



SAL

SILICON AUSTRIA LABS

VLSI FOR MACHINE LEARNING

PEDRO JULIAN, DIEGO GIGENA, NICOLÁS
RODRÍGUEZ

CAE 2023, Córdoba



SAL
SILICON AUSTRIA LABS

01

What is a Neuron

02

Deep Neural Networks

03

Activation Functions

04

Convolutional Neural Networks

2. BASIC BLOCKS OF DNN

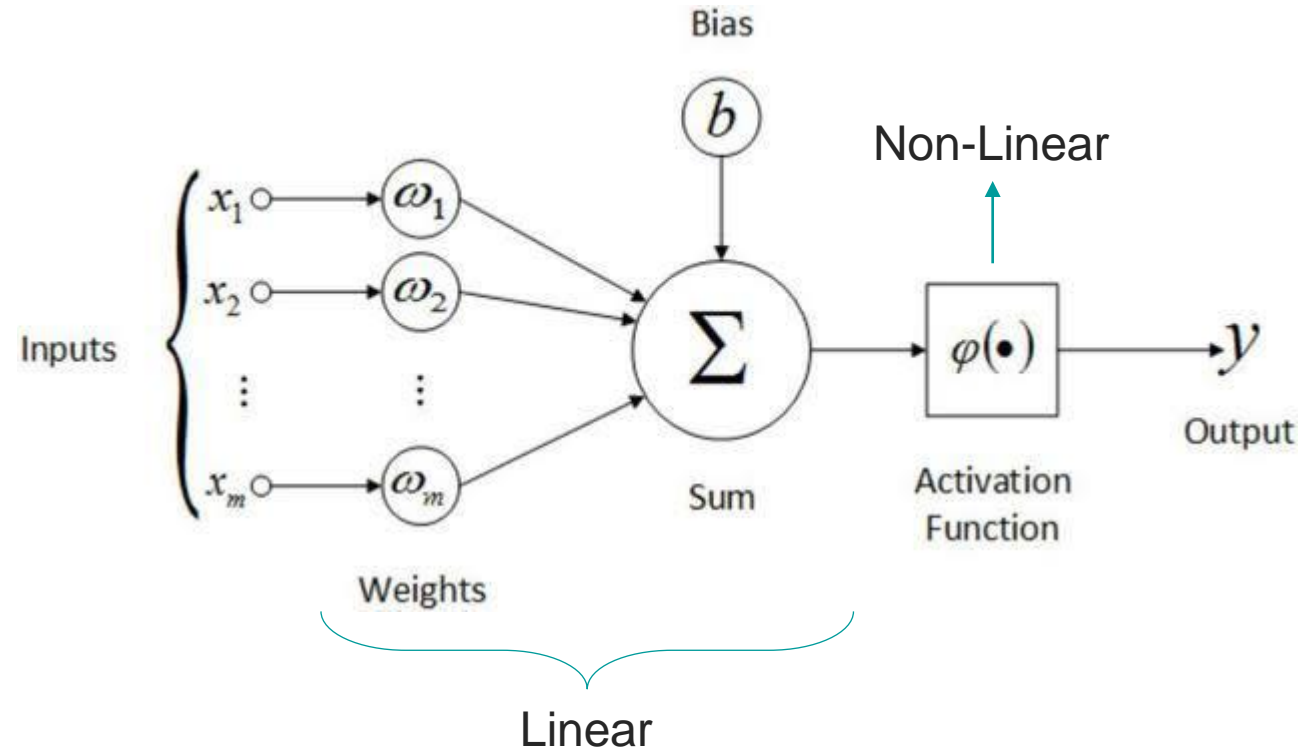


SAL
SILICON AUSTRIA LABS

WHAT IS A NEURON

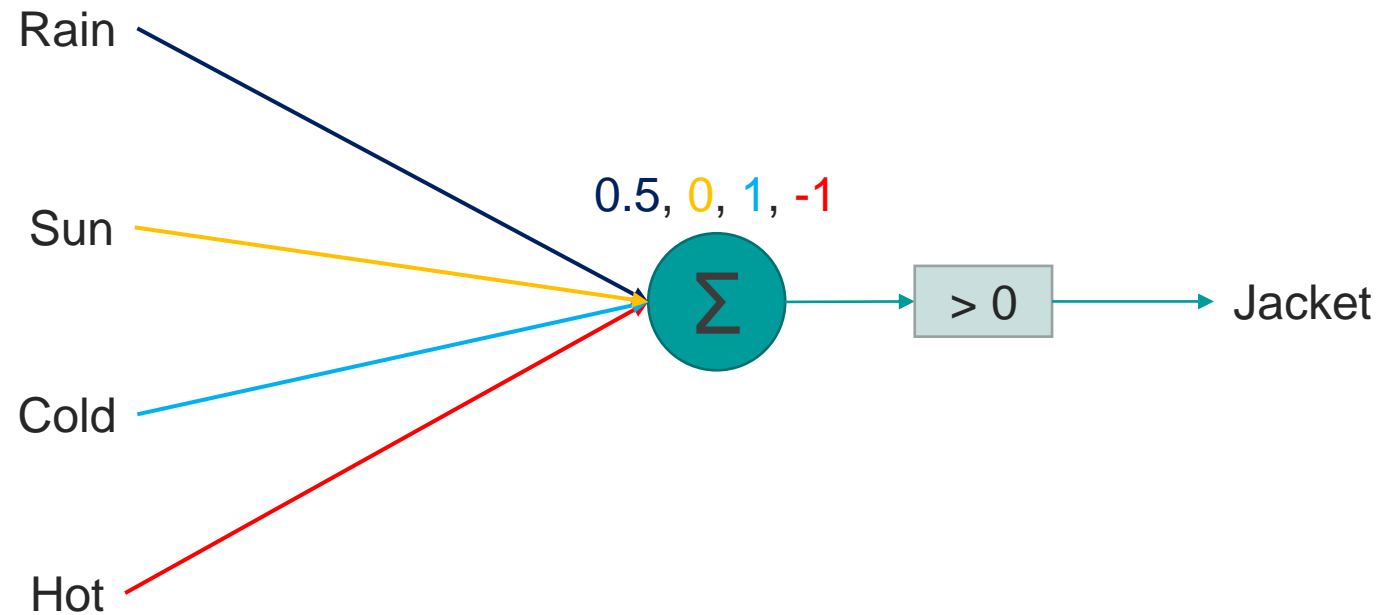
THE PERCEPTRON

- ≡ Simple model of a Neuron
- ≡ Foundational Element of NN
- ≡ Components:
 - ≡ Input Values
 - ≡ Weighted connections
 - ≡ Adder with Bias
 - ≡ Activation Function (AF)
- ≡ How does it work?
 - ≡ Addition of weighted values
 - ≡ Decision using AF



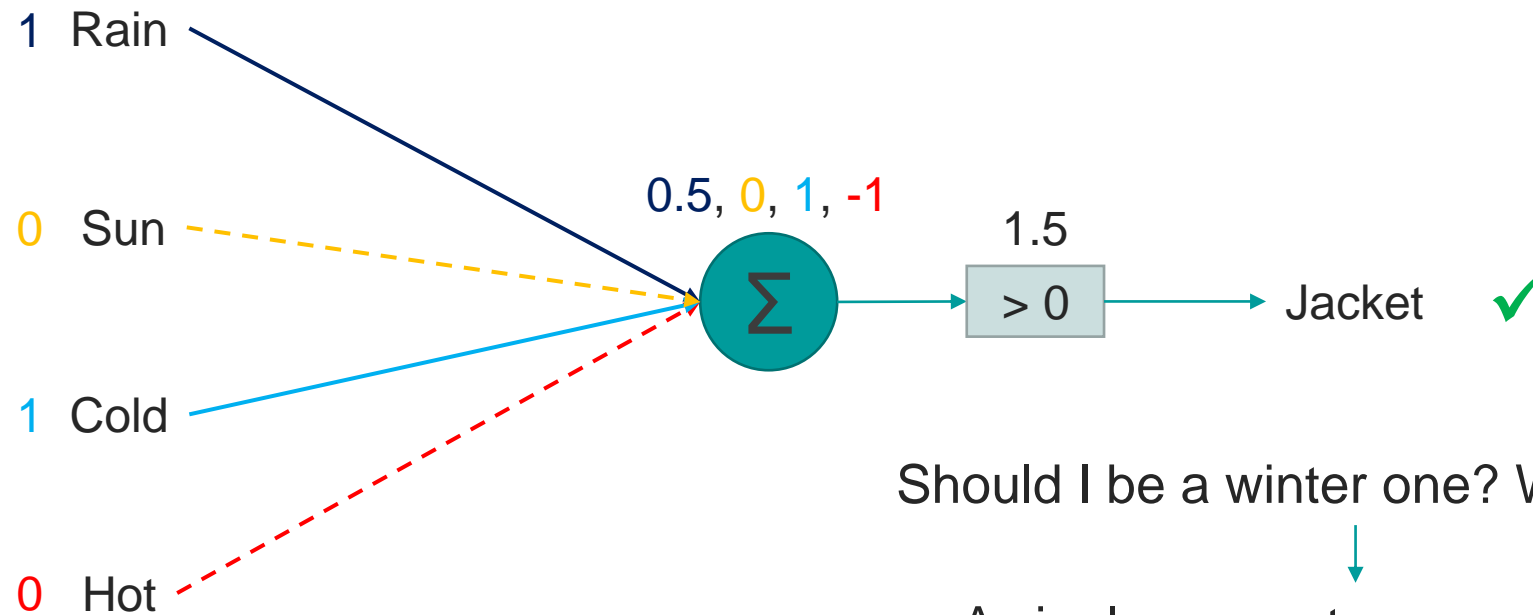
EXAMPLE

How's the weather like? → Should I wear a jacket?



