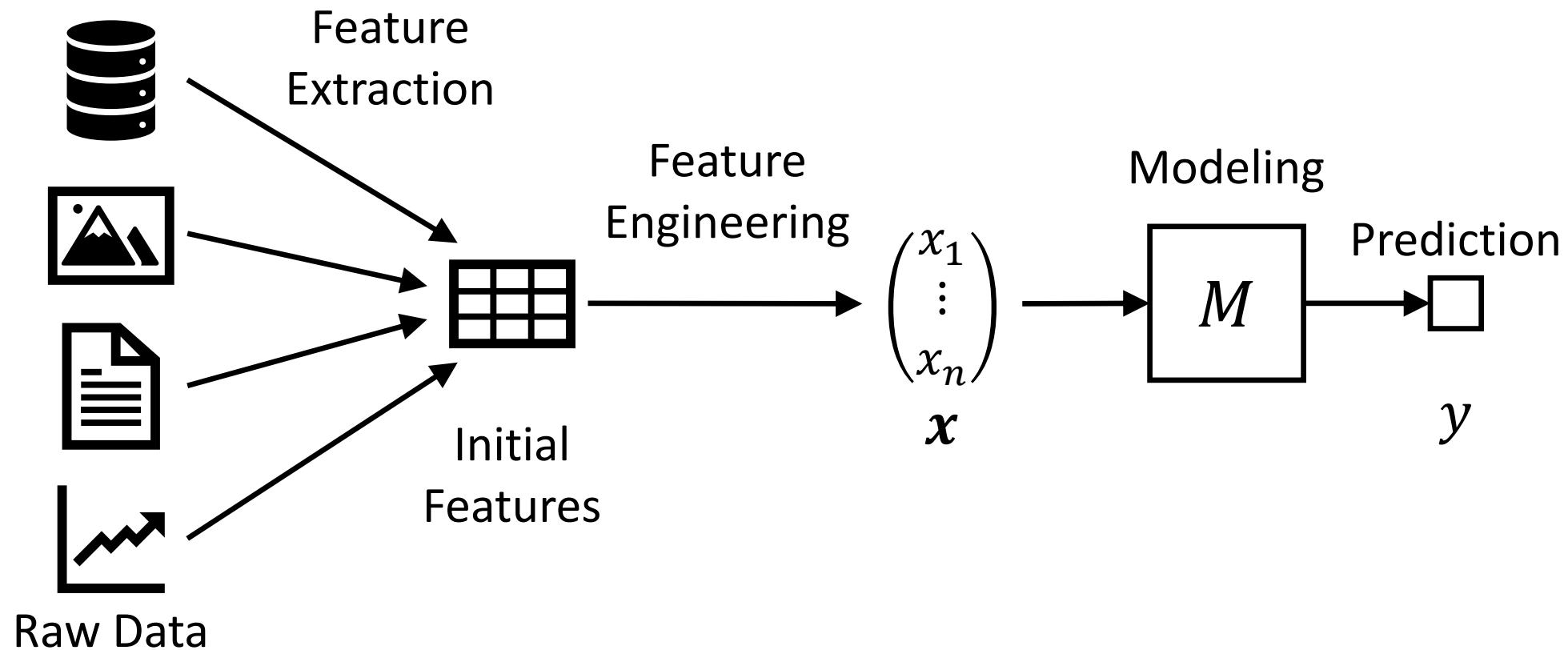


<https://9gag.com/gag/ax9Roon>

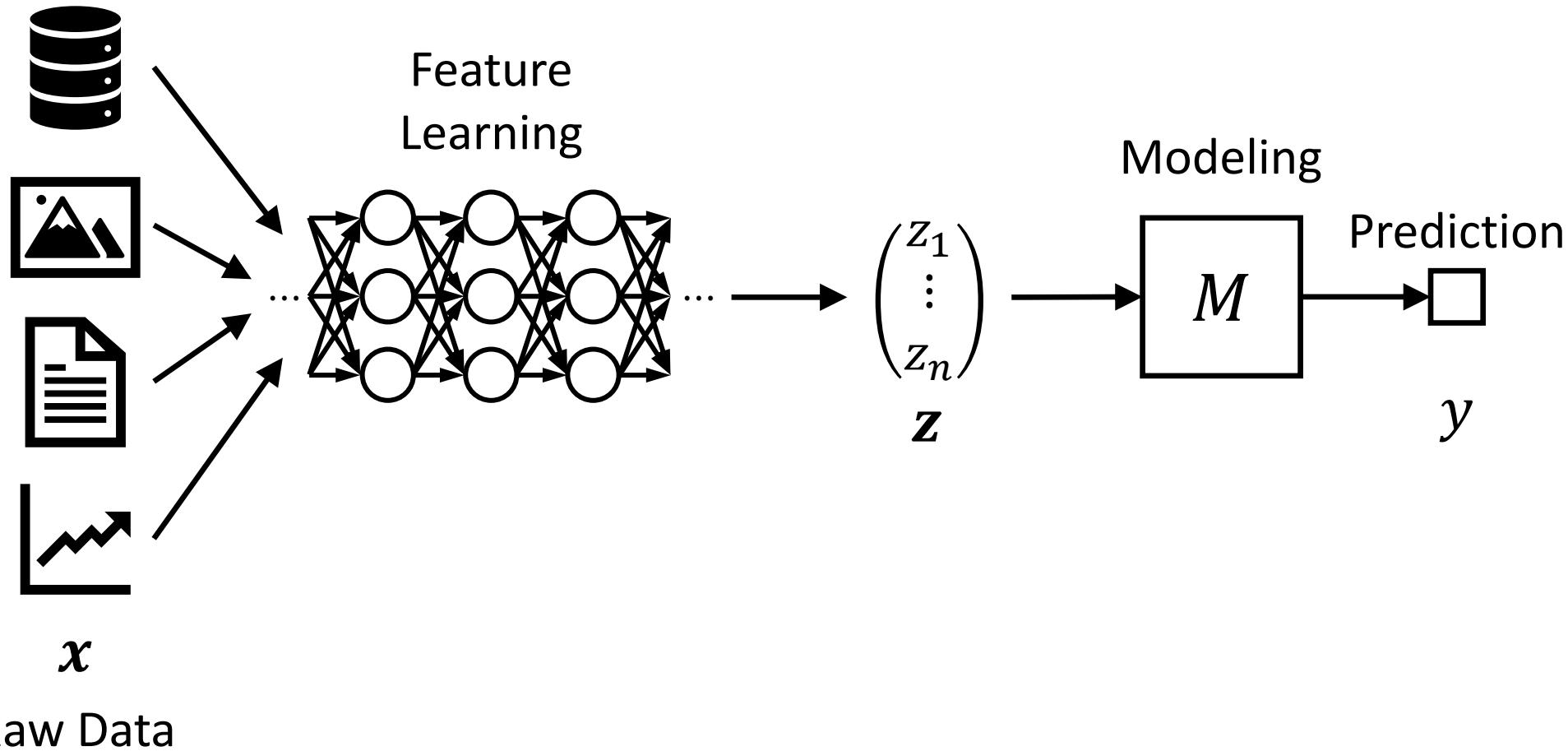


Andrew Ng <https://youtu.be/LcfLo7YP8O4>

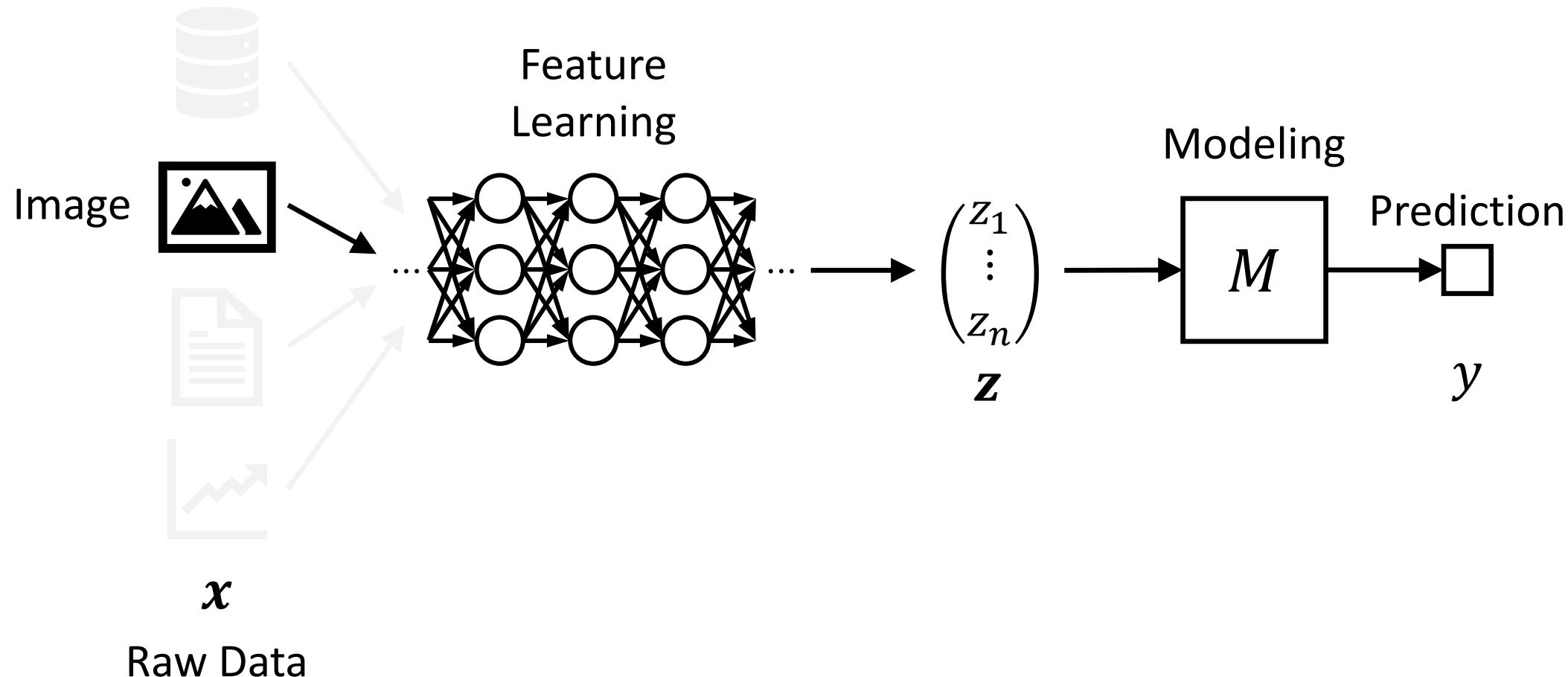
Feature Extraction/Engineering → Modeling



From Manual Feature Engineering To Automatic Feature Learning



From Manual Feature Engineering To Automatic Feature Learning



Convolutional Neural Networks (CNN)



Department of Computer Science
School of Computing



SCHOOL OF COMPUTING

Applications of CNN

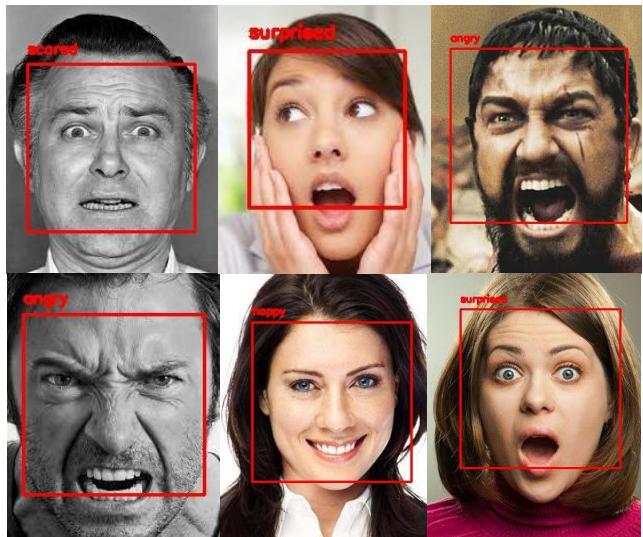
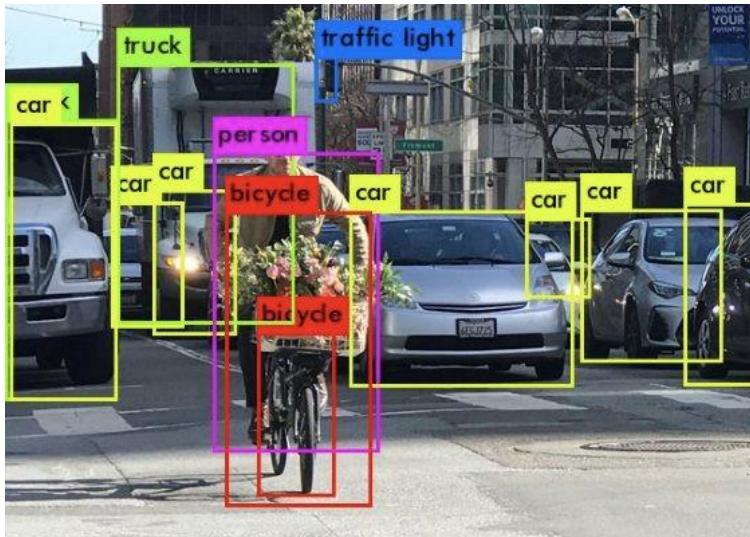


Image Classification
e.g., face emotions



Object Detection
e.g., self-driving cars

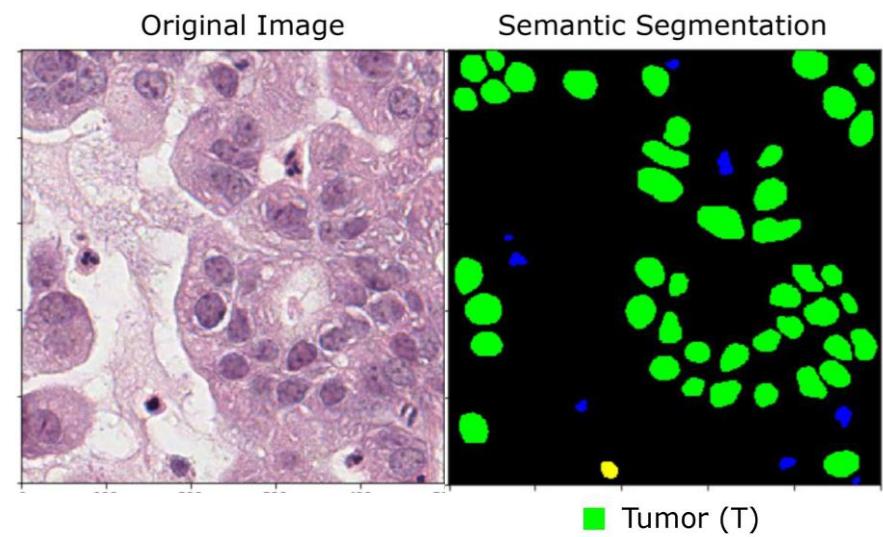


Image Segmentation
e.g., cancer cell detection

Image credit:

<https://monica-dommaraju.medium.com/analysis-of-deep-learning-based-object-detection-f14d5138148>

[https://ajp.amjpathol.org/article/S0002-9440\(18\)31121-0/fulltext](https://ajp.amjpathol.org/article/S0002-9440(18)31121-0/fulltext)

<https://appliedmachinelearning.blog/2018/11/28/demonstration-of-facial-emotion-recognition-on-real-time-video-using-cnn-python-keras/>





<https://foodai.org>

Try out our demo below or visit our developer portal for our API services.

To try our demo, you can **click** the upload icon to choose the image, or **copy and paste** the image or **drag and drop** the image from desktop or internet to the upload area.



Chicken rice

Boiled kampung chicken

Chicken porridge

Fish porridge

Fried fish porridge

Register for FoodAI API Free Trial

Backend Models

Food & Drink Recognition

It is a **dairy** type beverage called **strawberry milkshake** served in a transparent glass cup without logo, which is high in sugar and does not have alcohol in it.

It is a **beer** type beverage called **ale** served in a semi-transparent glass bottle with logo, which has no sugar but has alcohol in it.