EXAMPLE

How's the weather like? → Should I wear a jacket?



WHY ARTIFICIAL NEURONS

-
- ≡ Relatively simple models of a real neuron
 - ≡ Group of neurons as universal approximator
 - ≡ Can solve complex tasks
 - ≡ And most importantly...

They can learn!!!

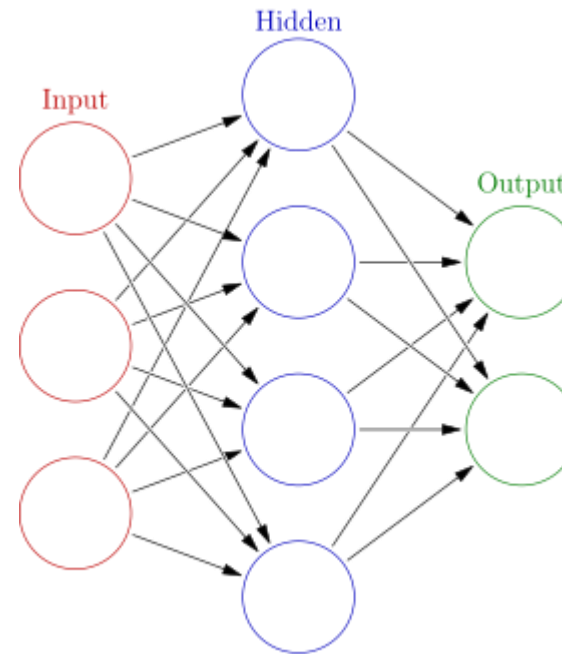


SAL
SILICON AUSTRIA LABS

DEEP NEURAL NETWORKS

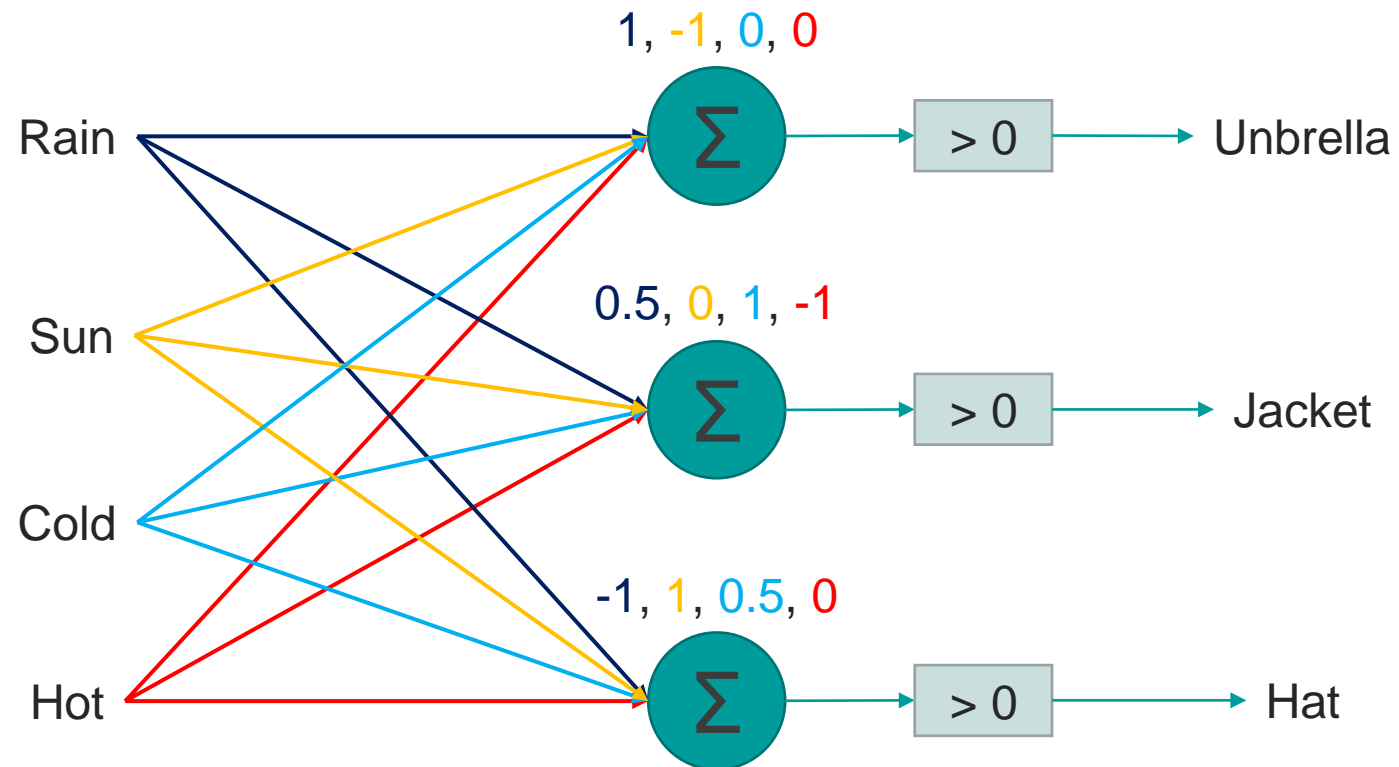
ARTIFICIAL NEURAL NETWORK

- ≡ Group of neurons to solve more complex tasks
- ≡ Groups with shared inputs, between input and output, compose a “Hidden Layer”
- ≡ When all inputs are shared:
 - ≡ Fully-Connected Layer
 - ≡ Dense Layer
- ≡ A single Hidden Layer is enough to be called ANN



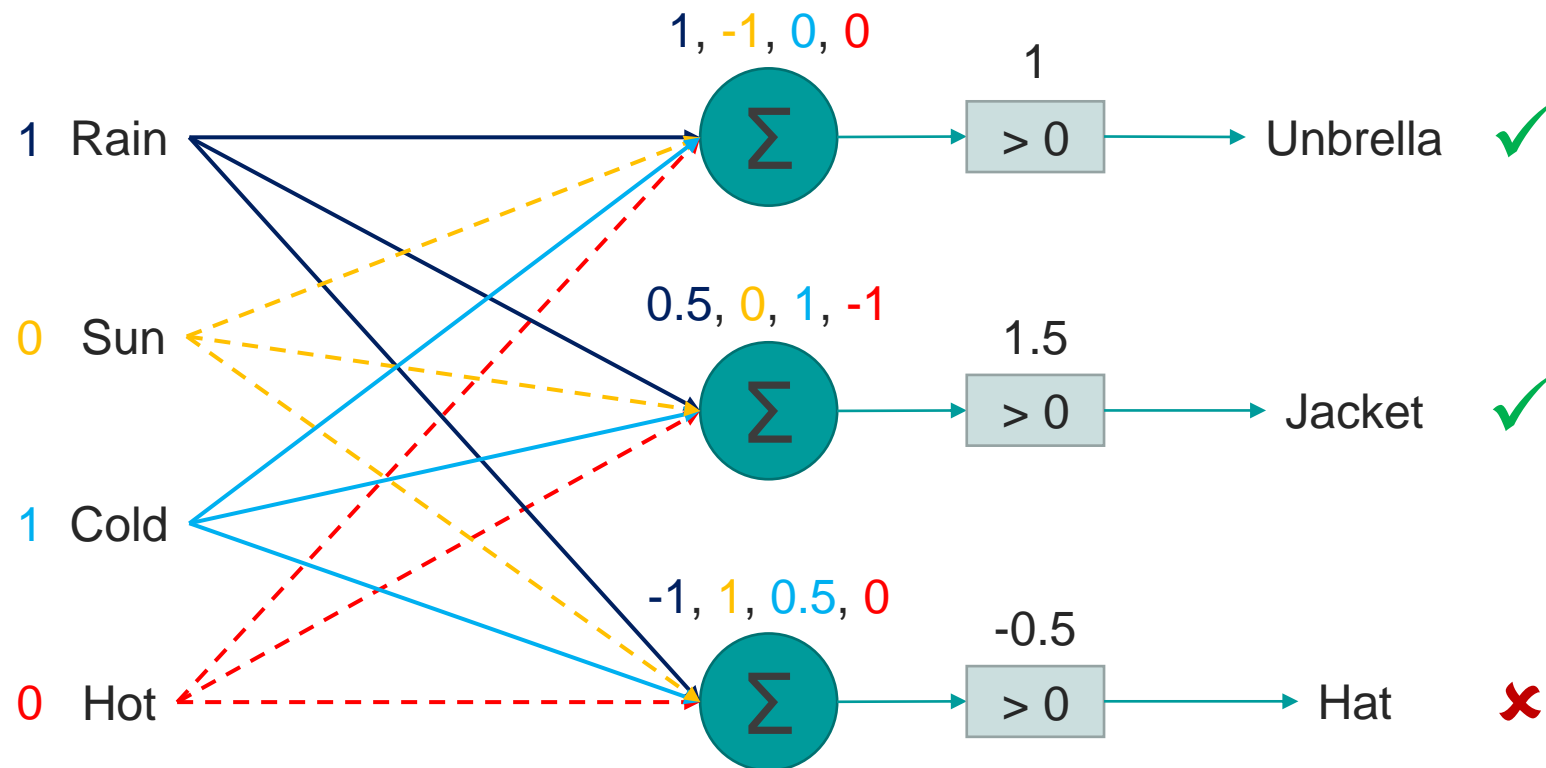
EXPANDED EXAMPLE

How's the weather like? \longrightarrow What should I wear?



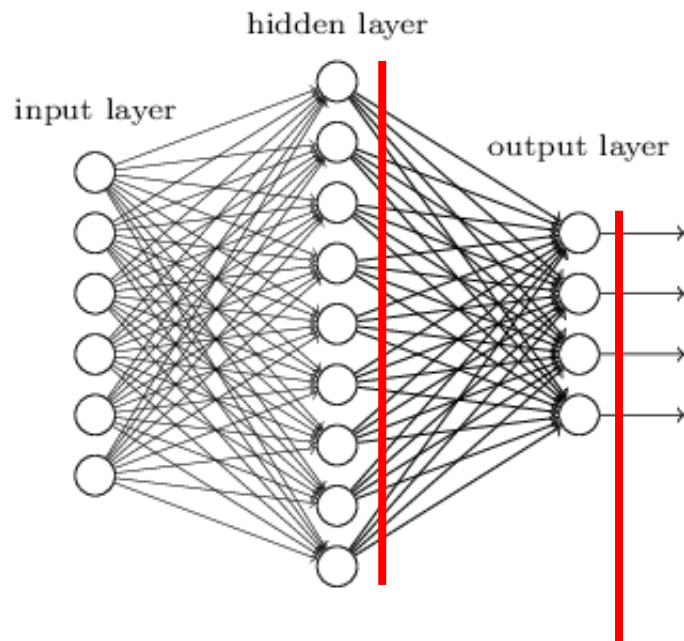
EXPANDED EXAMPLE

How's the weather like? \longrightarrow What should I wear?