The diagram illustrates a residual block-based neural network architecture. It shows two parallel paths: Model A and Model B. Both paths use 'Globally Shared Residual Blocks'. Model A consists of three residual blocks followed by a 'pool /2' layer. Model B consists of four residual blocks followed by a 'pool /2' layer. The final output is produced via 'avg pool'.



Frontend UI/UX

Family Food Logging

TableChat

Dinner time in office

Husband, F2

Adult daughter, F2

Husband, F2

Lychee

Challenge #3: Today, tell a family member how you care about their health

Mother, F2

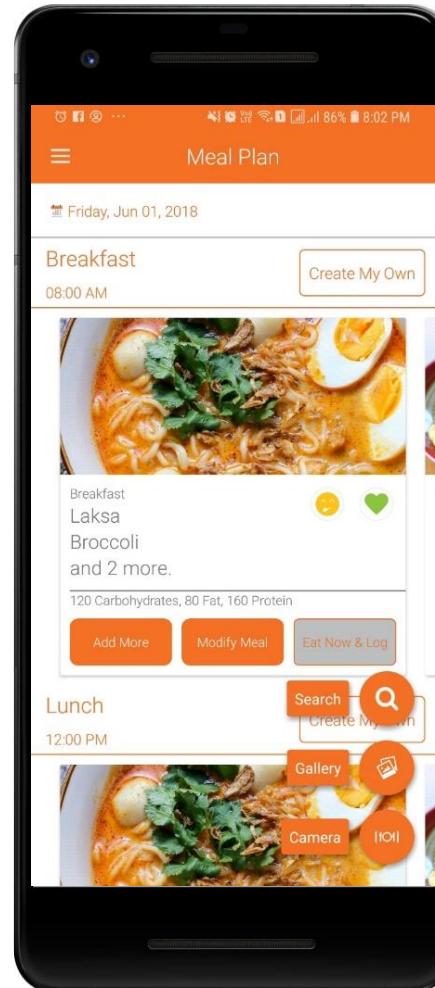
Husband, F3

Baby this is to tell you to drink more water and eat on time so your health can be improved ok? 😊

Food Recommendation

RecGAN

Generator $\mathbf{x}_t^u = \text{RELU}(\mathbf{W}_{xh}^u \mathbf{h}_{t-1}^u + \mathbf{W}_{xk}^u \mathbf{y}_t^u)$ $\mathbf{r}_t^u = \sigma(\mathbf{W}_{rh}^u \mathbf{h}_{t-1}^u + \mathbf{W}_{rk}^u \mathbf{y}_t^u)$ $\mathbf{m}_t^u = \tanh(\mathbf{W}_h(\mathbf{r}_t^u \cdot \mathbf{h}_{t-1}^u) + \mathbf{W}_{input} \mathbf{y}_t^u)$ $\mathbf{h}_t^u = (1 - \mathbf{x}_t^u) \cdot \mathbf{h}_{t-1}^u + \mathbf{x}_t^u \cdot \mathbf{m}_t^u$	Discriminator $\hat{\mathbf{x}}_t^u = \text{RELU}(\mathbf{V}_{xh}^u \hat{\mathbf{h}}_{t-1}^u + \mathbf{V}_{xk}^u \mathbf{y}_t^u)$ $\hat{\mathbf{r}}_t^u = \sigma(\mathbf{V}_{rh}^u \hat{\mathbf{h}}_{t-1}^u + \mathbf{V}_{rk}^u \mathbf{y}_t^u)$ $\hat{\mathbf{m}}_t^u = \tanh(\mathbf{V}_h(\hat{\mathbf{r}}_t^u \cdot \hat{\mathbf{h}}_{t-1}^u) + \mathbf{V}_{input} \mathbf{y}_t^u)$ $\hat{\mathbf{h}}_t^u = (1 - \hat{\mathbf{x}}_t^u) \cdot \hat{\mathbf{h}}_{t-1}^u + \hat{\mathbf{x}}_t^u \cdot \hat{\mathbf{m}}_t^u$
--	--



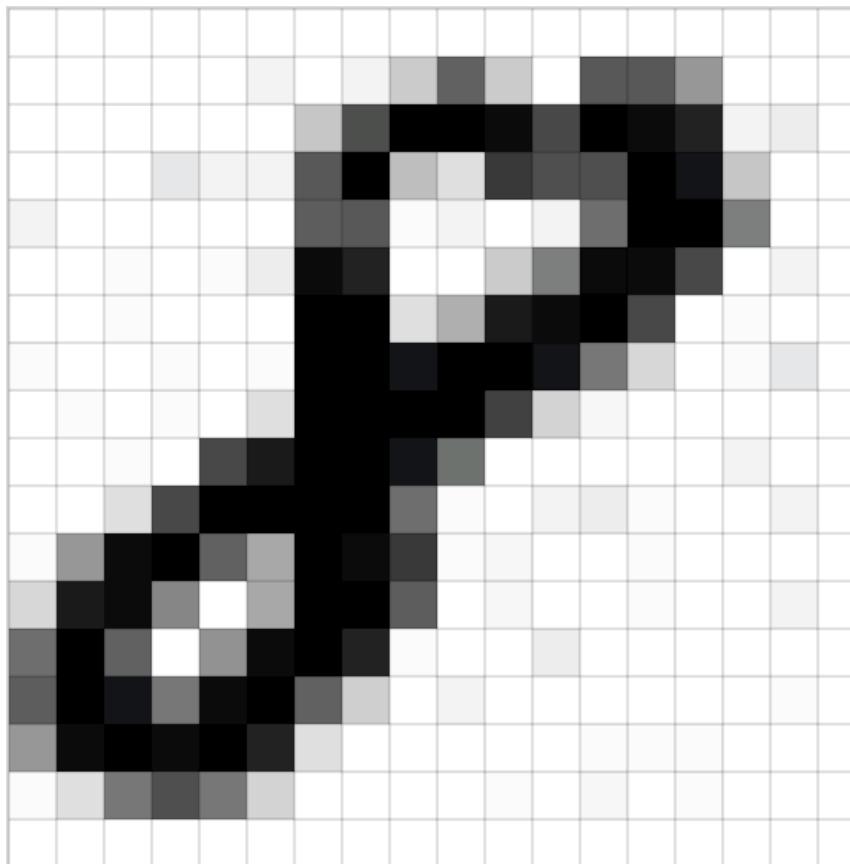
Multi-Attribute Sorting

(a) Imma Sort

(b) One-attribute Sort by Price

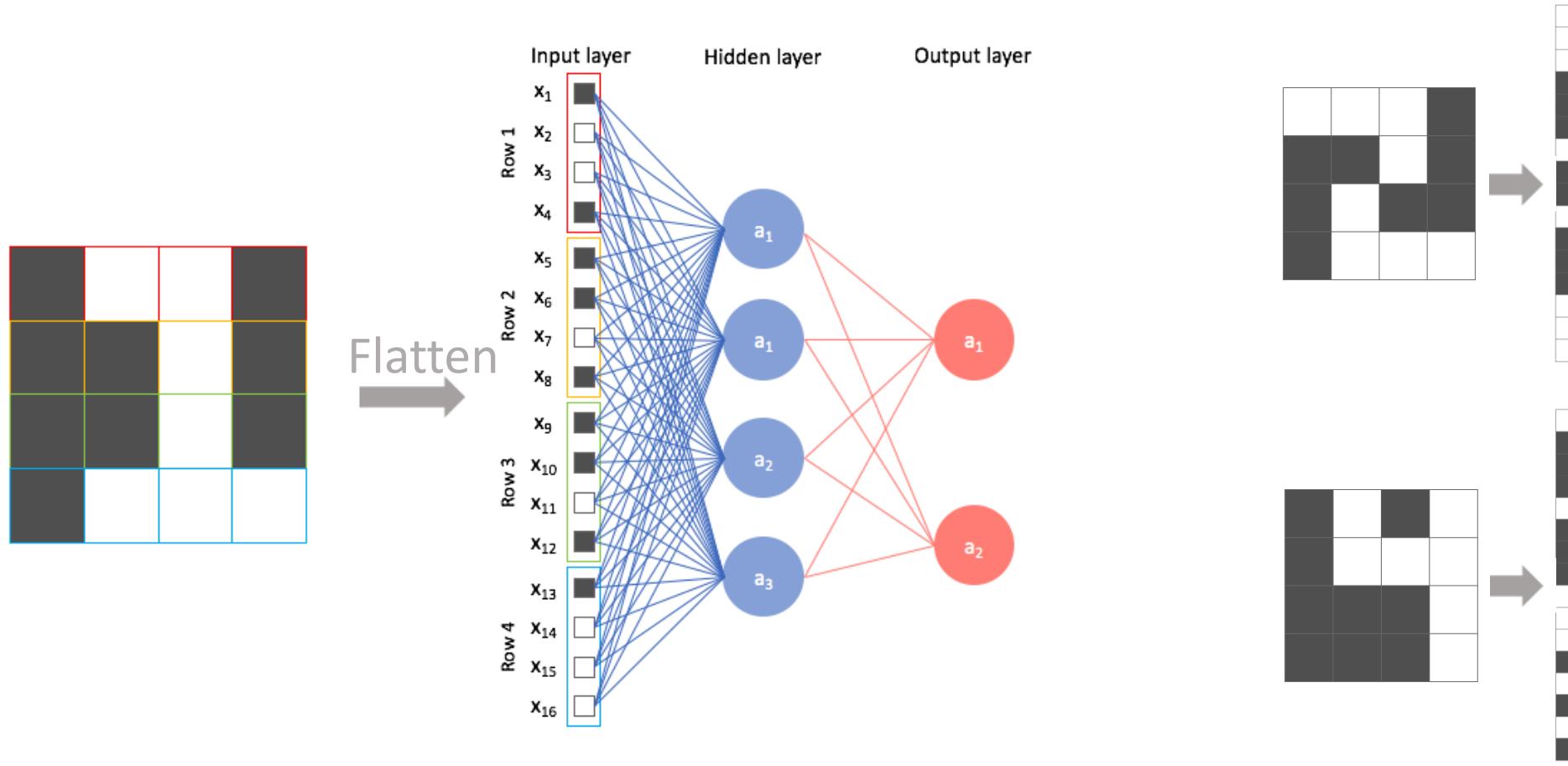
- Lyu, Y., Gao, F., Wu, I. S., & Lim, B. Y. 2020. Imma Sort by Two or More Attributes With Interpretable Monotonic Multi-Attribute Sorting. TVCG.
- Park, H., Bharadhwaj, H., and Lim, B. Y. 2019. Hierarchical Multi-Task Learning for Healthy Drink Classification. IJCNN.
- Bharadhwaj, H., Park, H., Lim, B. Y.. 2018. RecGAN: Recurrent Generative Adversarial Networks for Recommendation Systems. RecSys '18.
- Lukoff, K., Li, T., Zhuang, Y., & Lim, B. Y. 2018. TableChat: Mobile Food Journaling to Facilitate Family Support for Healthy Eating. CSCW '18.

Images as 2D matrices



[Image credit](#)

Image Feature Extraction with Fully Connected Neural Networks (Multi-Layer Perceptron)

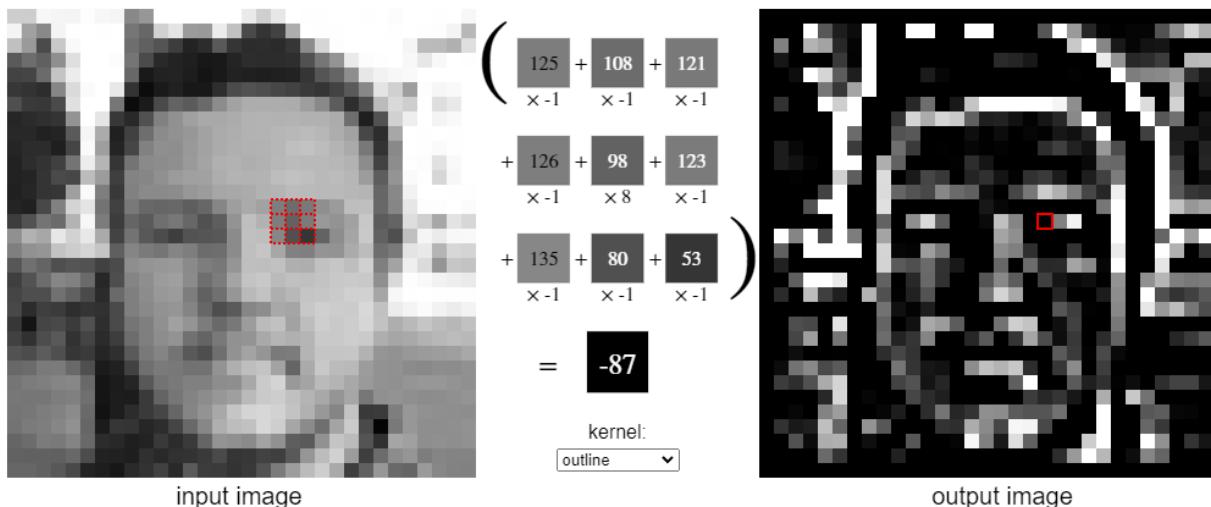


Reduce parameters for images: Exploit Spatial Relations with Convolutions

Let's walk through applying the following 3x3 **outline** kernel to the image of a face from above.

$$\text{outline} \quad \begin{pmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{pmatrix}$$

Below, for each 3x3 block of pixels in the image on the left, we multiply each pixel by the corresponding entry of the kernel and then take the sum. That sum becomes a new pixel in the image on the right. Hover over a pixel on either image to see how its value is computed.



Manually finding
good filters is
tedious

Further study:
<https://setosa.io/ev/image-kernels/>

High-level Feature Detection



Eyes, Nose, Mouth
Facial Hair



Wheels, Headlights,
Bonnet/Hood



Fish, Rice,
Vegetables

How to automatically learn these features?

Feature Detectors: Intuition of Neuron Kernels in Layers

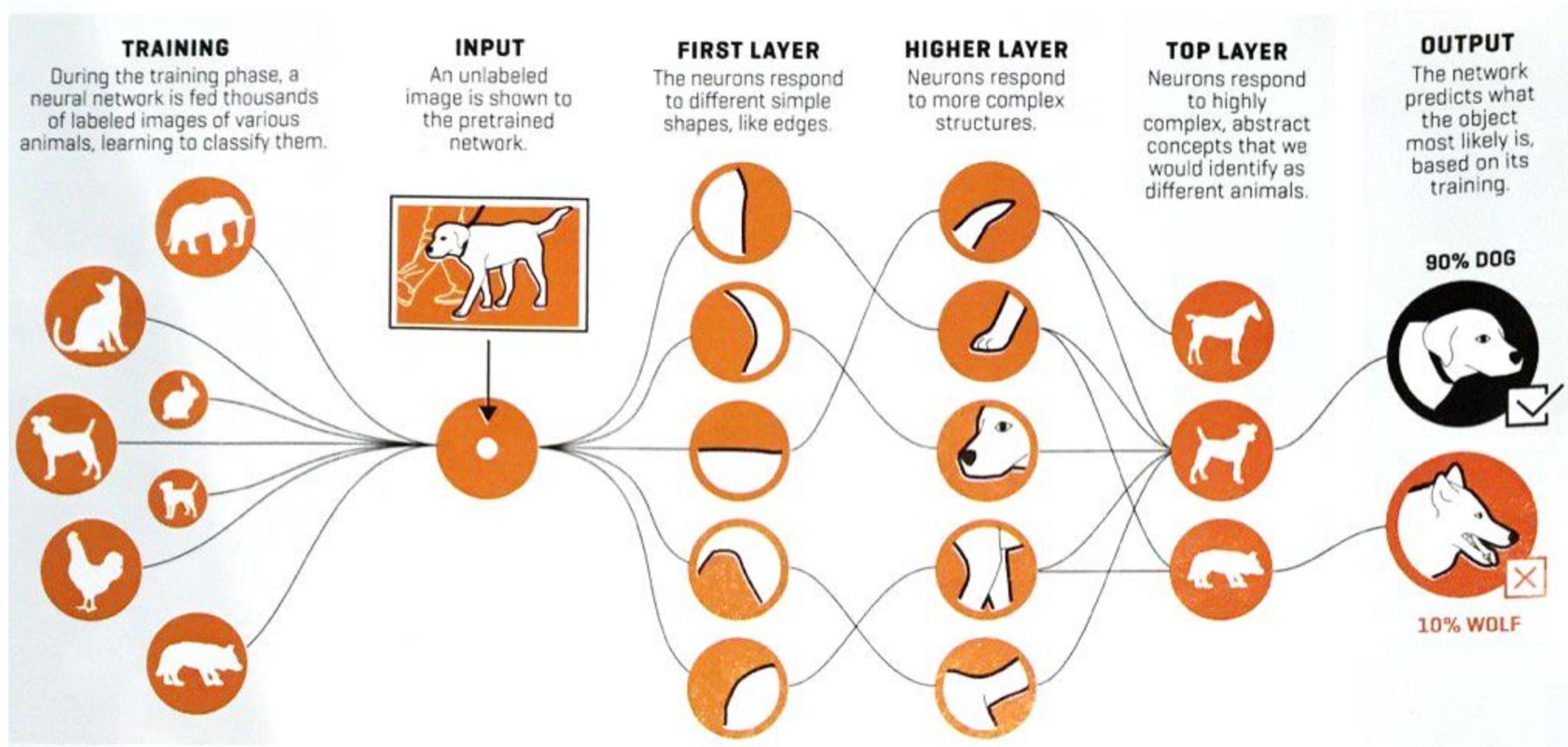


Image credit: <https://fortune.com/longform/ai-artificial-intelligence-deep-machine-learning/>

Analogy: activations of different filters learned by CNNs is like seeing the image through different lens filters



Image credit: <https://www.amazon.com/Godefa-Samsung-Andriod-Smartphone-Universal/dp/B07RQRLQYH>
<https://www.yankodesign.com/2020/02/17/this-retro-inspired-camera-records-dreamy-looking-gifs-that-replicate-vintage-8mm-film/>

Convolutions: Kernel Size, Stride, Padding

$$W = \begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$$

$$x = \begin{pmatrix} 9 & 9 & 3 & 3 & 4 \\ 9 & 3 & 3 & 4 & 5 \\ 9 & 3 & 3 & 5 & 5 \\ 9 & 3 & 3 & 4 & 5 \\ 9 & 9 & 3 & 3 & 4 \end{pmatrix}$$

$$W * x = \begin{pmatrix} -6 - 6 - 6 & -6 + 1 + 2 & 1 + 2 + 2 \\ -6 - 6 - 6 & 1 + 2 + 1 & 2 + 2 + 2 \\ -6 - 6 - 6 & 2 + 1 - 6 & 2 + 2 + 1 \end{pmatrix}$$

Stride $s \neq 1$

$$W = \begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$$

$$x = \begin{pmatrix} 9 & 9 & 3 & 3 & 4 \\ 9 & 3 & 3 & 4 & 5 \\ 9 & 3 & 3 & 5 & 5 \\ 9 & 3 & 3 & 4 & 5 \\ 9 & 9 & 3 & 3 & 4 \end{pmatrix}$$

$$W * x = \begin{pmatrix} -6 - 6 - 6 & 1 + 2 + 2 \\ -6 - 6 - 6 & 2 + 2 + 1 \end{pmatrix}$$

Padding $p \neq 0$

$$W = \begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$$

$$x = \begin{pmatrix} 3 & 3 & 4 \\ 3 & 3 & 5 \\ 3 & 3 & 4 \end{pmatrix}$$

$$W * x = \begin{pmatrix} -6 - 6 - 6 & 1 + 2 + 2 \\ -6 - 6 - 6 & 2 + 2 + 1 \end{pmatrix}$$

Convolutional Layer

What's the **differences** between the left and right expressions?

Convolution of each **Feature Kernel** $\mathbf{W}^{[l]}$ produces a **Feature Map** = $\mathbf{W}^{[l]} * \mathbf{A}^{[l-1]}$

$\mathbf{a}^{[l]} = g^{[l]} \left((\mathbf{W}^{[l]})^T \mathbf{a}^{[l-1]} \right)$ vs. $\mathbf{A}^{[l]} = g^{[l]}(\mathbf{W}^{[l]} * \mathbf{A}^{[l-1]})$

Activation function can be the same

Weights = Kernels
(3D matrix)

2D Convolution
(height \times width)

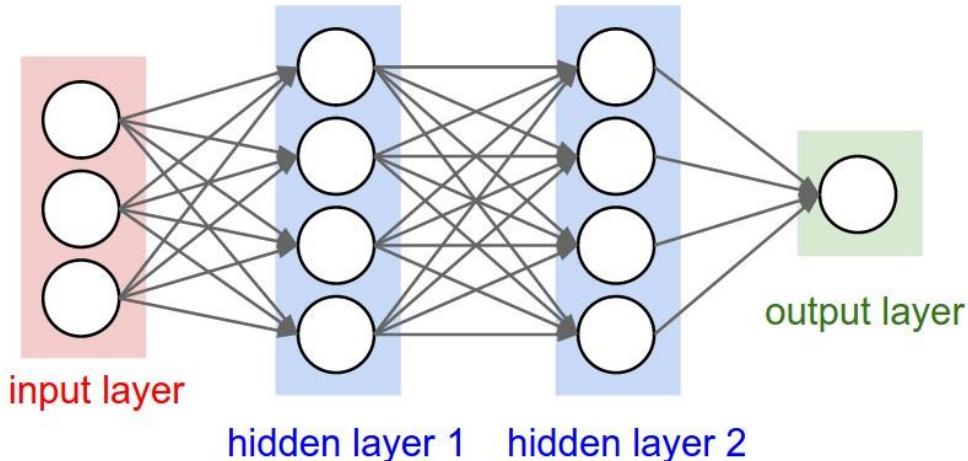
Activation is a 3D Matrix!
Concatenation of Feature Maps

The diagram illustrates the difference between two ways of representing a convolutional layer. On the left, the expression $\mathbf{a}^{[l]} = g^{[l]} \left((\mathbf{W}^{[l]})^T \mathbf{a}^{[l-1]} \right)$ shows a linear transformation (matrix multiplication of $\mathbf{W}^{[l]}$ and $\mathbf{a}^{[l-1]}$) followed by an activation function $g^{[l]}$. On the right, the expression $\mathbf{A}^{[l]} = g^{[l]}(\mathbf{W}^{[l]} * \mathbf{A}^{[l-1]})$ shows a convolution operation ($\mathbf{W}^{[l]}$ convolved with $\mathbf{A}^{[l-1]}$) followed by an activation function $g^{[l]}$. A curved arrow from the left equation points to the right equation, indicating that the convolution step in the right equation corresponds to the linear transformation in the left equation. A callout box at the top left states: "Convolution of each Feature Kernel $\mathbf{W}^{[l]}$ produces a Feature Map = $\mathbf{W}^{[l]} * \mathbf{A}^{[l-1]}$ ". To the right of the equations, a vertical line labeled "Weights = Kernels (3D matrix)" connects to the $\mathbf{W}^{[l]}$ term in both equations. A curved arrow labeled "2D Convolution (height \times width)" connects the $\mathbf{W}^{[l]}$ term in the right equation to the $\mathbf{W}^{[l]}$ term in the left equation. Below the equations, the text "Activation function can be the same" is centered. To the right of the right equation, the text "Activation is a 3D Matrix! Concatenation of Feature Maps" is shown.

Convolutional Layers

Fully Connected Layers

- Each layer has multiple **neurons** ○
- Neuron output: **0D scalar activation**
- Neuron input: **1D vector of activations**
 - Each *element* is a different neuron
- Each layer is a **1D vector**



Remember: each kernel is like a different lens filter



Convolutional Layers

- Each layer has multiple **ernels** □
- Kernel output: **2D matrix feature map**
- Kernel input: **3D matrix of feature maps**
 - Each *depth position* is a different kernel
 - Analogy: filters are “stacked” together
- Each layer is a **3D matrix**

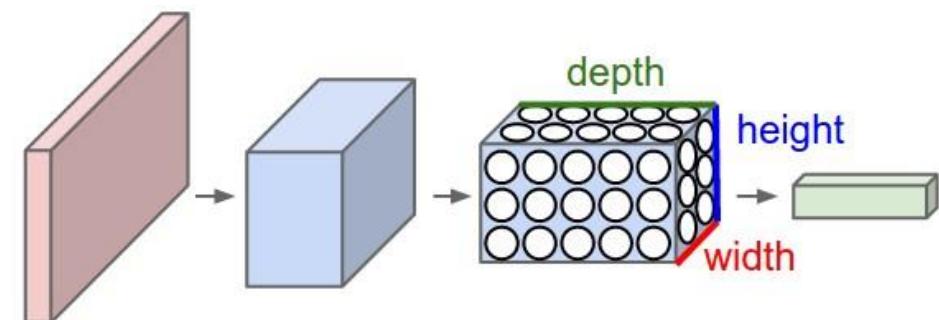


Image credit: <https://cs231n.github.io/convolutional-networks/>

Convolutional Layer: Feature Kernels & Feature Maps

$$X^{[0]} \rightarrow g^{[1]}(W^{[1]} * X^{[0]}) = A^{[1]} \xrightarrow{\text{Pooling}} g^{[2]}(W^{[2]} * X^{[1]}) = A^{[2]}$$

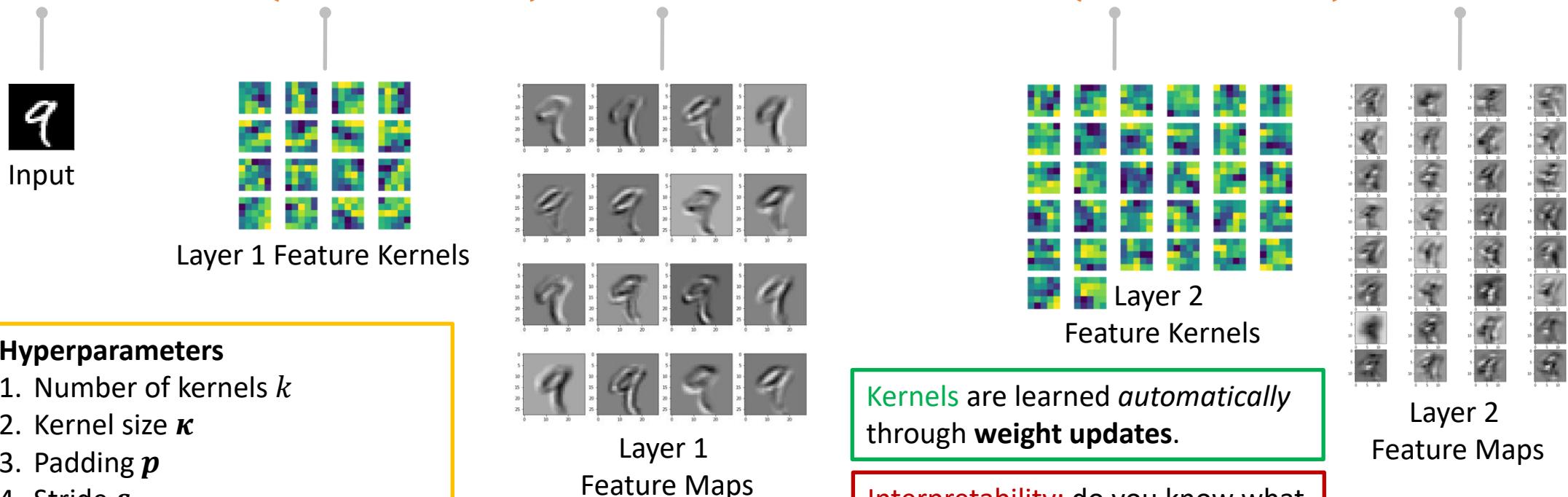


Image credit: <https://medium.com/dataseries/visualizing-the-feature-maps-and-filters-by-convolutional-neural-networks-e1462340518e>

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

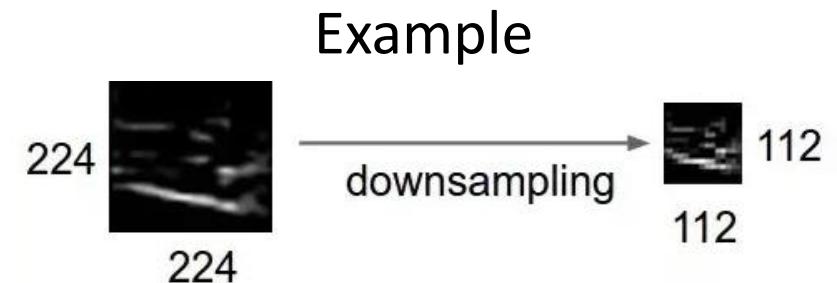
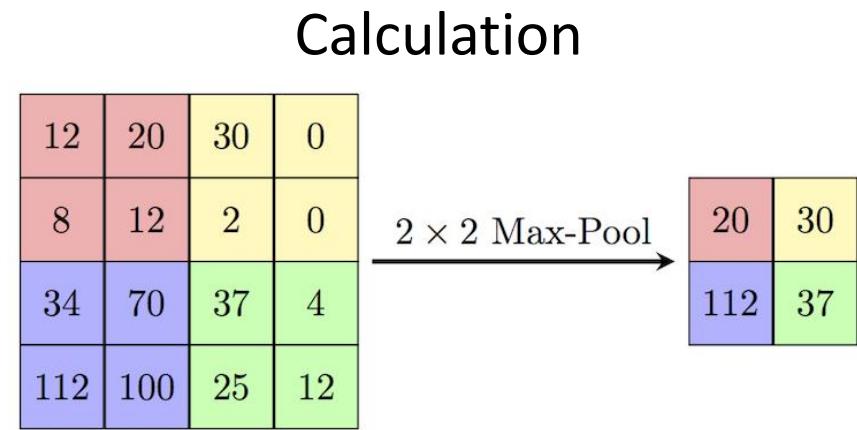
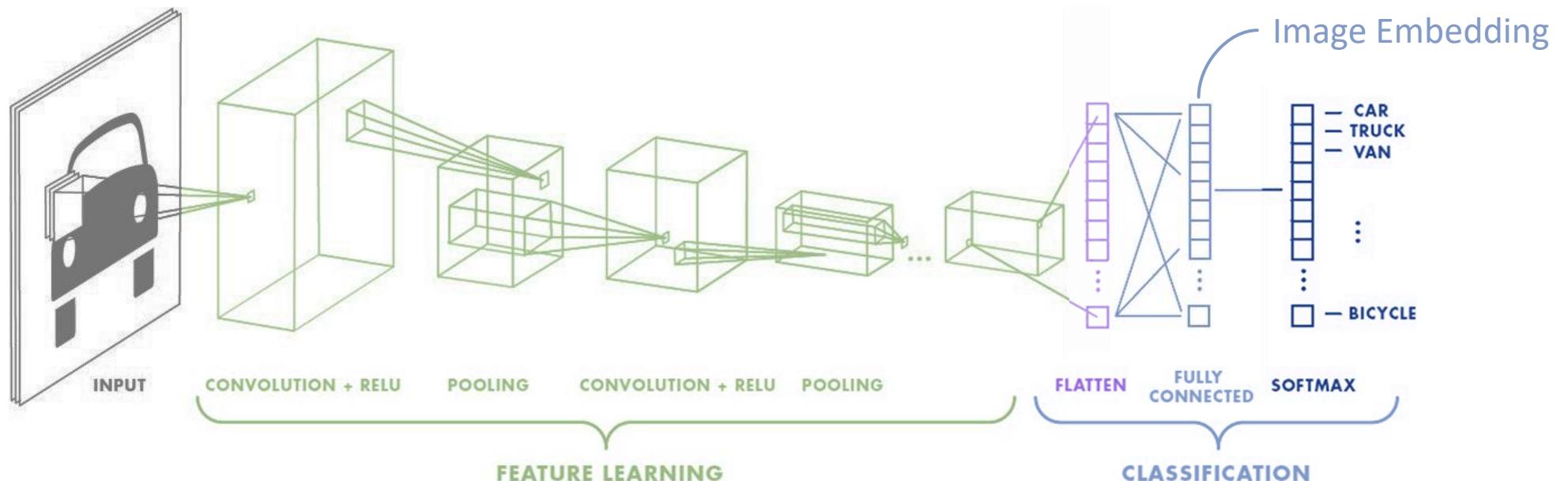