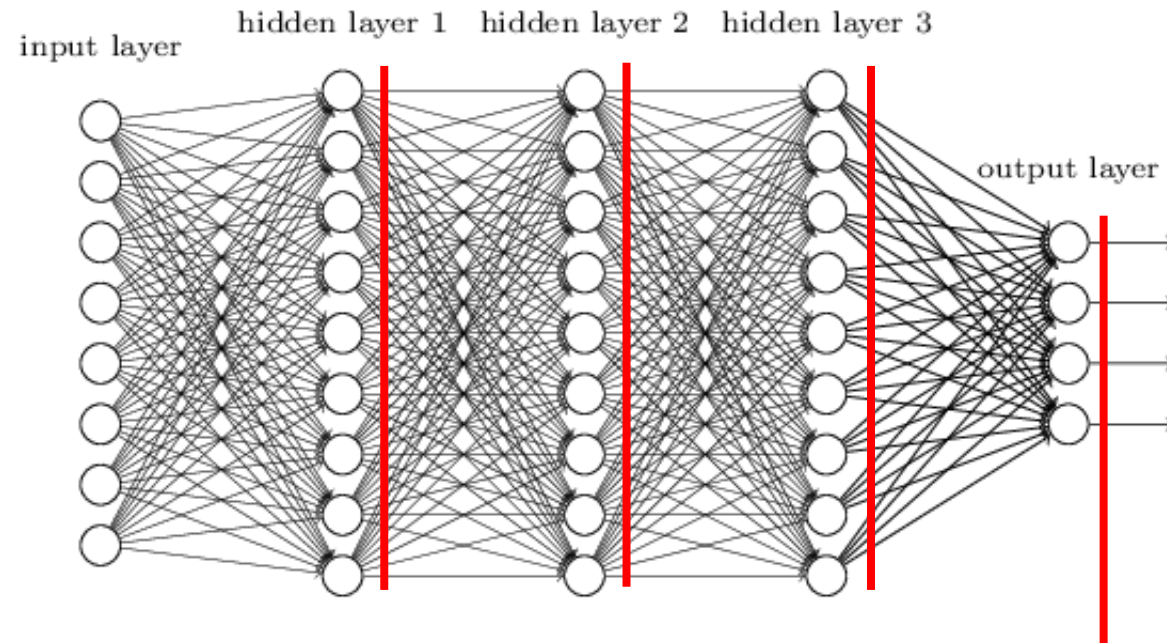
WHAT IS DEEP

"Non-deep" feedforward
neural network



Umbrella, Jacket and/or Hat

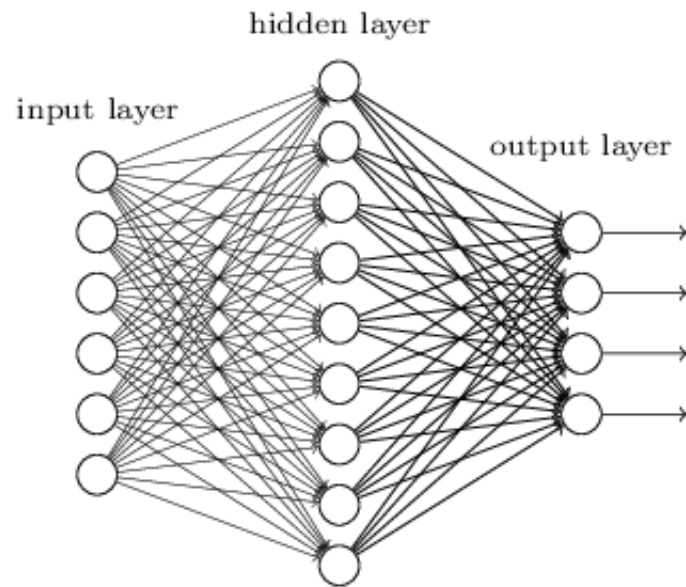
Deep neural network



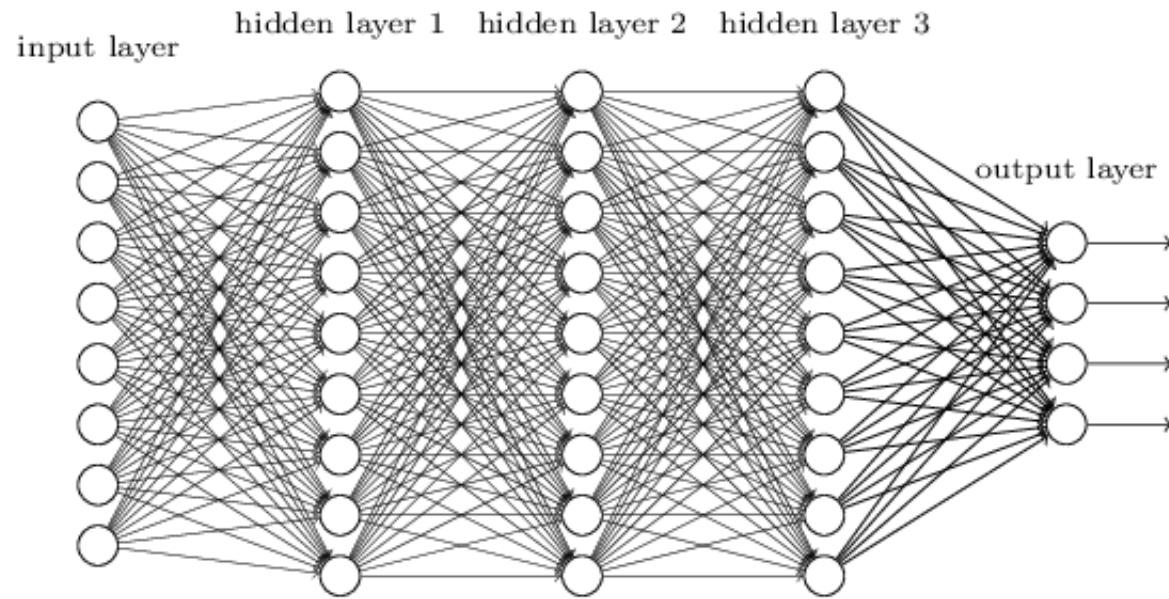
Winter waterproof Jacket

WHAT IS DEEP

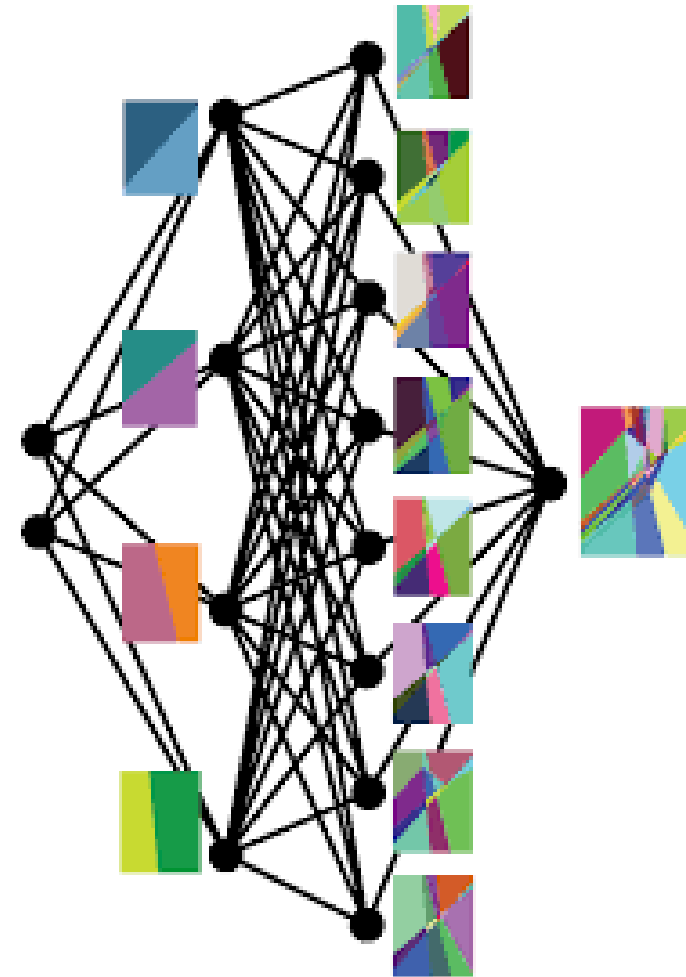
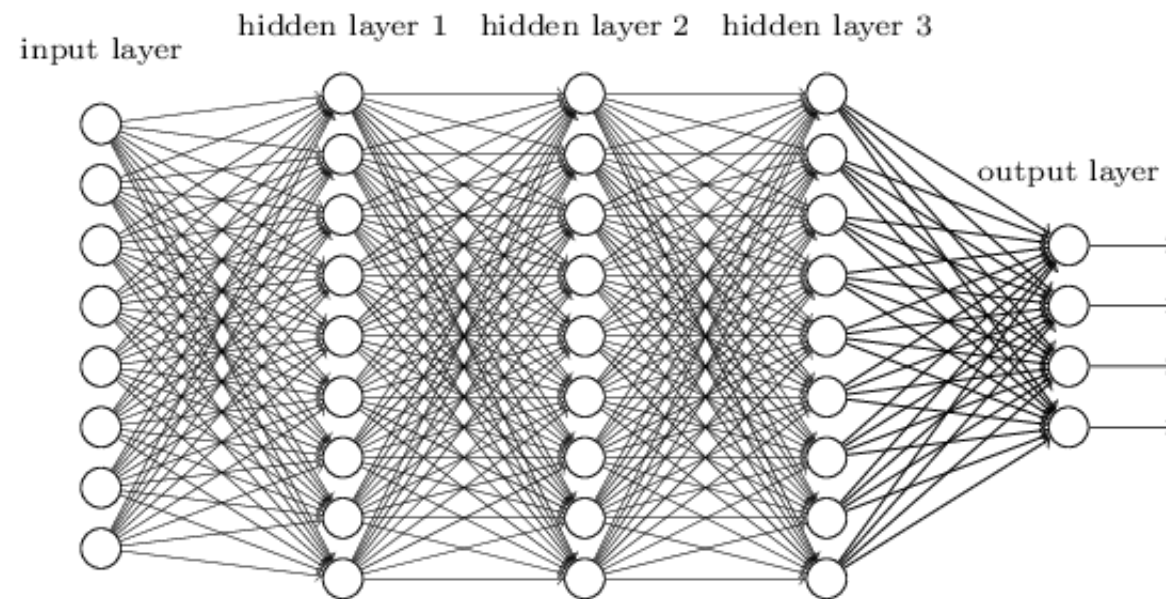
≡ No Activation Functions?