Image credit: https://computersciencewiki.org/index.php/Max-pooling_/_Pooling

Convolutional Neural Network



Key concepts

① Learn Spatial Feature

- Series of multiple convolution + pooling layers
- Progressively learn more diverse and higher-level features

② Flattening

- Convert to fixed-length 1D vector

③ Learn Nonlinear Features

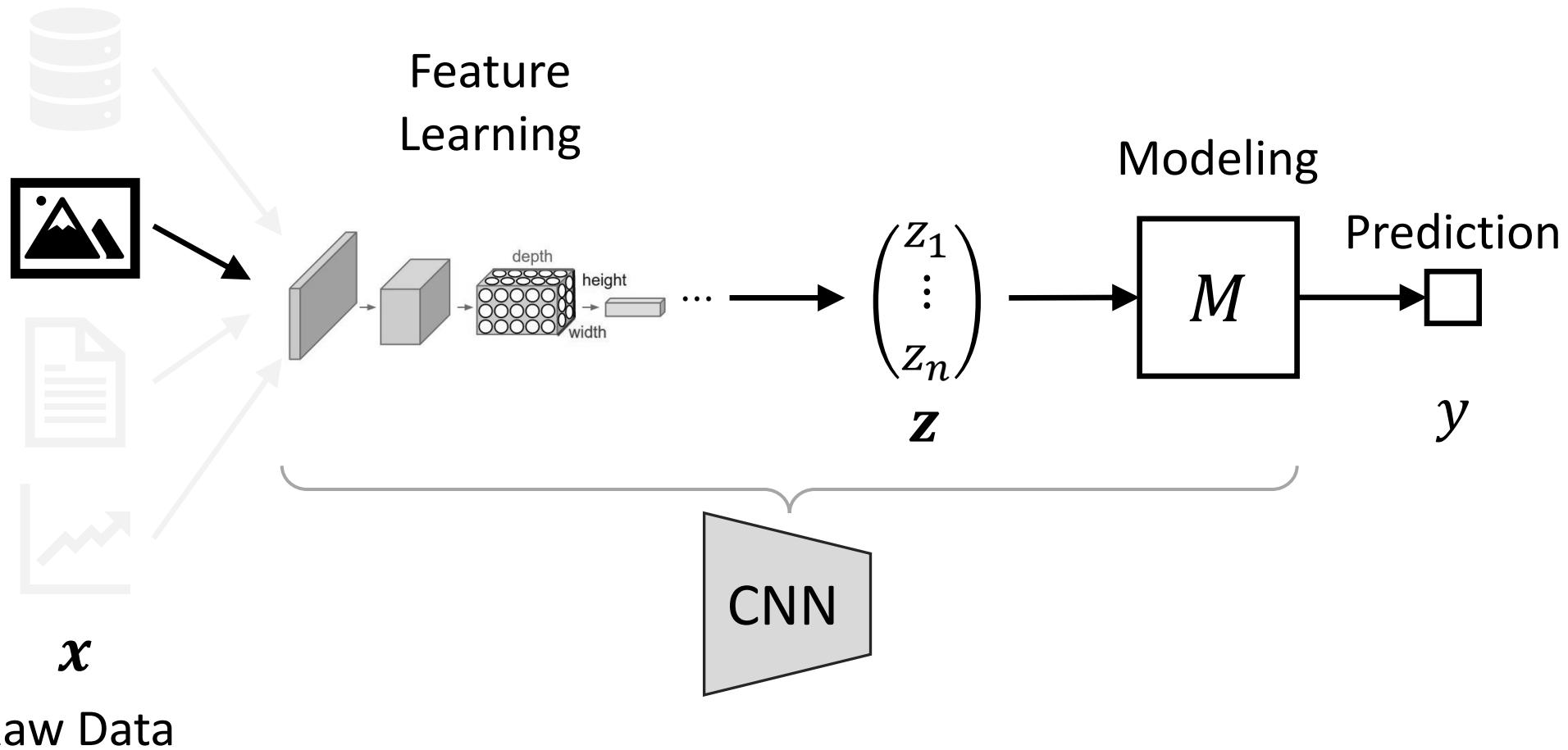
- With fully connected layers (regular neurons)
- Learns nonlinear relations with multiple layers

④ Classification

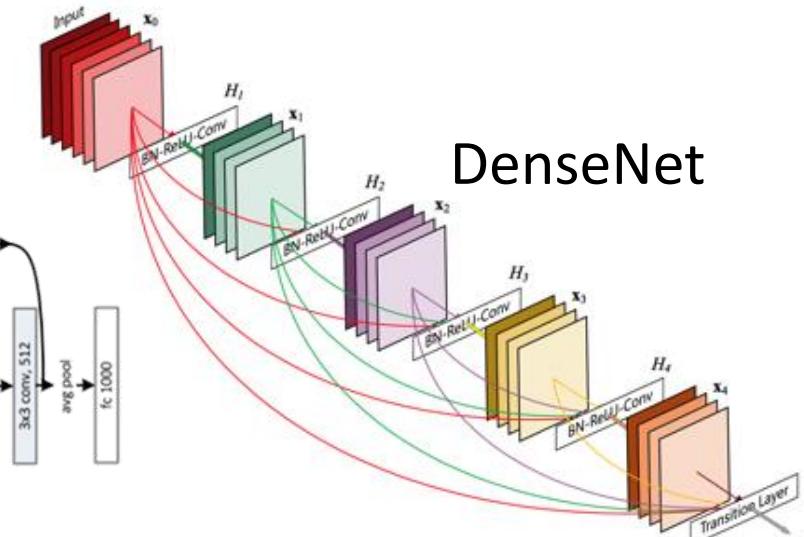
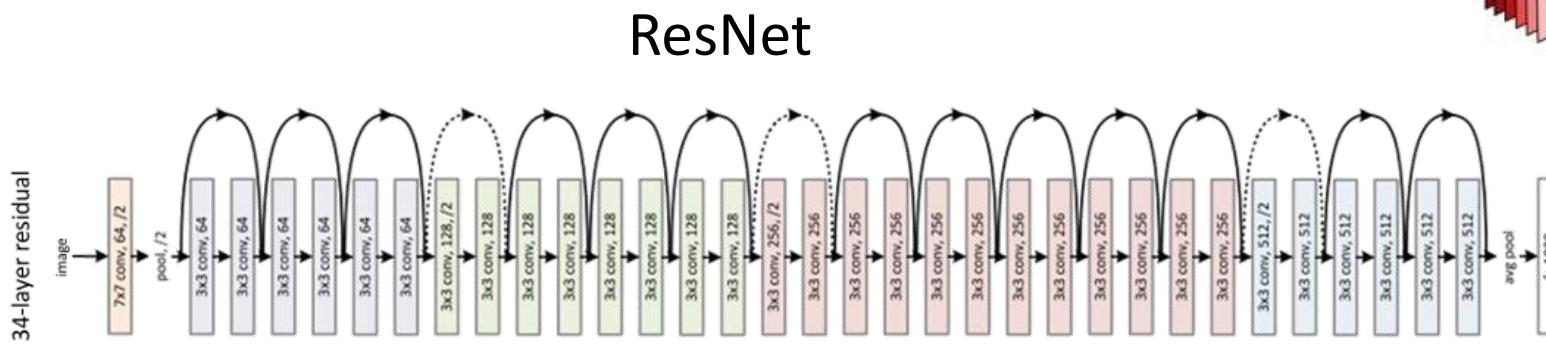
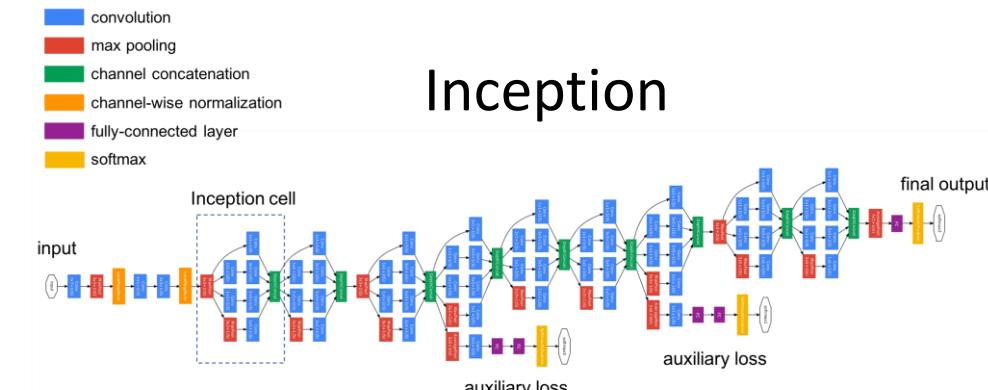
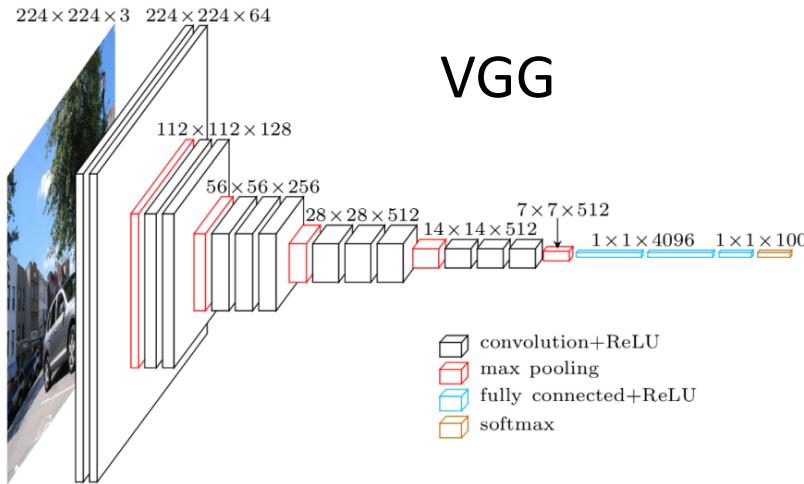
- Softmax := Multiclass Logistic Regression
- Feature input = image embedding vector (typically large vector)

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

From Manual Feature Engineering To Automatic Feature Learning



Other popular CNN architectures



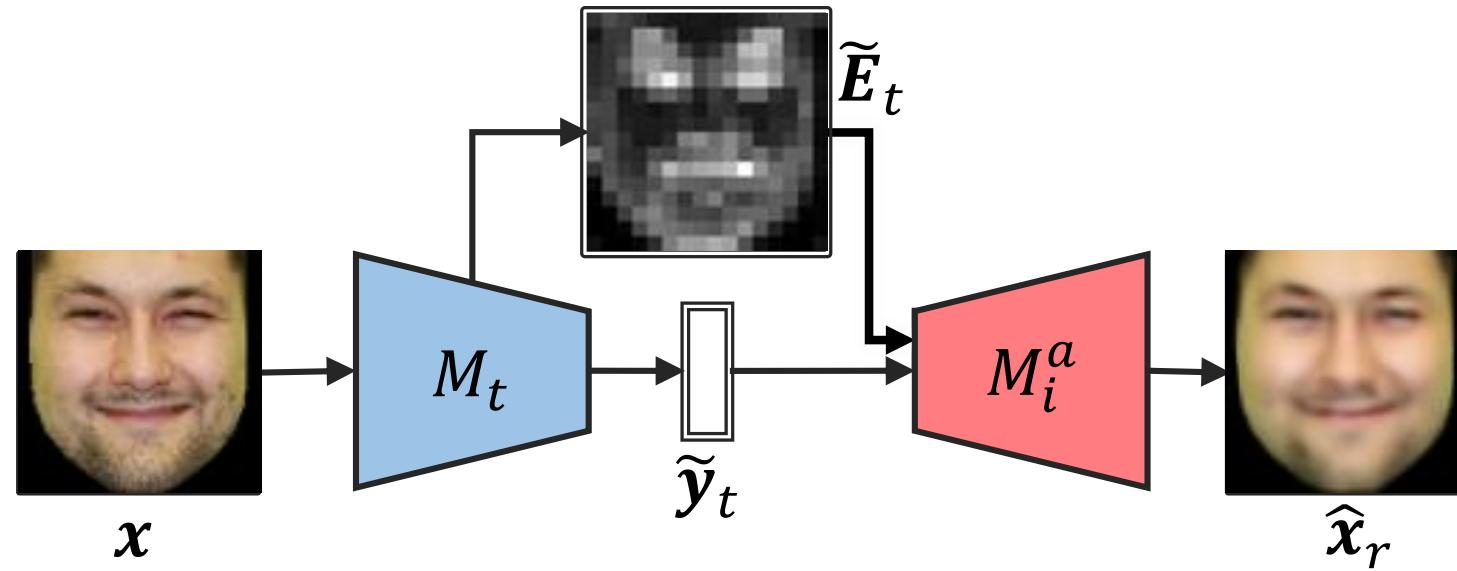
Further reading: <https://www.jeremyjordan.me/convnet-architectures/>



Questions!



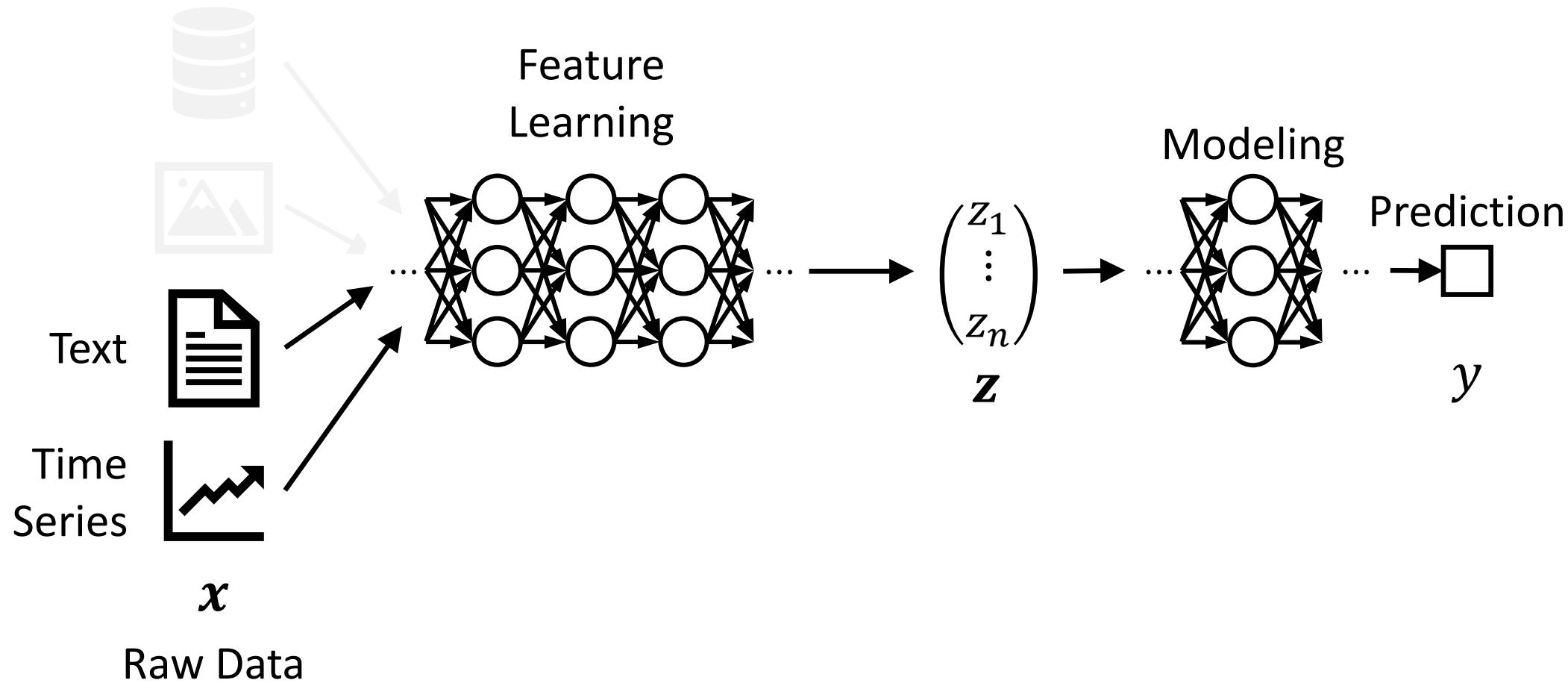
Model Inversion Attack: Predicting Images from Classification Vector



Model **inversion attacks** can reconstruct **private** face photos from prediction vectors only.

Model **explanations** can **worsen** model **inversion attacks**.

From Manual Feature Engineering To Automatic Feature Learning





Recurrent Neural Networks (RNN)



NUS
National University
of Singapore

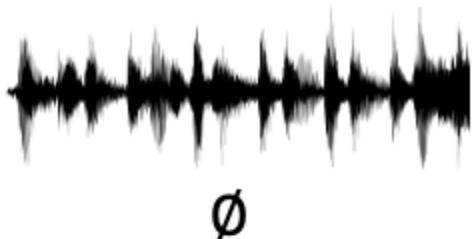
Department of Computer Science
School of Computing



SCHOOL OF COMPUTING

Applications of RNN

Speech recognition



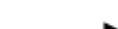
“The quick brown fox jumped
over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like
in this movie.”



DNA sequence analysis

AGCCCCCTGTGAGGAACCTAG



AG~~CCCCCTGTGAGGAACCTAG~~

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/>

 info  random  clear

Model: 



start drawing face.

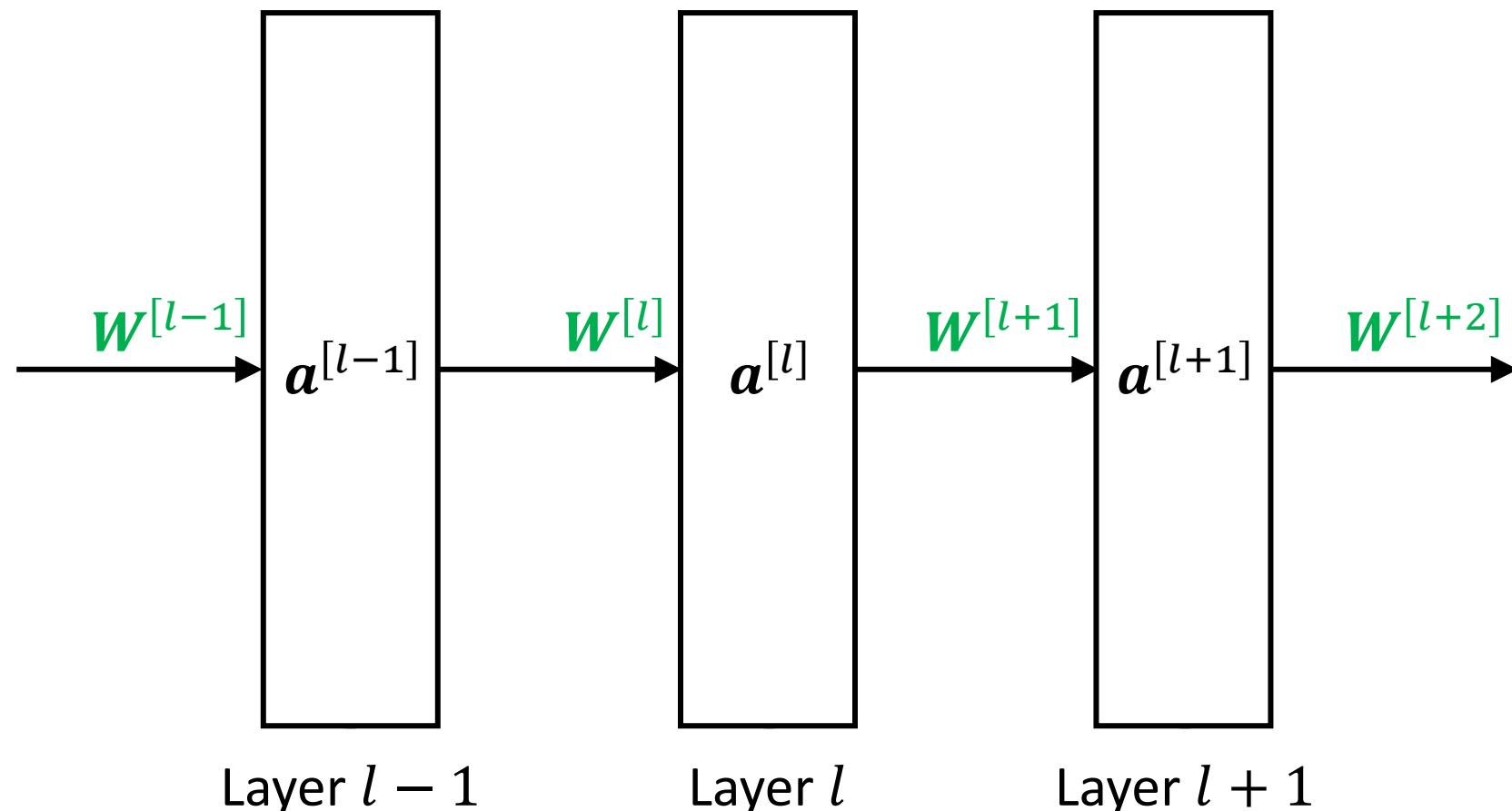
Try yourself:

https://magenta.tensorflow.org/assets/sketch_rnn_demo/index.html

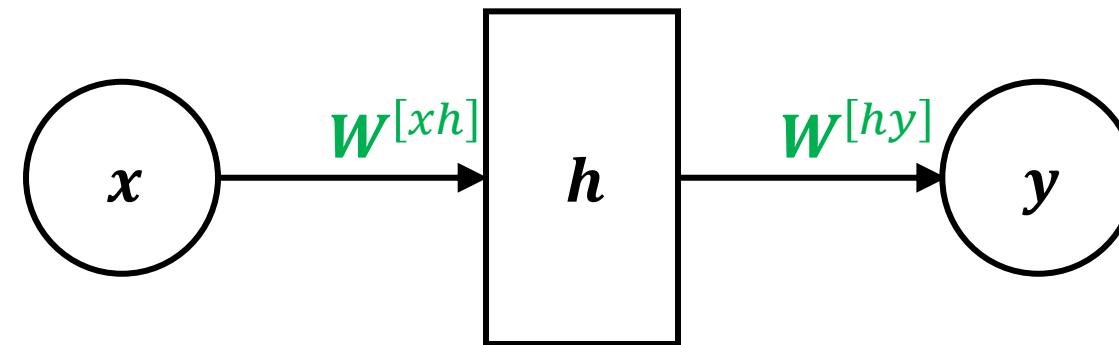
Play

- Visit:
https://magenta.tensorflow.org/assets/sketch_rnn_demo/index.html
- Draw something and see what the AI will continue drawing
- Take screenshot and post your results
 - No obscenities please (we will have to take disciplinary action)
- Emote
 - Up vote those you like
 - Down vote those with mistakes
- Discuss how you think the model predicted what to draw next

Feedforward Neural Network



Neural Network (simplified 1-hidden layer)

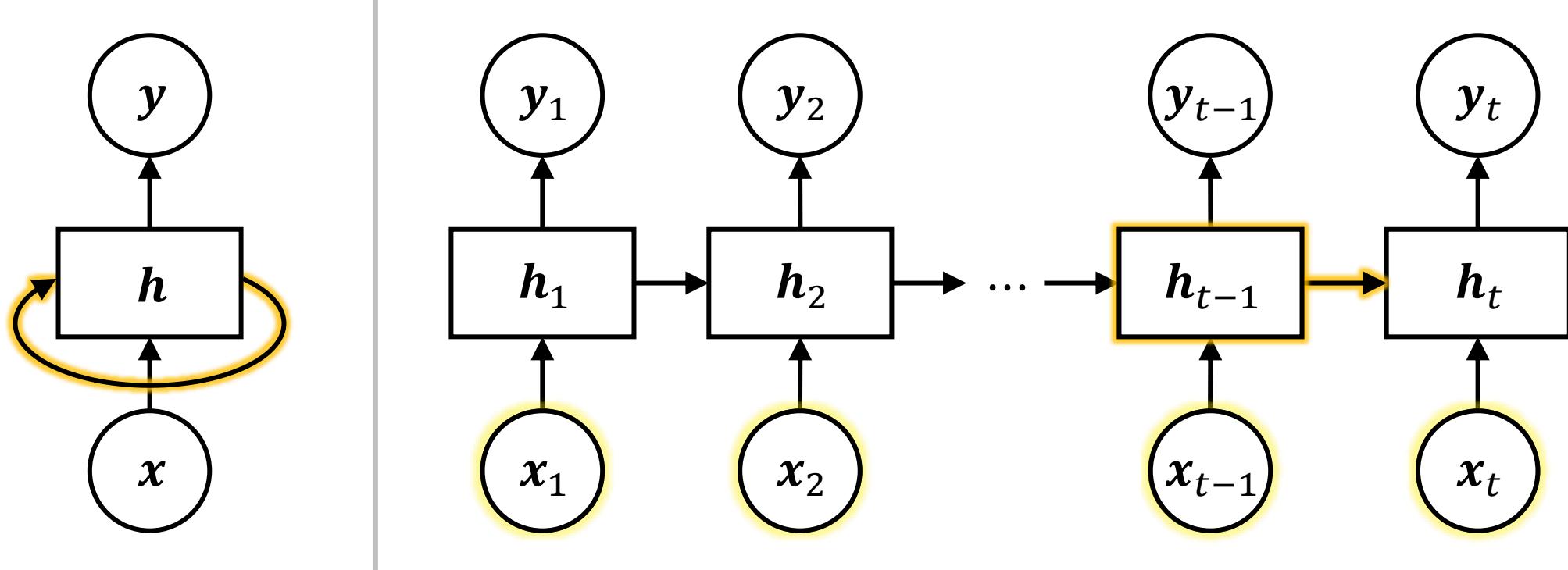


Input Layer

Hidden Layer

Output Layer

Neurons with Recurrence

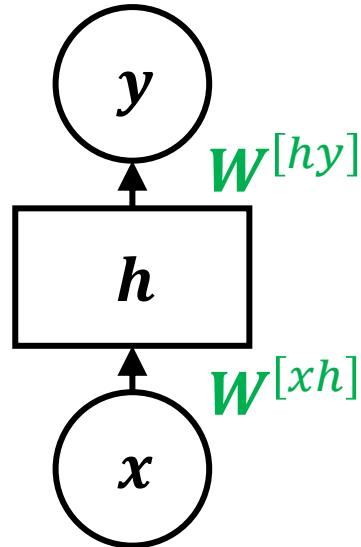


$$\hat{y} = g^{[y]}(g^{[h]}(x_t, \mathbf{h}_{t-1}))$$

Recurrent Neural Network (RNN)

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

RNN Weights

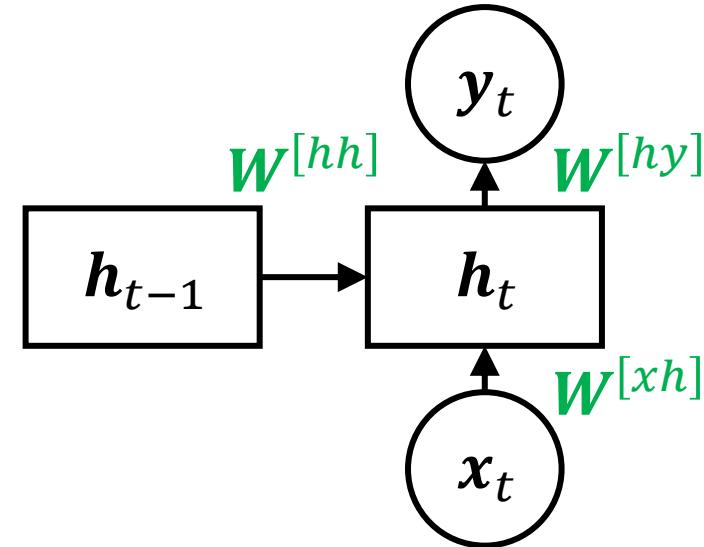


Feedforward Neural Network

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

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

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



Recurrent Neural Network

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

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

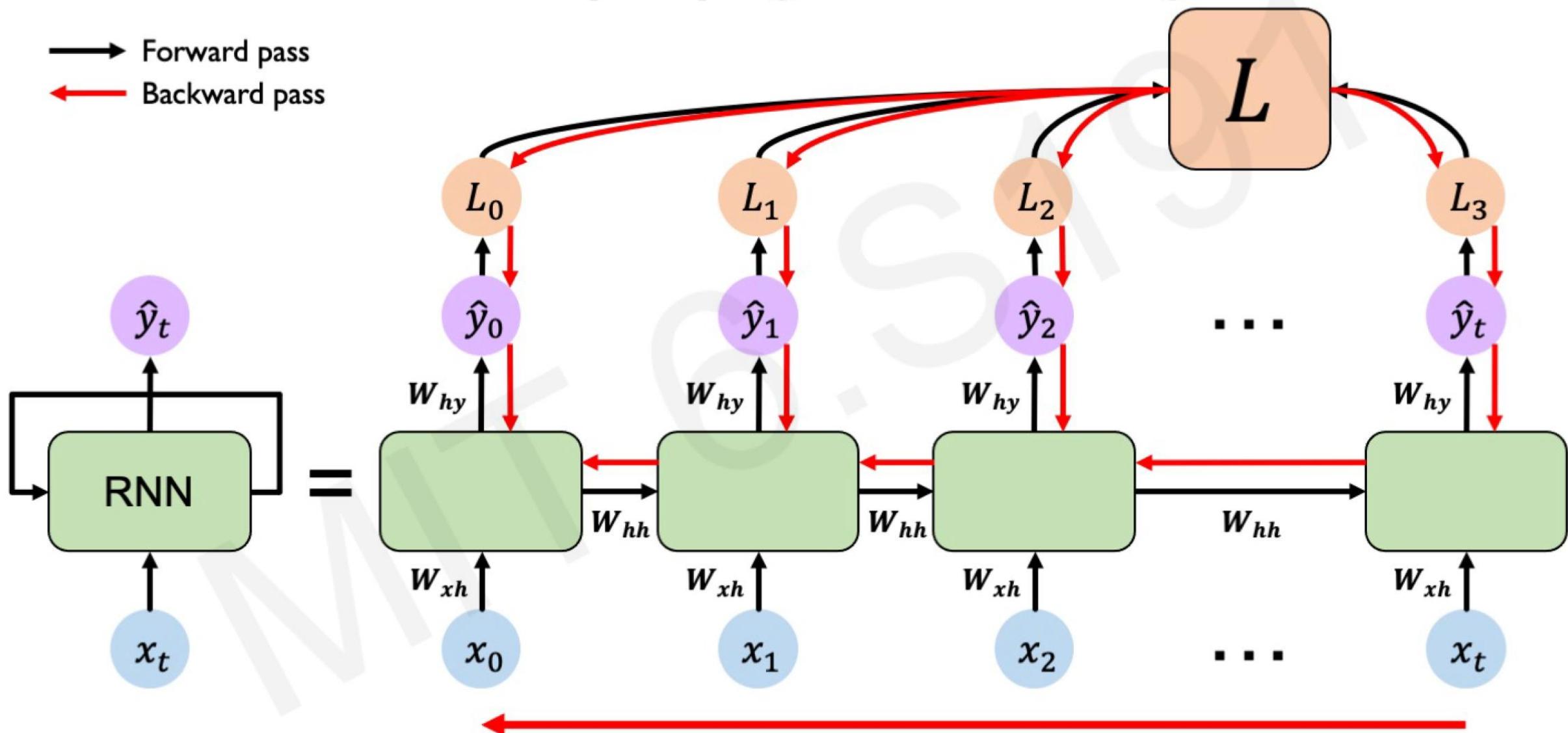
$$\mathbf{h}_t = g^{[hh]} \left((\mathbf{W}^{[xh]} \oplus \mathbf{W}^{[hh]})^T (\mathbf{x}_t \oplus \mathbf{h}_{t-1}) \right)$$