≡

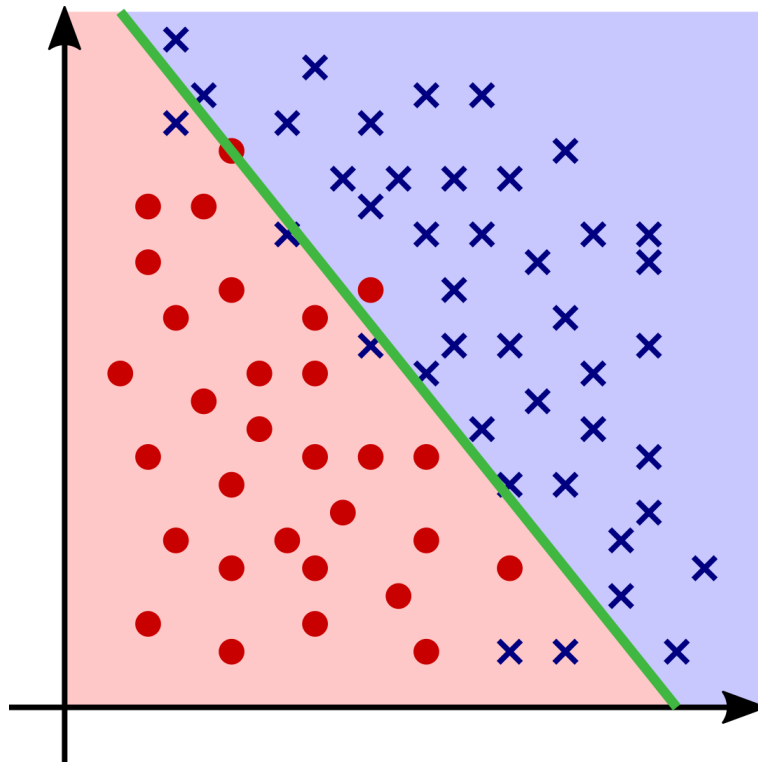


LINEAR REGIONS

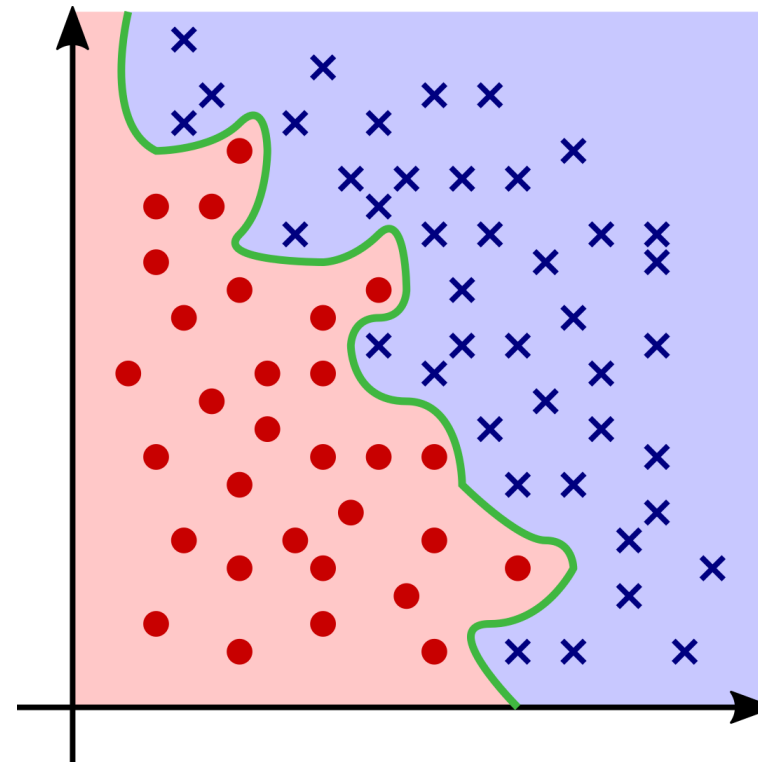


ADDING LAYERS

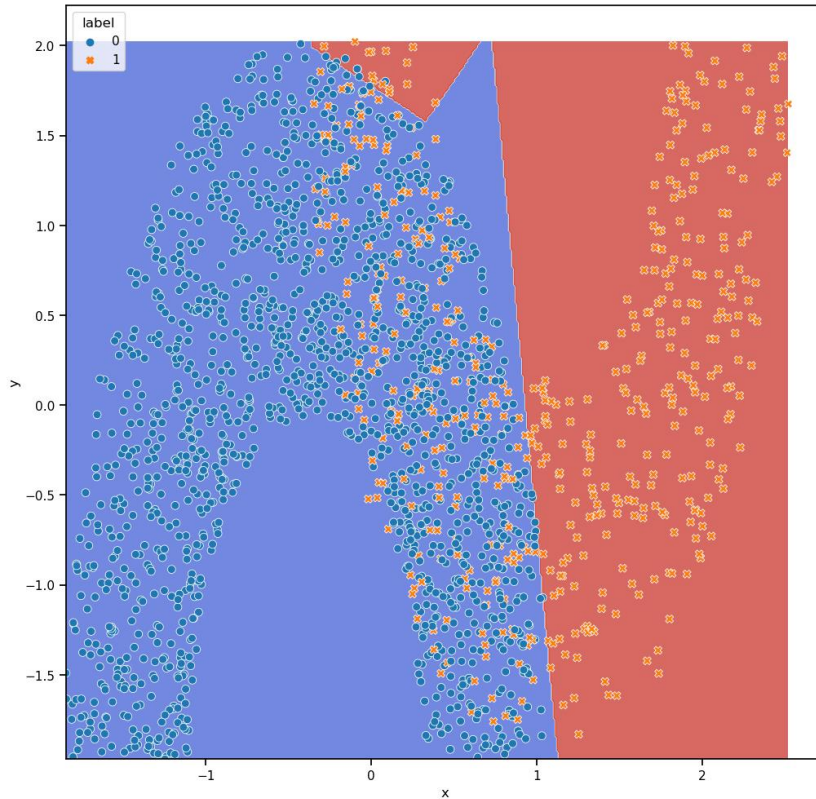
Single Neuron



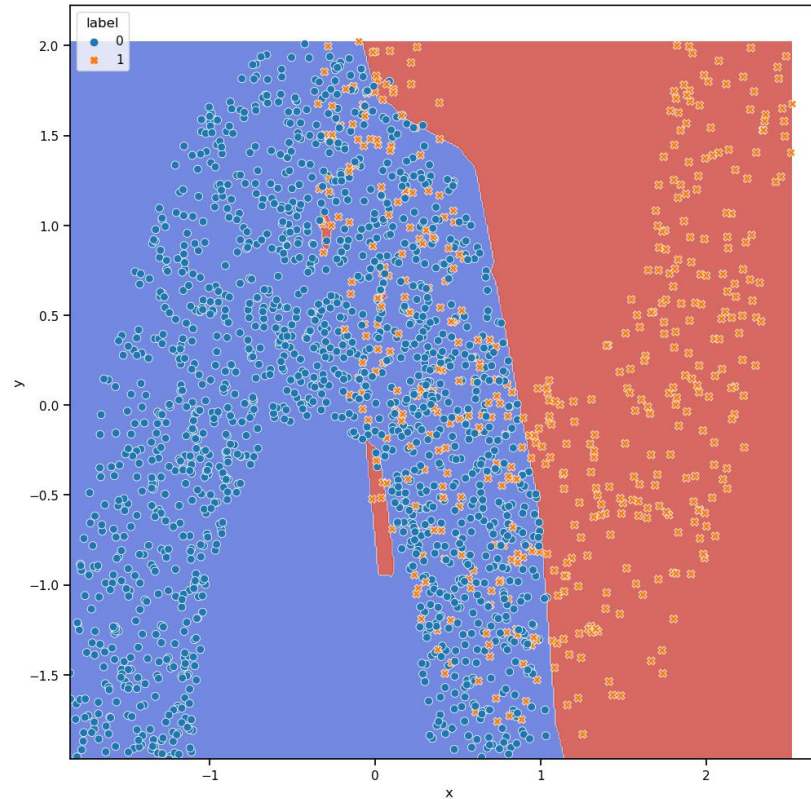
Multiple Neurons/Layers



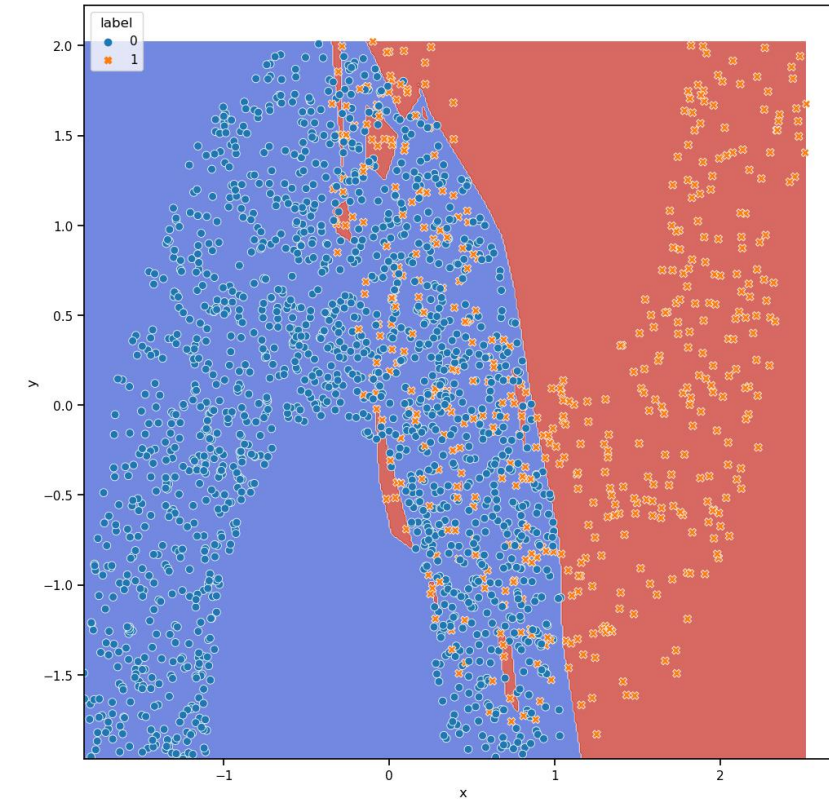
ADDING LAYERS



Single Hidden Layer, 10 neurons



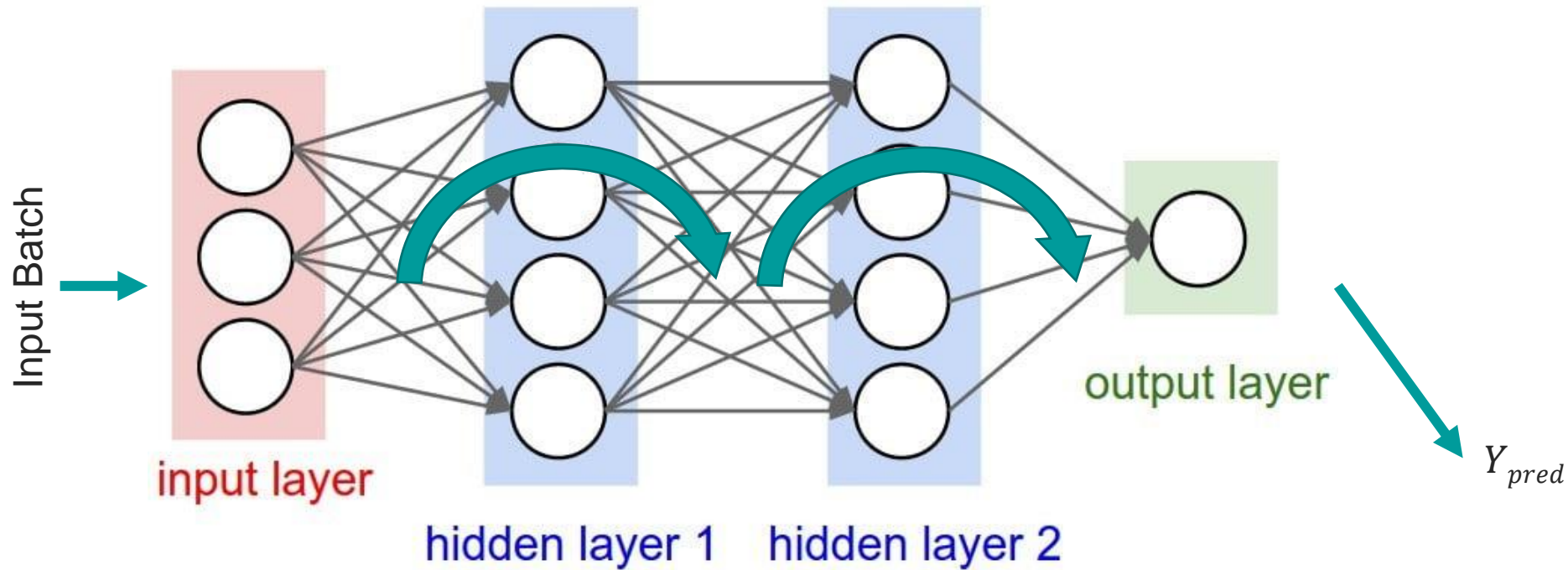
4 Hidden Layers, 10 neurons each



8 Hidden Layers, 10 neurons each

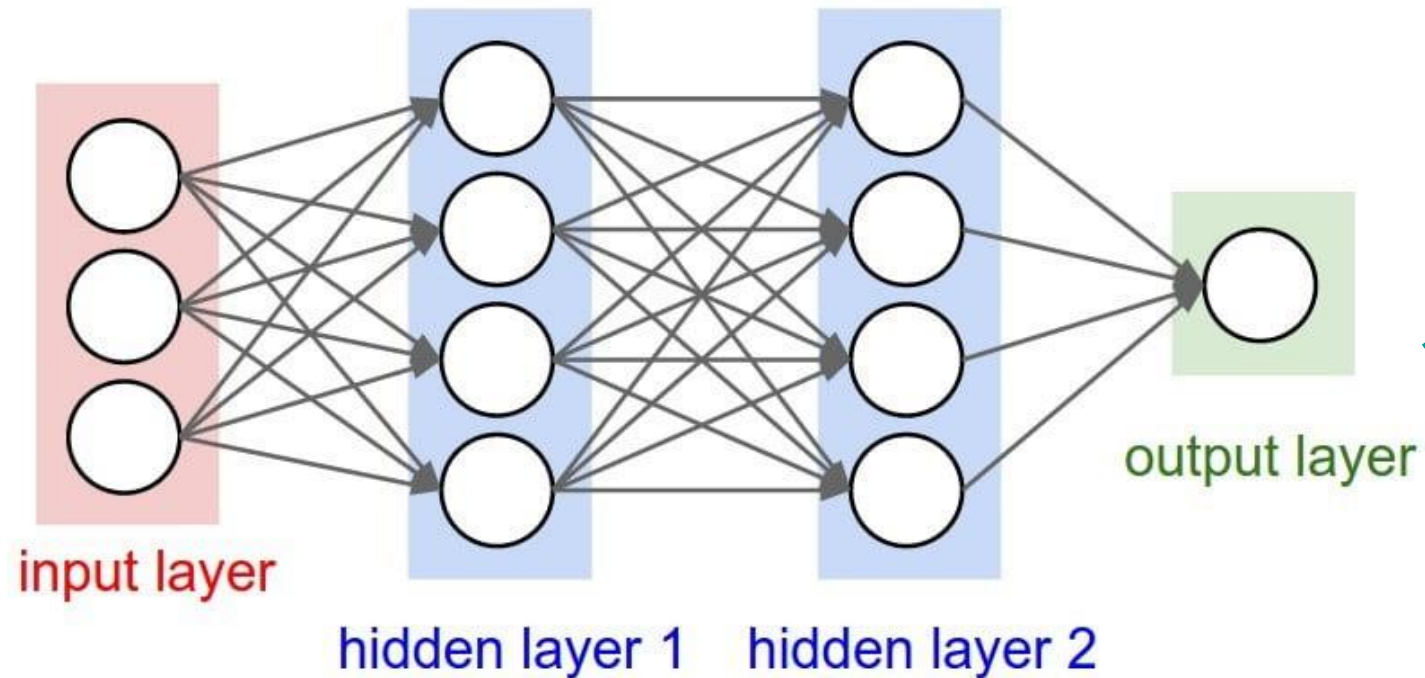
BACKPROPAGATION

Forward Pass