HIT

2007



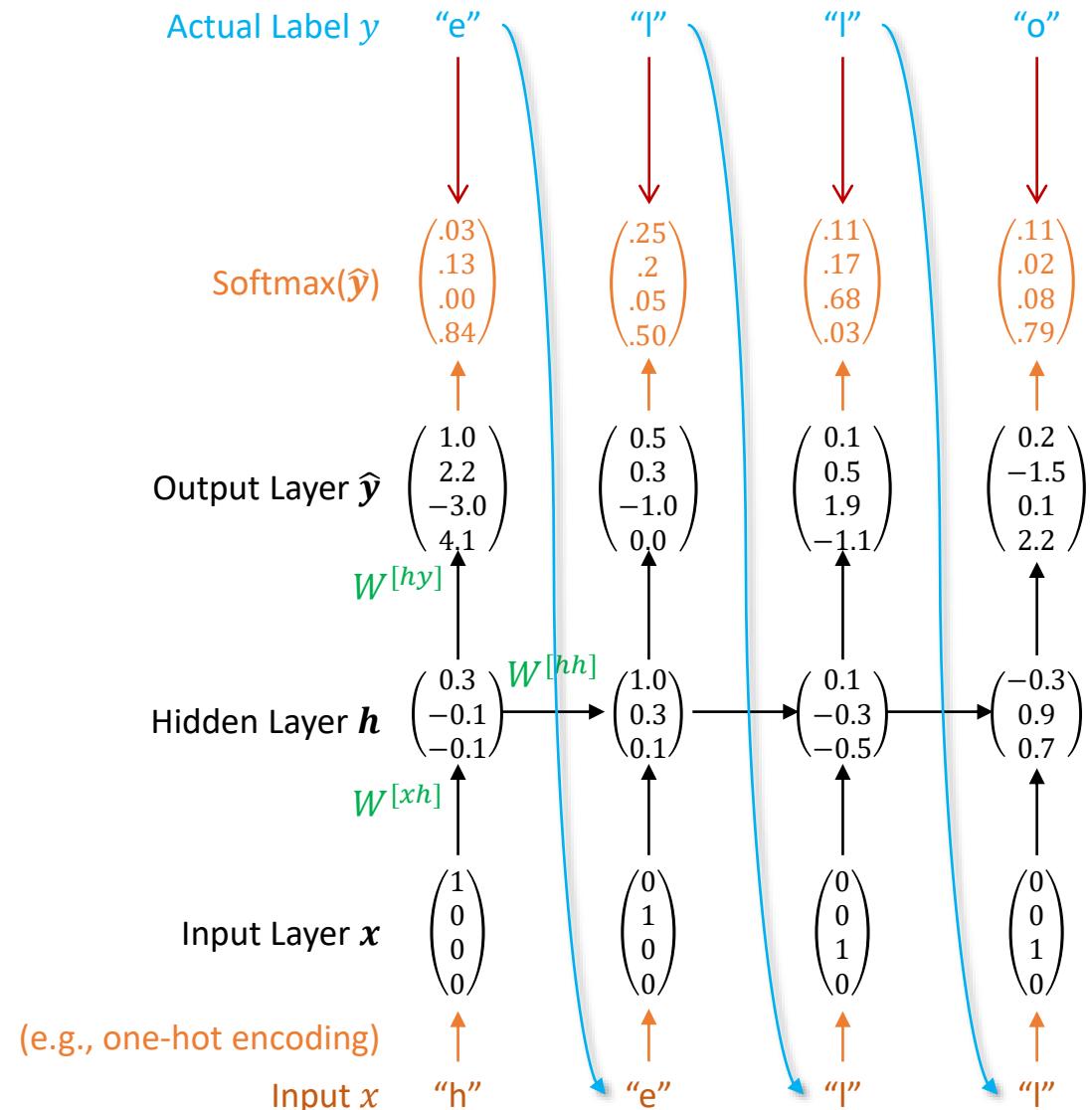
RNNs: Backpropagation Through Time



Example RNN

Text character prediction

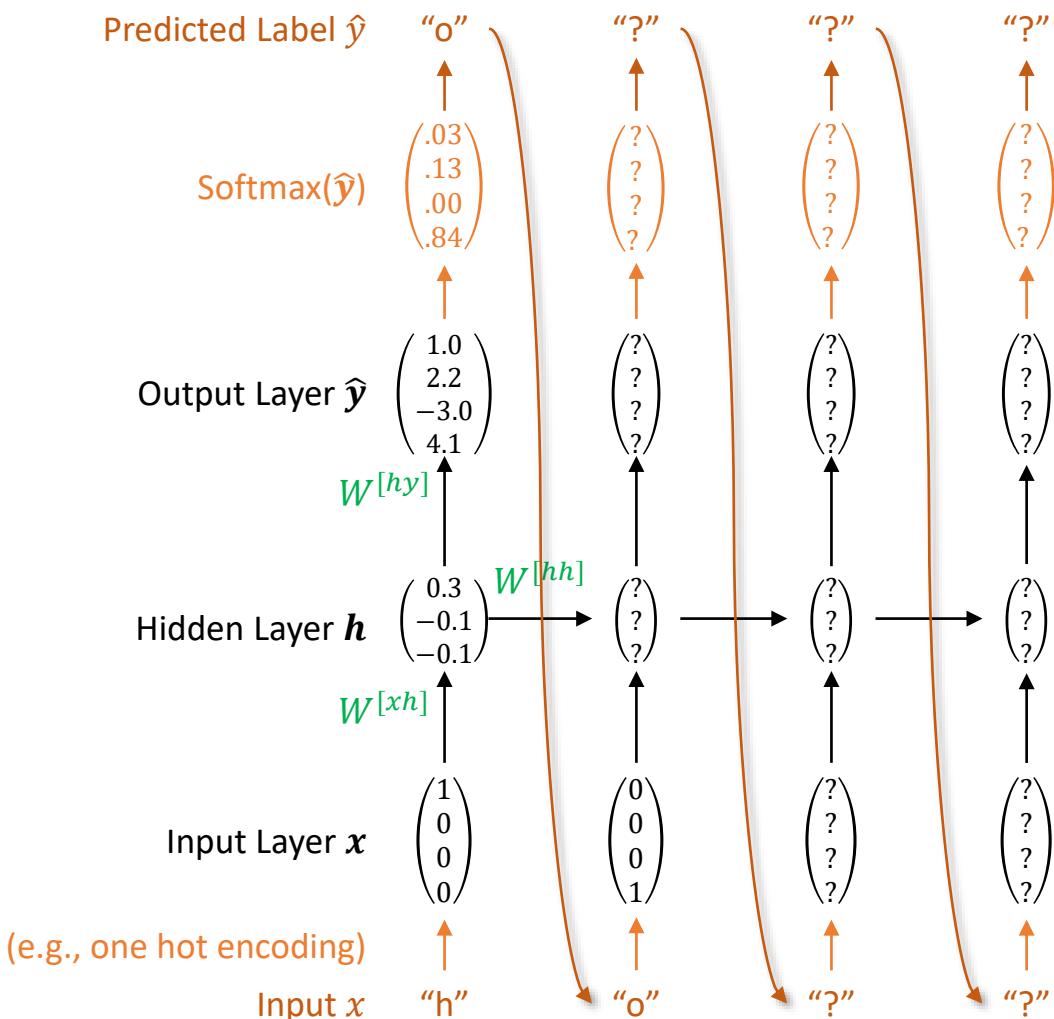
- Dictionary
 - [h, e, l, o]
- Training: sequence “hello”
- Encoding and Decoding chars
 - One-hot encoding (e.g., BOW)
 - Softmax classification
- At training time,
 - $x_t = y_{t-1}$
 - Loss (Error) is calculated as cross-entropy loss between \hat{y}_t and y_t



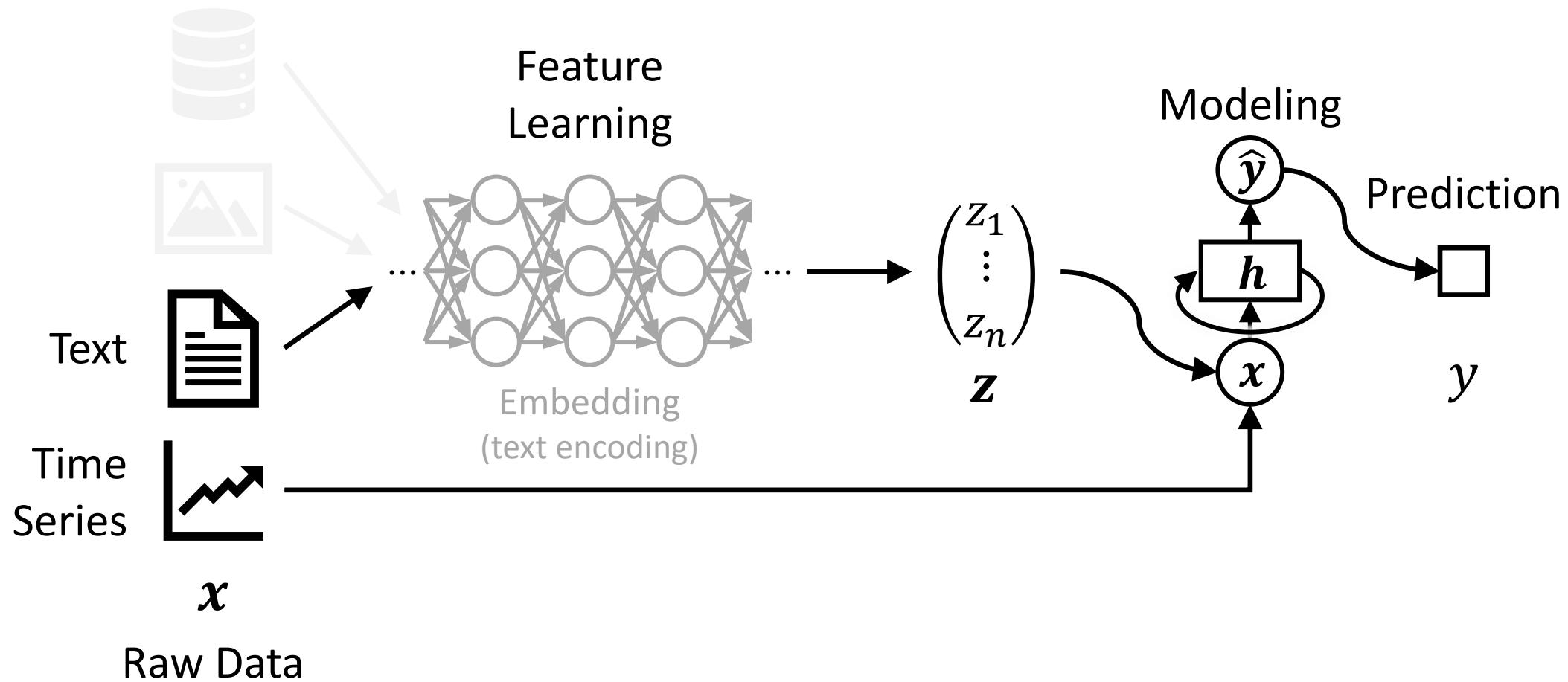
Example RNN

Text character prediction

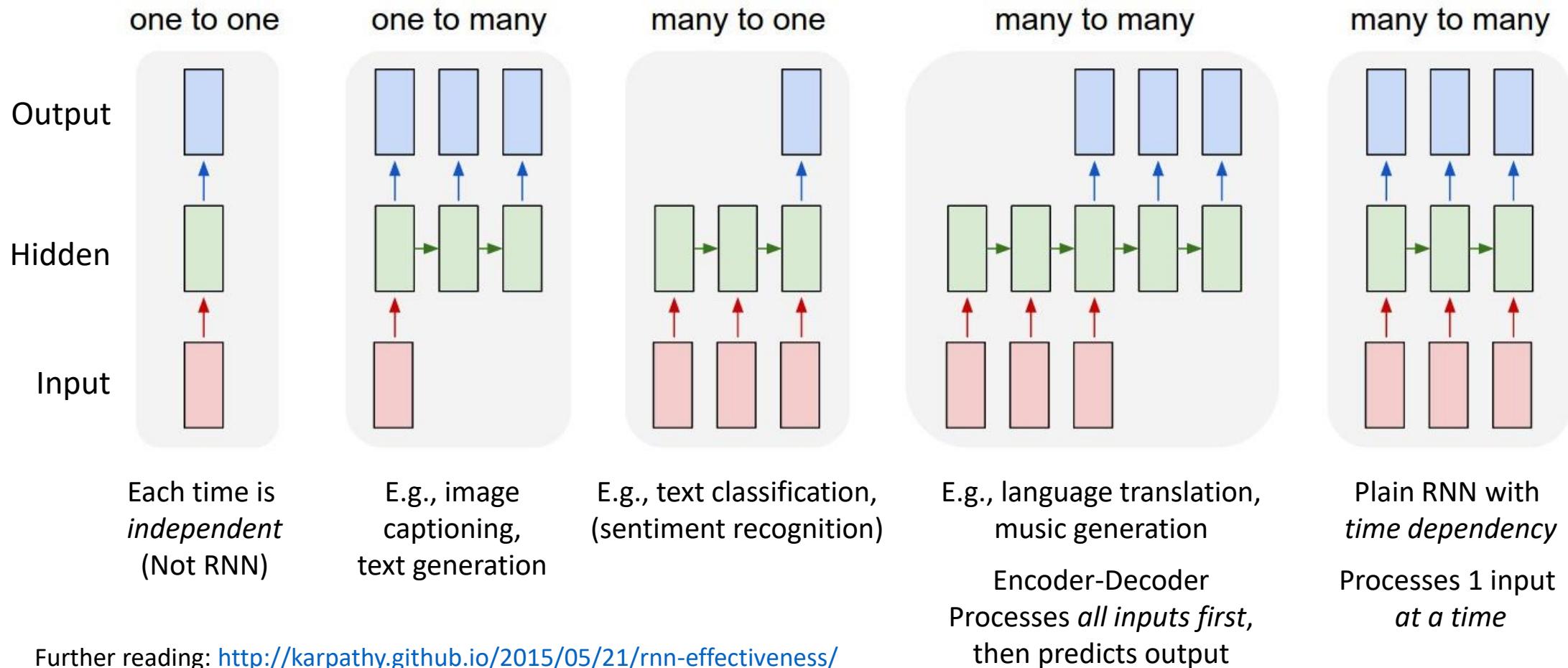
- Bag-of-Words dictionary
 - [h, e, l, o]
- At prediction time,
 - $x_t = \hat{y}_{t-1}$
 - Then propagate calculating activations to generate sequence of characters



From Manual Feature Engineering To Automatic Feature Learning

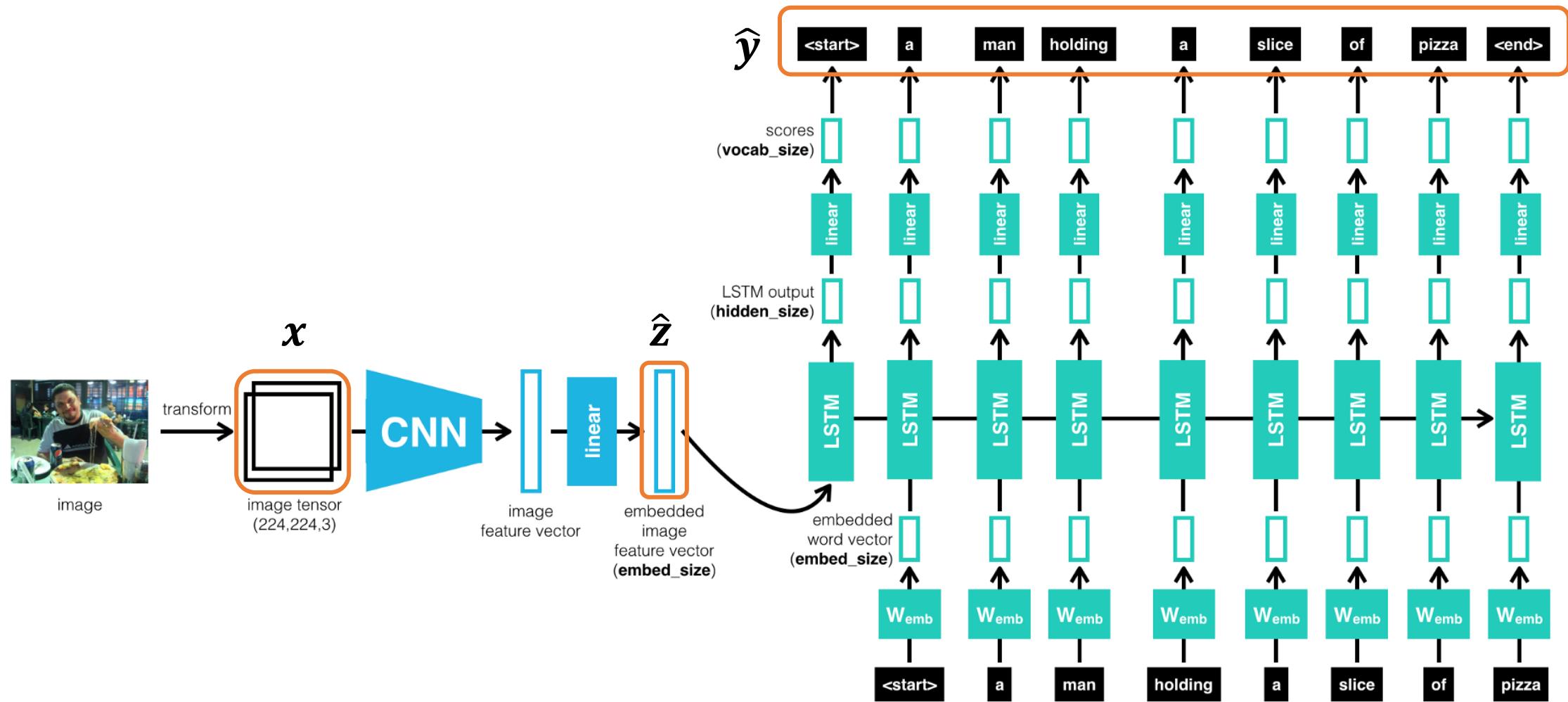


Sequence Modeling Applications



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

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





Questions!



Deep Learning Training Issues



NUS
National University
of Singapore

Department of Computer Science
School of Computing



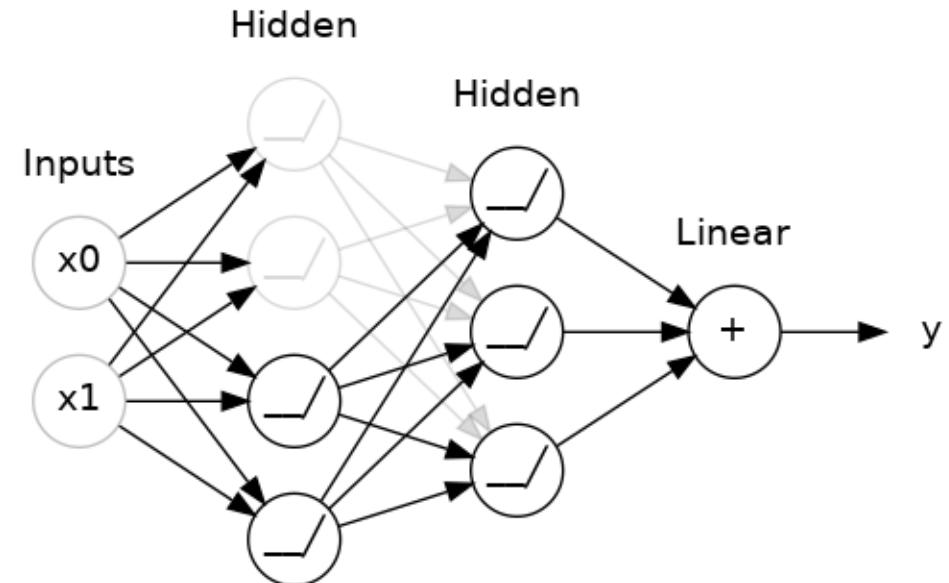
SCHOOL OF COMPUTING

Deep Learning Training Issues

- Overfitting
- Saturating Gradient Problem
- Vanishing Gradient Problem

Overfitting in deep neural networks

- Recall: what is overfitting?
- Why can deep learning overfit?
 - Too **many** parameters!
- 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

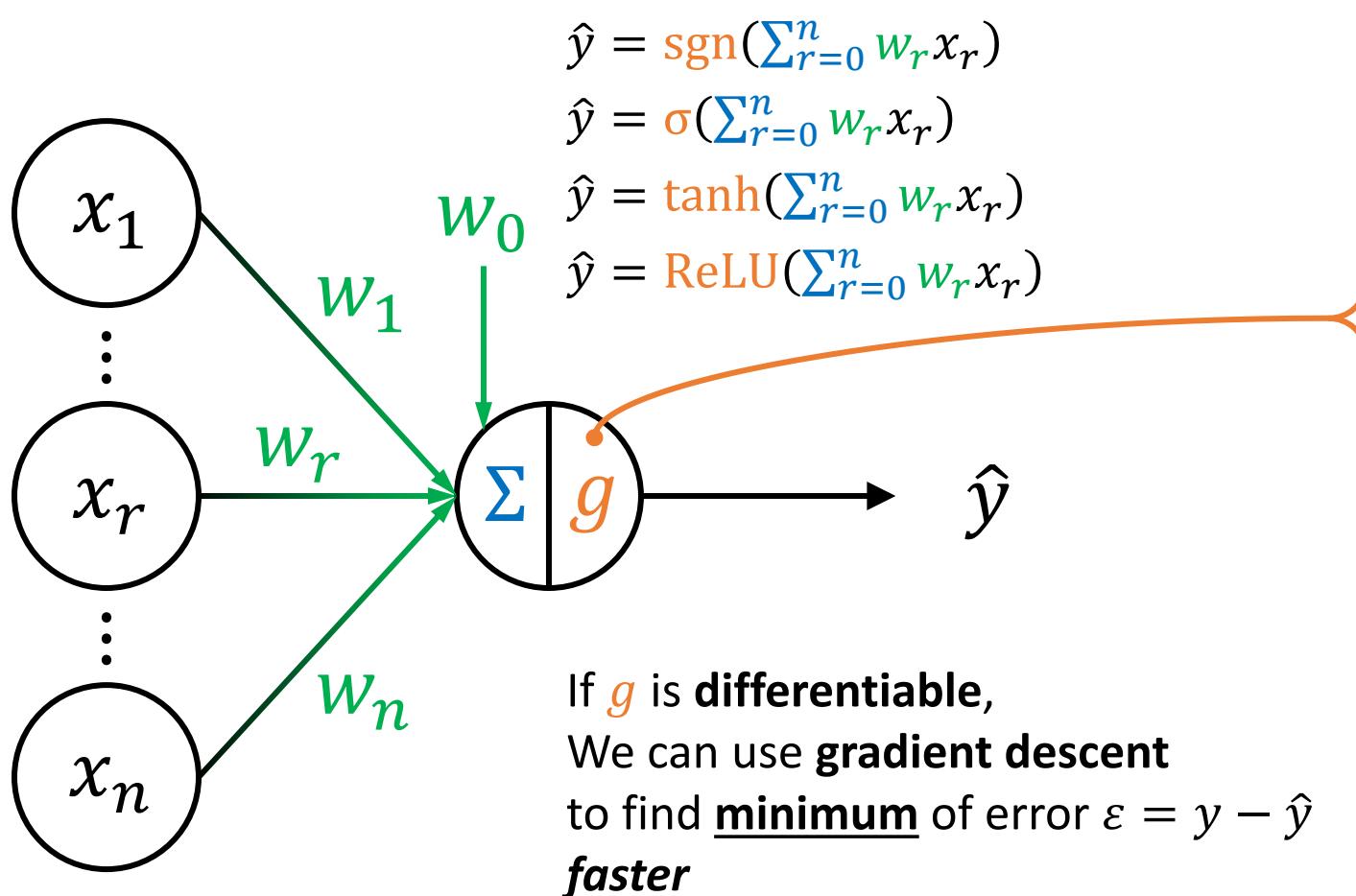


Further reading: <https://towardsdatascience.com/12-main-dropout-methods-mathematical-and-visual-explanation-58cdc2112293>

Deep Learning Training Issues

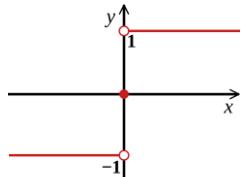
- Overfitting
- Saturating Gradient Problem
- Vanishing Gradient Problem

Differentiable Activation Functions



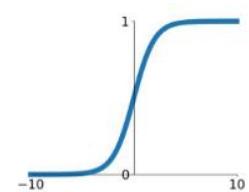
Step

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



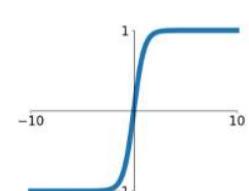
Sigmoid

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



tanh

$$\tanh(x)$$



ReLU

$$\max(0, x)$$

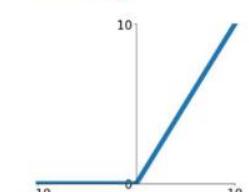


Image credit.

https://miro.medium.com/max/1400/0*slJ-gbjlz0rz8lb.png

Gradient Descent Weight Update

Reference

$$\mathbf{a}^{[l]} = g^{[l]}(\mathbf{f}^{[l]})$$

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

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \nabla \varepsilon$$

MSE error

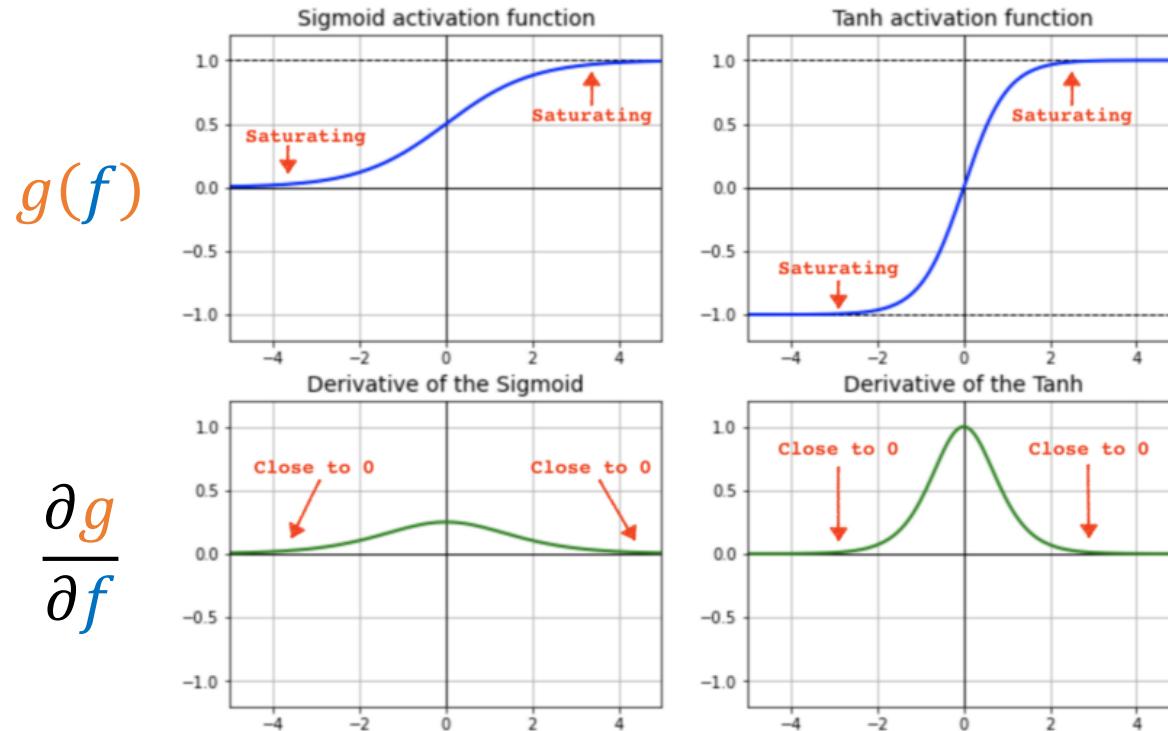
$$\varepsilon = \frac{1}{2} (\hat{y} - y)^2$$

Gradient of error

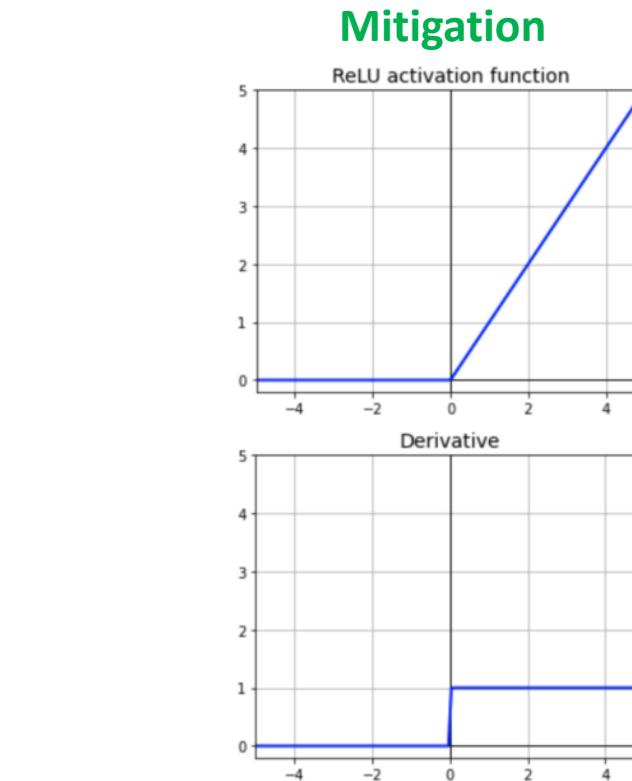
$$\nabla \varepsilon = \frac{\partial \varepsilon}{\partial \mathbf{W}} = \frac{\partial \varepsilon}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \mathbf{W}}$$

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

Saturating Gradient Problem due to activation functions Mitigate with ReLU activation function



When x value far from 0, gradient $\rightarrow 0$ (saturating)
When gradient ≈ 0 , then weights don't update much



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

Vanishing Gradient Problem



$$\hat{y}'(\mathbf{W}^{[1]}) = \frac{\partial \mathbf{g}^{[L]}}{\partial \mathbf{W}^{[1]}} = \frac{\partial \mathbf{f}^{[1]}}{\partial \mathbf{W}^{[1]}} \frac{\partial \mathbf{g}^{[1]}}{\partial \mathbf{f}^{[1]}} \dots \frac{\partial \mathbf{g}^{[l]}}{\partial \mathbf{f}^{[l]}} \frac{\partial \mathbf{f}^{[l+1]}}{\partial \mathbf{g}^{[l]}} \frac{\partial \mathbf{g}^{[l+1]}}{\partial \mathbf{f}^{[l+1]}} \dots \frac{\partial \mathbf{f}^{[L]}}{\partial \mathbf{g}^{[L-1]}} \frac{\partial \mathbf{g}^{[L]}}{\partial \mathbf{f}^{[L]}}$$

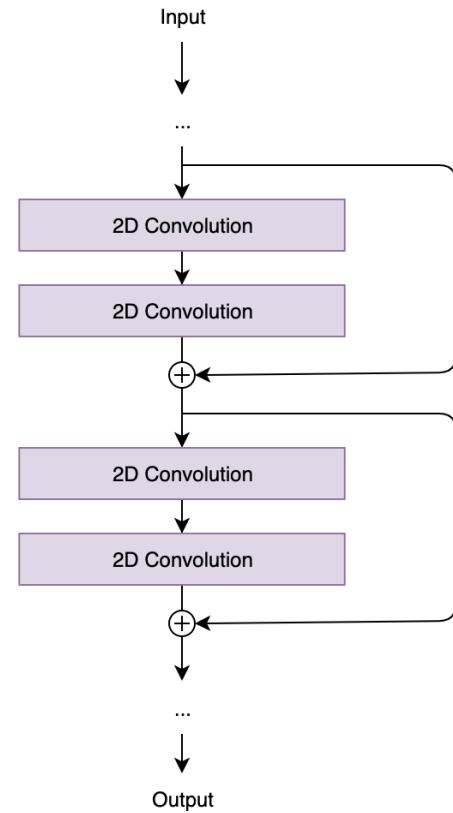
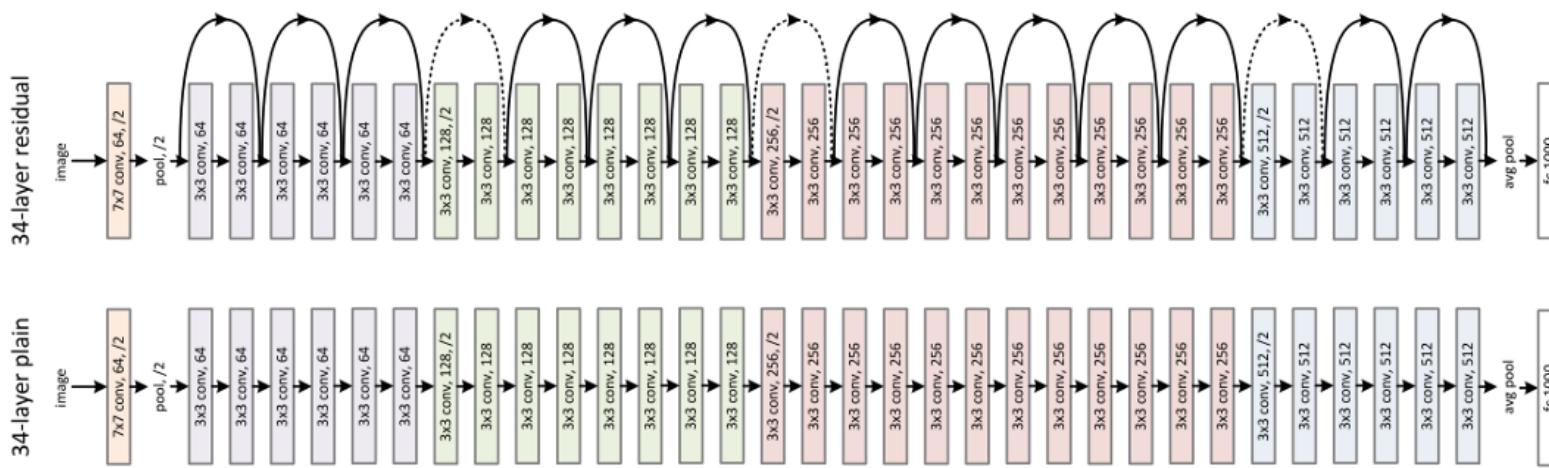
If some gradients are small (< 1),
multiplying many small numbers equals a very small number.