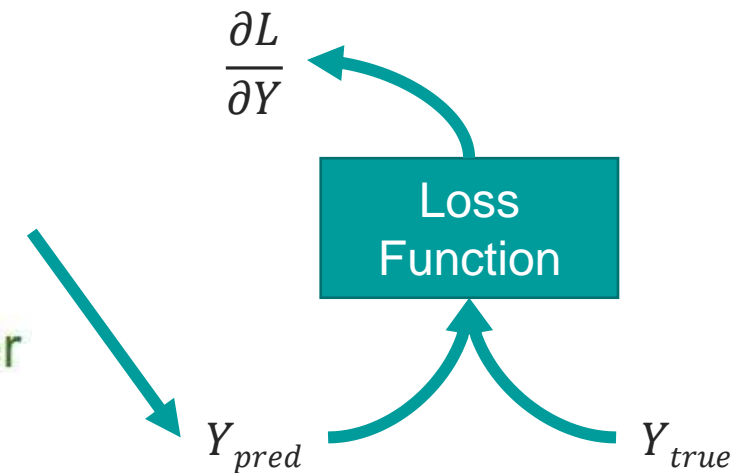


BACKPROPAGATION

Loss Computation

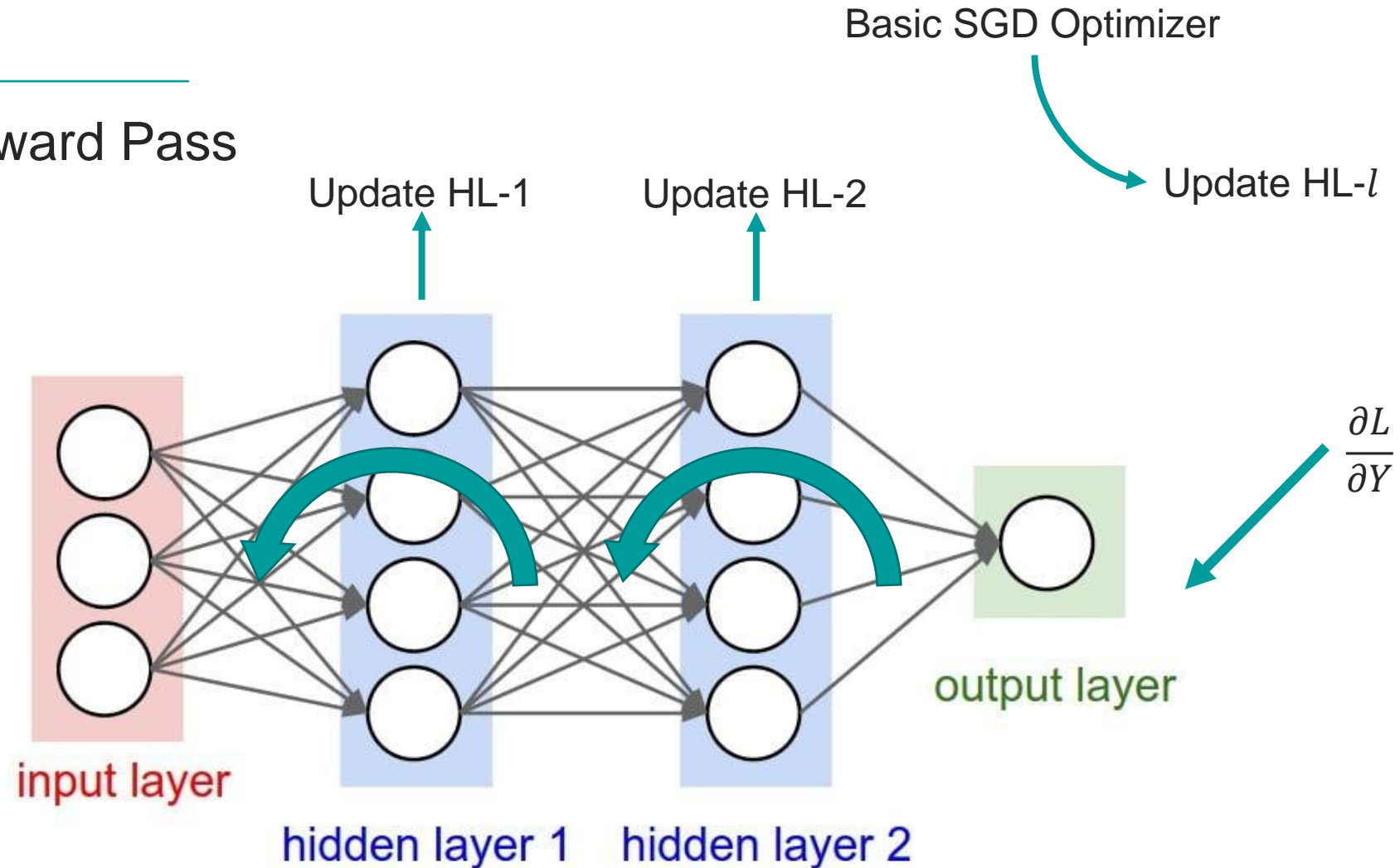


$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_{true,i} - Y_{pred,i})^2$$



BACKPROPAGATION

Backward Pass



$$\begin{cases} w_k^l = w_{k-1}^l - \eta \frac{\partial L}{\partial w^l} \\ b_k^l = b_{k-1}^l - \eta \frac{\partial L}{\partial y^l} \end{cases}$$



SAL
SILICON AUSTRIA LABS

ACTIVATION FUNCTIONS

TYPE OF ACTIVATIONS

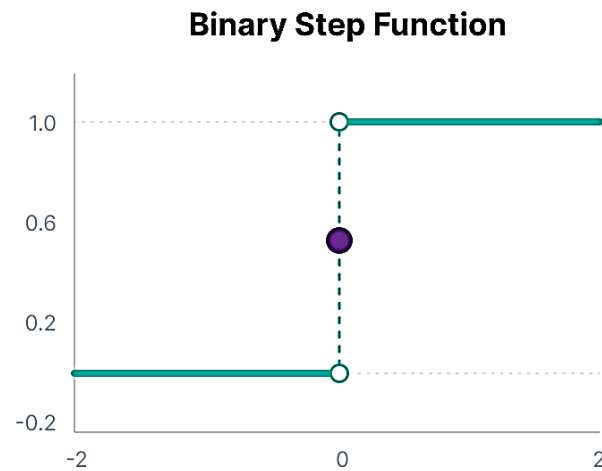
≡ AF makes the neuron's decision

≡ Types of AF:

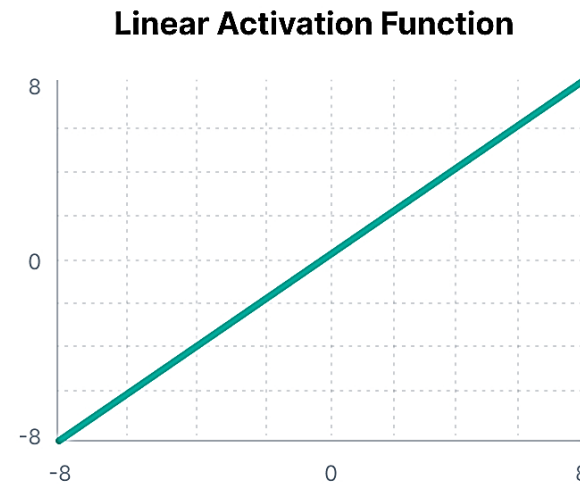
≡ Binary Step

≡ Linear

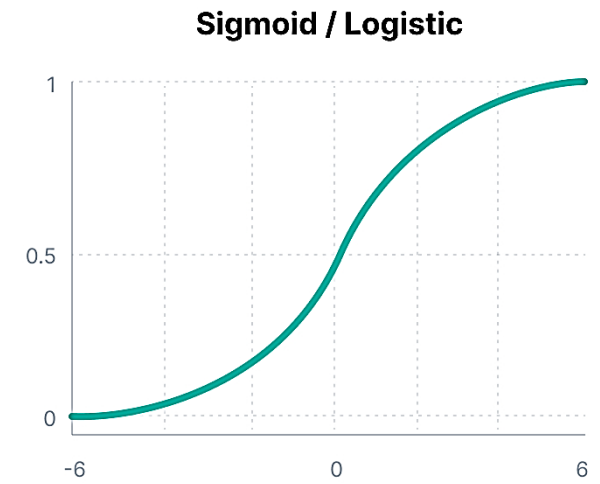
≡ Sigmoid



$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$



$$f(x) = x$$

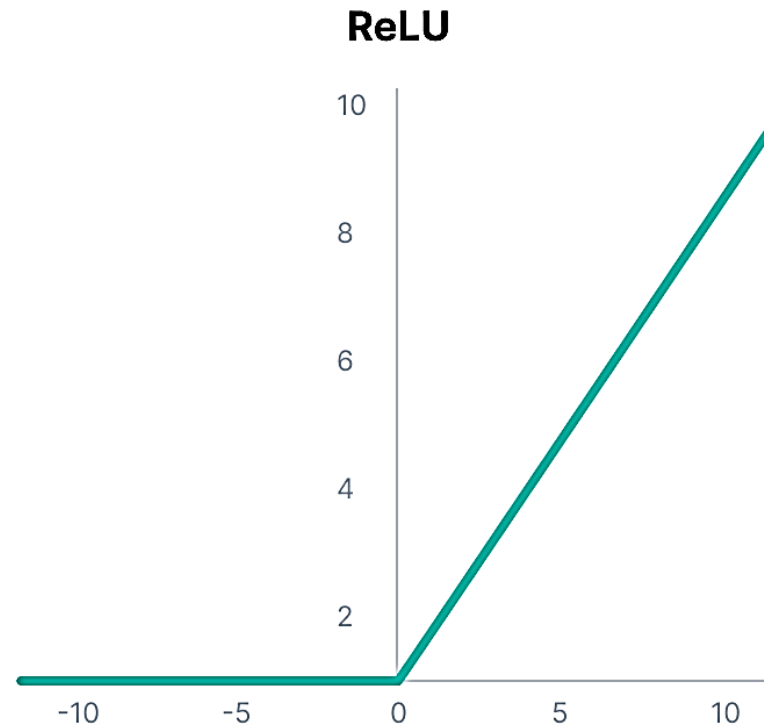


$$f(x) = \frac{1}{1 + e^{-x}}$$

RELU

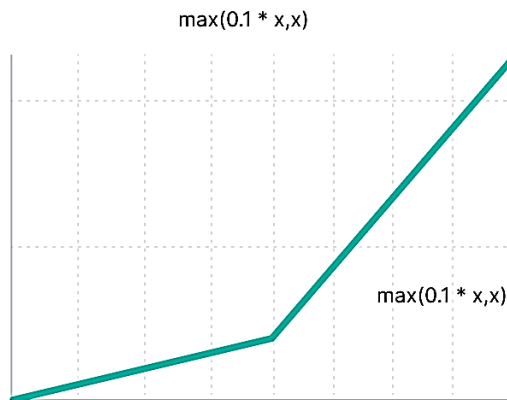
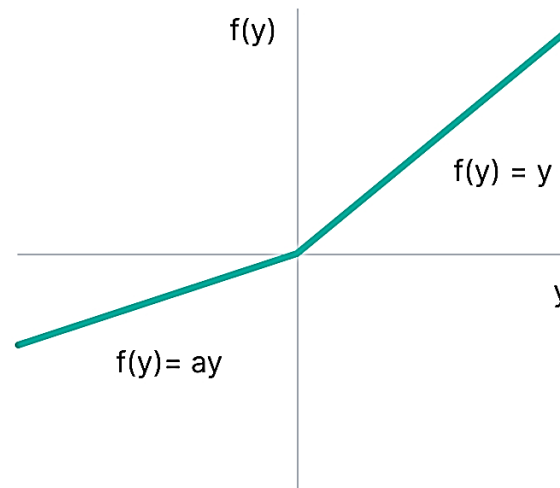
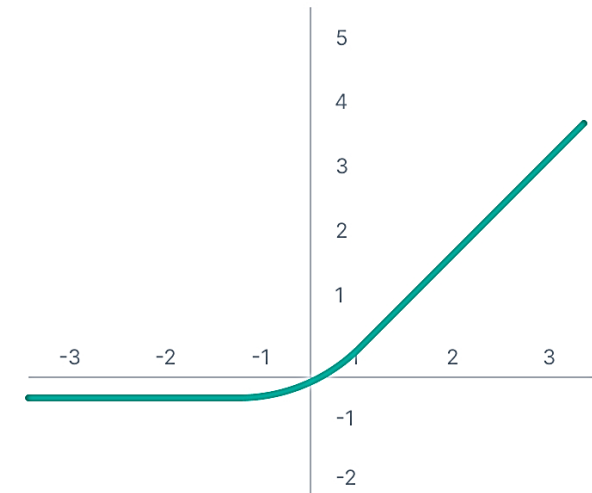
- ≡ Most popular Activation Function
- ≡ Rectified Linear Unit
- ≡ Advantages:
 - ≡ Make NN efficient as only activates some neurons
 - ≡ Accelerates Training due to its linear, non-saturating property
- ≡ Disadvantage:
 - ≡ Dying ReLU problem