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

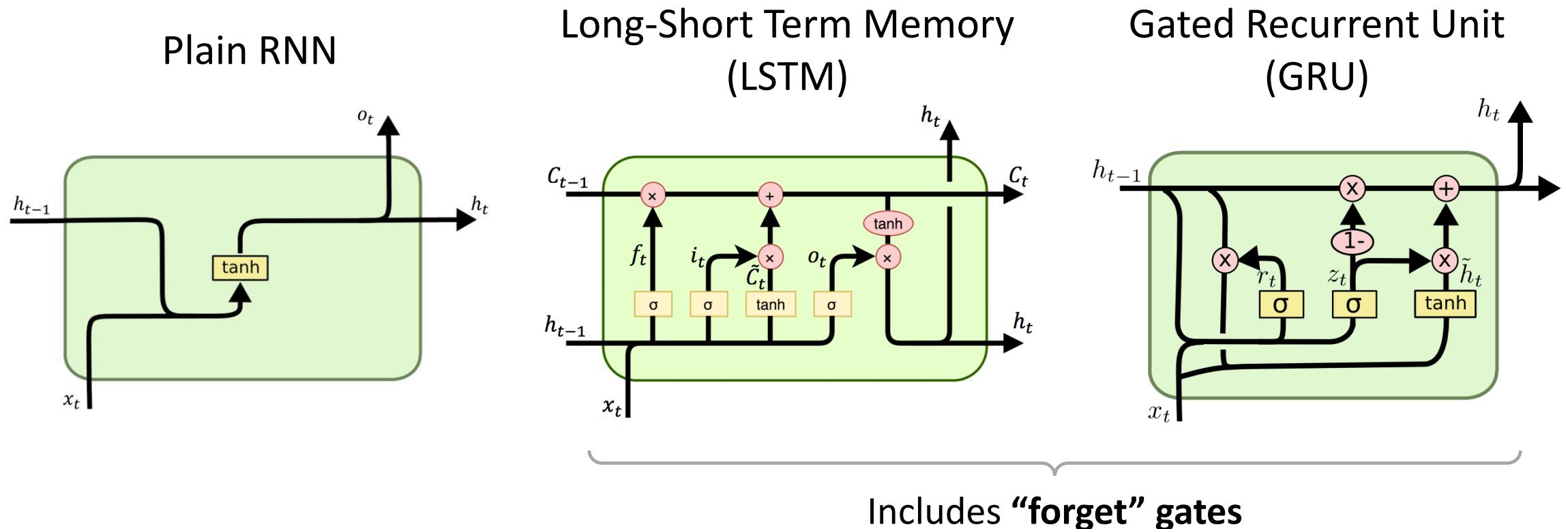


Image Credit: <http://dprogrammer.org/rnn-lstm-gru>

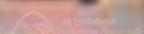
Further reading: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Wrapping Up



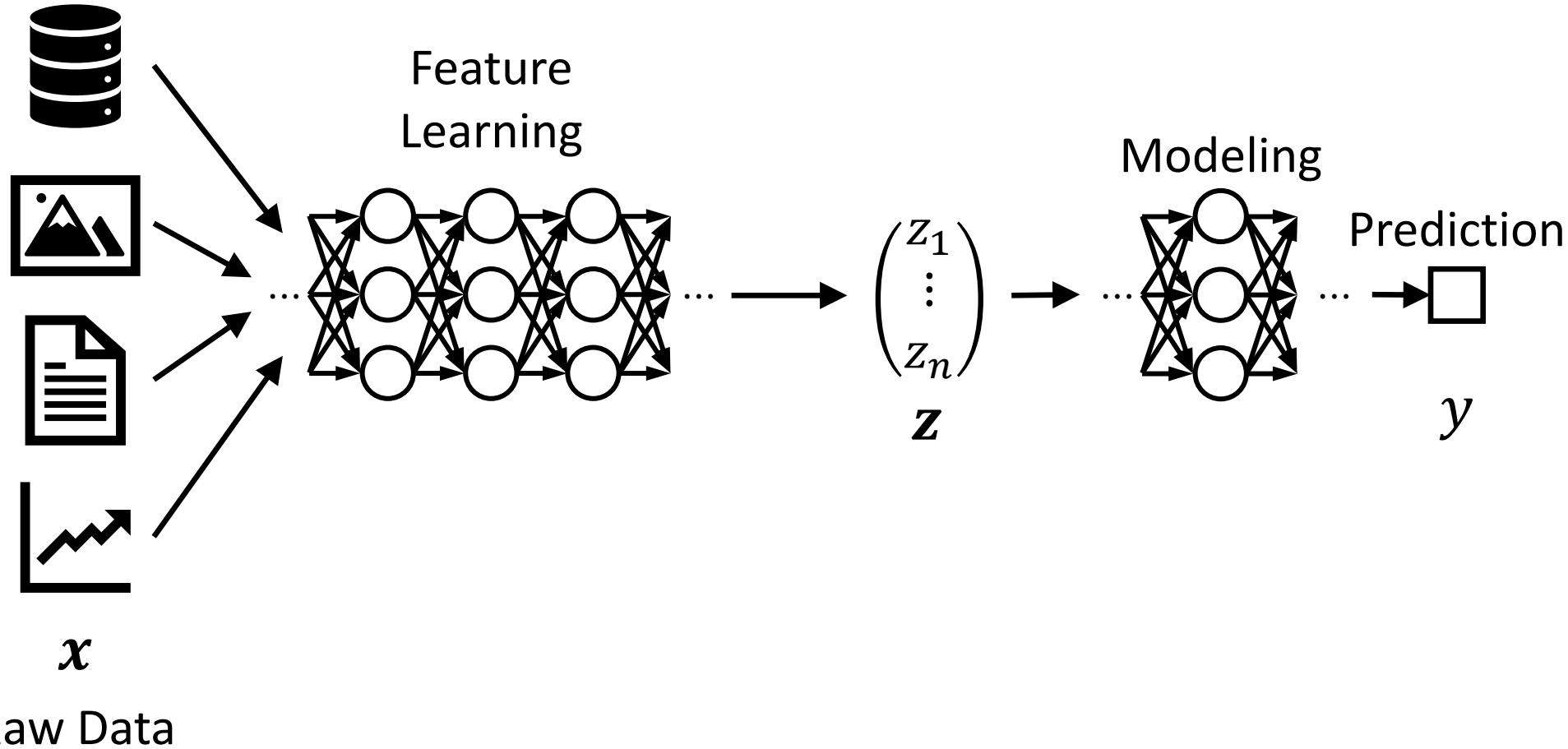
NUS
National University
of Singapore

Department of Computer Science
School of Computing



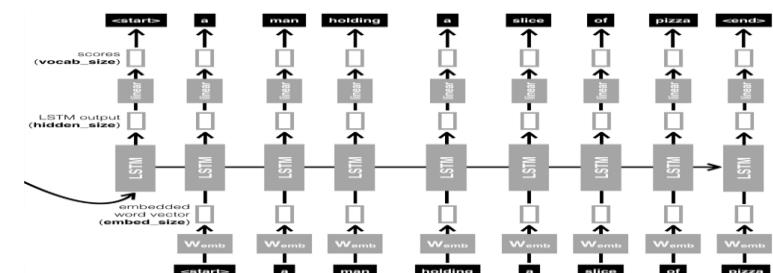
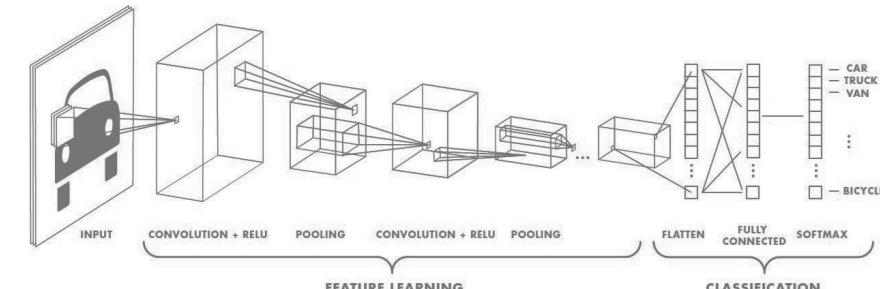
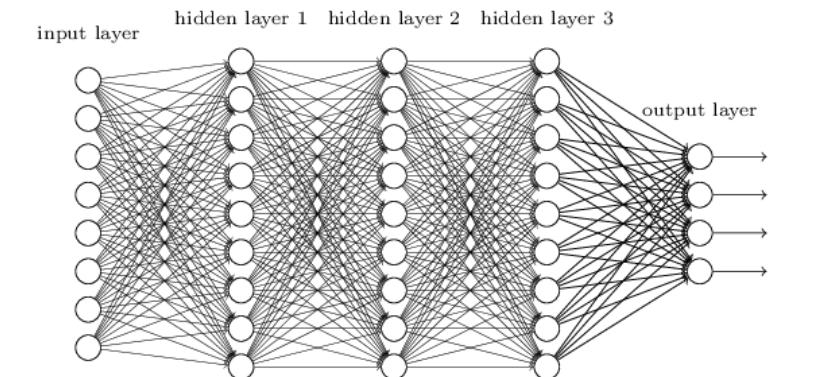
SCHOOL OF COMPUTING

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**



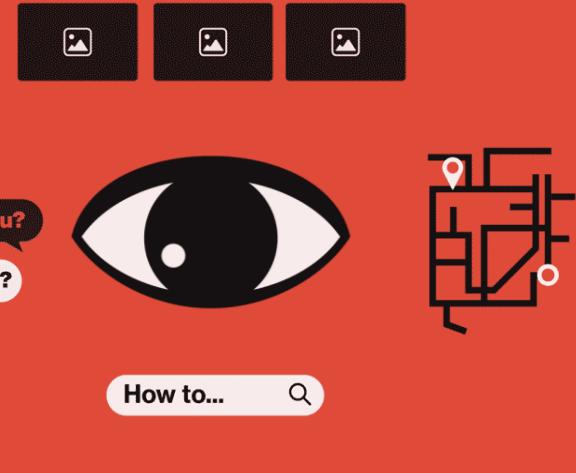
Grand issues with AI (Deep Learning)



Lack of **Explainability**
[W10b]



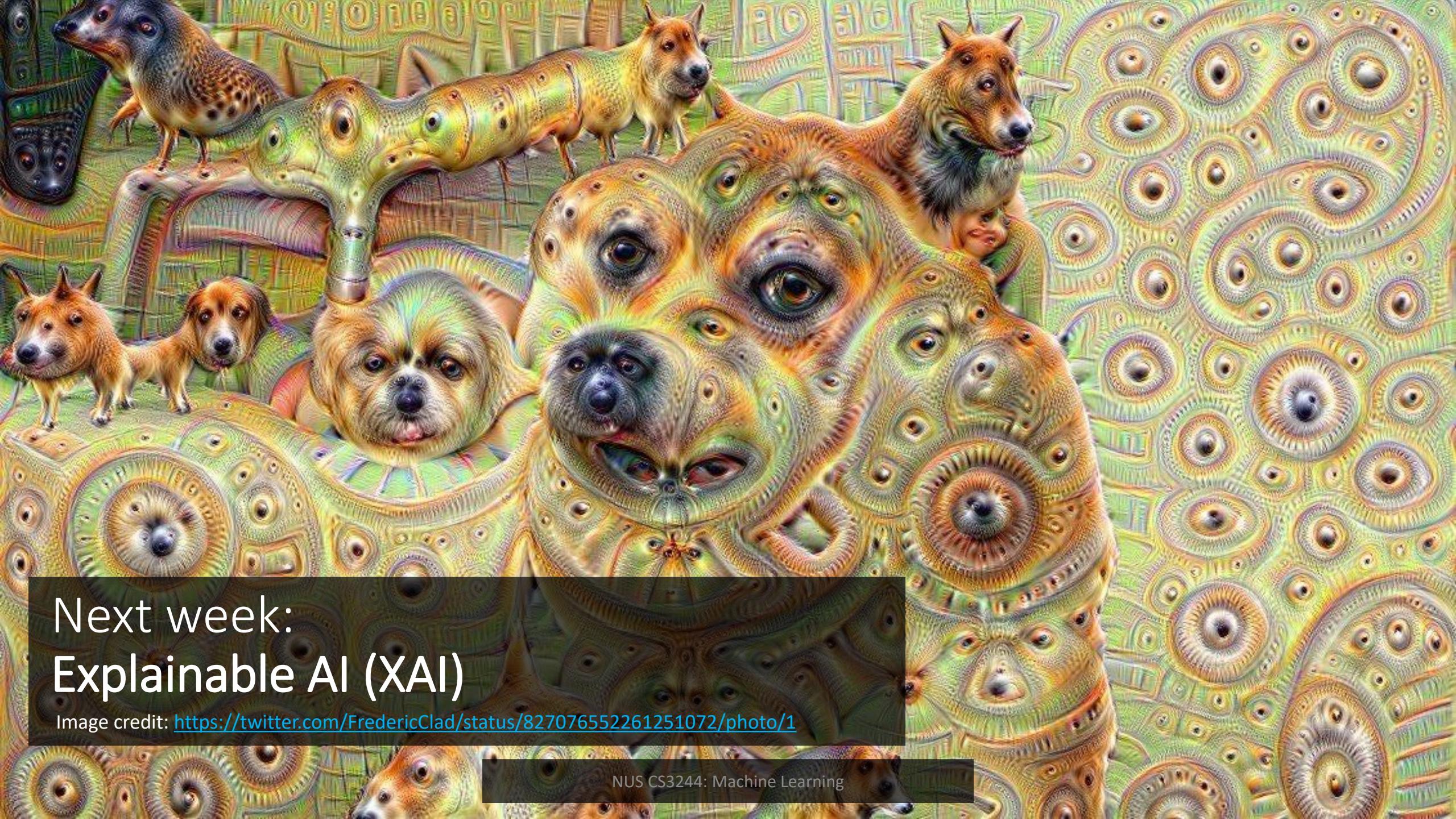
Algorithmic Bias (Societal)
[W12a]



Data Privacy

Image credits:

https://miro.medium.com/max/2000/1*H4cW-RCyHpu5FNtVaAPoQ.gif
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>