$$f(x) = \max(0, x)$$



RELU

- ≡ Variants to solve Dying ReLU while keep linear and non-saturation property

Leaky ReLU**Parametric ReLU****ELU**



SAL
SILICON AUSTRIA LABS

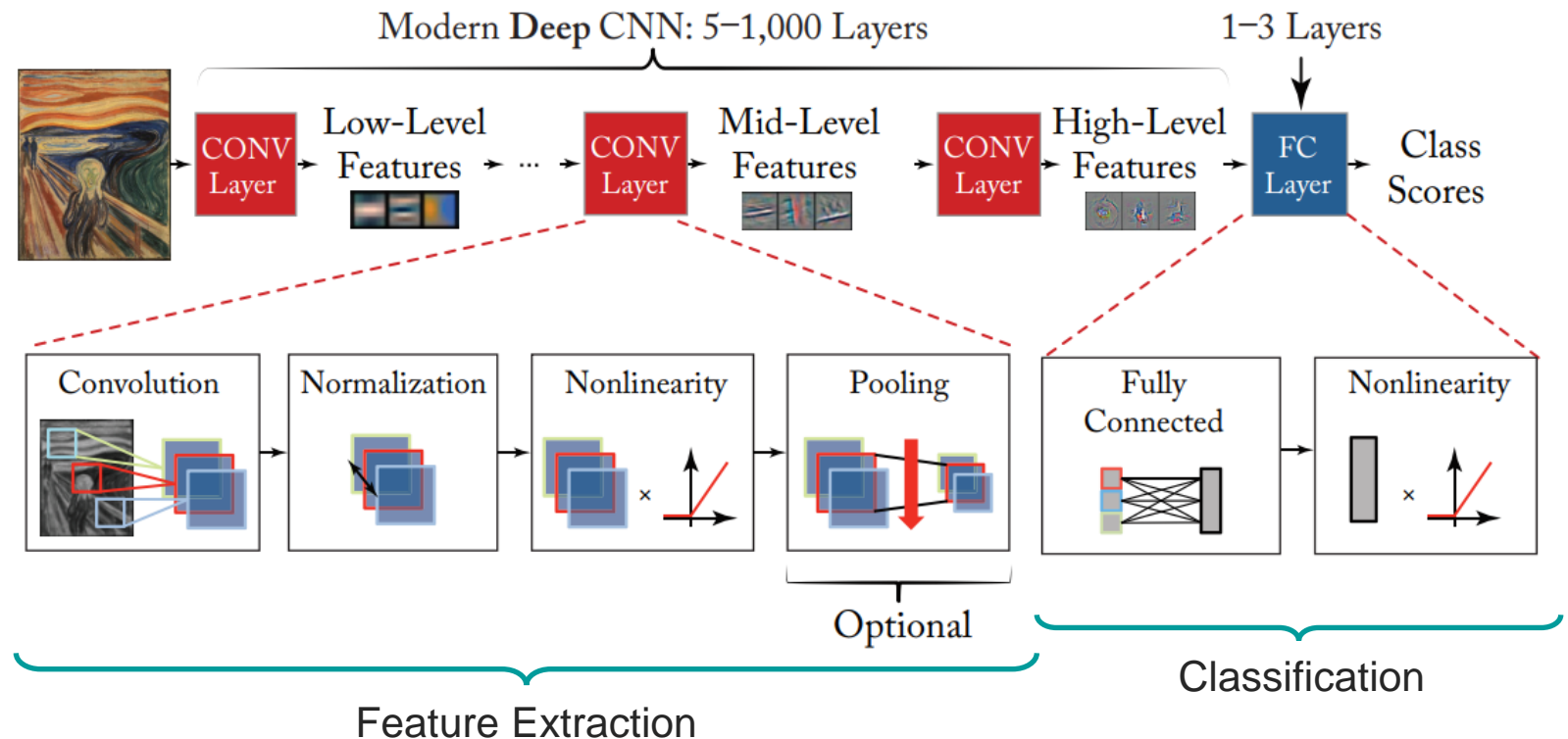
CONVOLUTIONAL NEURAL NETWORKS

WHY USE A CNN

- ≡ Problem too complex
- ≡ DNN requires many neurons
- ≡ CNN Layers:

Many parameters and operations

- ≡ Conv. Filters
- ≡ Normalization
- ≡ Down-sample
- ≡ Fully-Connected



CONV LAYER

≡ Input tensor (I_h, I_w, I_{ch})

≡ N x Kernels (K_h, K_w, K_{ch})

Kernel

1	0	1
0	1	0
1	0	1

1	2	1	0	2
2	0	0	1	0
1	0	2	1	0
0	1	0	2	1
0	2	1	0	2

Input tensor

$$\begin{aligned}
 5 &= 1 * 1 + 2 * 0 + 1 * 1 \\
 &+ 2 * 0 + 0 * 1 + 0 * 0 \\
 &+ 1 * 1 + 0 * 0 + 2 * 1
 \end{aligned}$$

5	3	6
2	6	2
5	3	7

Output tensor

CONV LAYER

≡ Input tensor (I_h, I_w, I_{ch})

≡ N x Kernels (K_h, K_w, K_{ch})

Kernel

1	0	1
0	1	0
1	0	1

1	2	1	0	2
2	0	0	1	0
1	0	2	1	0
0	1	0	2	1
0	2	1	0	2

Input tensor

5	3	6
2	6	2
5	3	7

Output tensor

CONV LAYER

≡ Input tensor (I_h, I_w, I_{ch})

≡ N x Kernels (K_h, K_w, K_{ch})

Kernel

1	0	1
0	1	0
1	0	1

1	2	1	0	2
2	0	0	1	0
1	0	2	1	0
0	1	0	2	1
0	2	1	0	2

Input tensor

5	3	6
2	6	2
5	3	7

Output tensor

CONV LAYER

≡ Input tensor (I_h, I_w, I_{ch})

≡ N x Kernels (K_h, K_w, K_{ch})

Kernel

1	0	1
0	1	0
1	0	1

1	2	1	0	2
2	0	0	1	0
1	0	2	1	0
0	1	0	2	1
0	2	1	0	2

Input tensor

5	3	6
2	6	2
5	3	7

Output tensor

CONV LAYER

≡ Input tensor (I_h, I_w, I_{ch})

≡ N x Kernels (K_h, K_w, K_{ch})

Kernel

1	0	1
0	1	0
1	0	1

1	2	1	0	2
2	0	0	1	0
1	0	2	1	0
0	1	0	2	1
0	2	1	0	2

Input tensor

5	3	6
2	6	2
5	3	7

Output tensor

CONV LAYER

≡ Input tensor (I_h, I_w, I_{ch})

≡ N x Kernels (K_h, K_w, K_{ch})

Kernel

1	0	1
0	1	0
1	0	1

1	2	1	0	2
2	0	0	1	0
1	0	2	1	0
0	1	0	2	1
0	2	1	0	2

Input tensor

5	3	6
2	0	2
5	3	7

Output tensor

CONV LAYER

$$O_{h/w} = I_{h/w} - (K_{h/w} - 1)$$

≡ Input tensor (I_h, I_w, I_{ch})

≡ N x Kernels (K_h, K_w, K_{ch})

≡ 2D Conv:

$$\equiv K_{ch} = I_{ch}$$

$$\equiv N = O_{ch}$$

Kernel

1	0	1
0	1	0
1	0	1

1	2	1	0	2
2	0	0	1	0
1	0	2	1	0
0	1	0	2	1
0	2	1	0	2

Input tensor

5	3	6
2	6	2
5	3	7

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

1	4	2	0	3
4	5	3	6	1
2	2	6	2	3
2	5	3	7	2
1	2	4	1	4

Output tensor

CONV LAYER

$$O_{h/w} = I_{h/w} - (K_{h/w} - 1) + 2 * P_{h/w}$$

≡ Padding (P_h, P_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

1	4	2	0	3
4	5	3	6	1
2	2	6	2	3
2	5	3	7	2
1	2	4	1	4

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

1	2	3
2	6	3
1	4	4

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

1	2	3
2	6	3
1	4	4

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

1	2	3
---	---	---

2	6	3
1	4	4

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

1	2	3
2	6	3
1	4	4

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

Kernel

1	0	1
0	1	0
1	0	1

$$O_{h/w} = \left\lfloor \frac{I_{h/w} - (K_{h/w} - 1) + 2 * P_{h/w} - 1}{S_{h/w}} + 1 \right\rfloor$$

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

1	2	3
2	6	3
1	4	4

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

≡ Dilation (D_h, D_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

2	2
2	2

Output tensor

CONV LAYER

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

≡ Dilation (D_h, D_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

2	2
2	2

Output tensor

CONV LAYER

$$O_{h/w} = \left\lfloor \frac{I_{h/w} - D_{h/w} * (K_{h/w} - 1) + 2 * P_{h/w} - 1}{S_{h/w}} + 1 \right\rfloor$$

≡ Padding (P_h, P_w)

≡ Stride (S_h, S_w)

≡ Dilation (D_h, D_w)

Kernel

1	0	1
0	1	0
1	0	1

0	0	0	0	0	0	0
0	1	2	1	0	2	0
0	2	0	0	1	0	0
0	1	0	2	1	0	0
0	0	1	0	2	1	0
0	0	2	1	0	2	0
0	0	0	0	0	0	0

Input tensor

2	2
2	2

Output tensor

POOL LAYERS

- ≡ Used for down-sampling
- ≡ Generally used with stride as kernel size: $S_{h/w} = K_{h/w}$
- ≡ Common Types:
 - ≡ MaxPool
 - ≡ AvgPool

Max Pooling

29	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2
pool size

100	184
12	45

Average Pooling

31	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2
pool size

36	80
12	15

BATCH NORM

- ≡ Range of layer's output differs from sample to sample
 - ≡ Original input data has samples with different distribution/range
 - ≡ Weights update after every batch
- ≡ Scales each layer's output features
 - ≡ Running feature mean: $E(x)$
 - ≡ Running feature variance: $Var(x)$
 - ≡ Trainable feature weight: γ
 - ≡ Trainable feature bias: β
- ≡ Stabilizes and speeds-up training

$$y = \frac{x - E(x)}{\sqrt{Var(x) + \varepsilon}} * \gamma + \beta$$

BATCH NORM FUSION

Conv/FC $y = b + \sum w \times x$

BatchNorm $z = \frac{y - \mu}{\sqrt{\sigma^2 + \varepsilon}} * \gamma + \beta$

$$z = \frac{(b + \sum w \times x) - \mu}{\sqrt{\sigma^2 + \varepsilon}} * \gamma + \beta$$

$$z = \underbrace{\left[\gamma \times \frac{b - \mu}{\sqrt{\sigma^2 + \varepsilon}} + \beta \right]}_{\hat{b}} + \sum \underbrace{\left[\frac{\gamma}{\sqrt{\sigma^2 + \varepsilon}} \times w \right]}_{\hat{w}} \times x$$

Constant after training

$$E(y) \rightarrow \mu$$

$$Var(y) \rightarrow \sigma^2$$

$$z = \hat{b} + \sum \hat{w} \times x$$



SAL
SILICON AUSTRIA LABS

QUESTIONS ?



SAL
SILICON AUSTRIA LABS

UNFOLD THE FUTURE