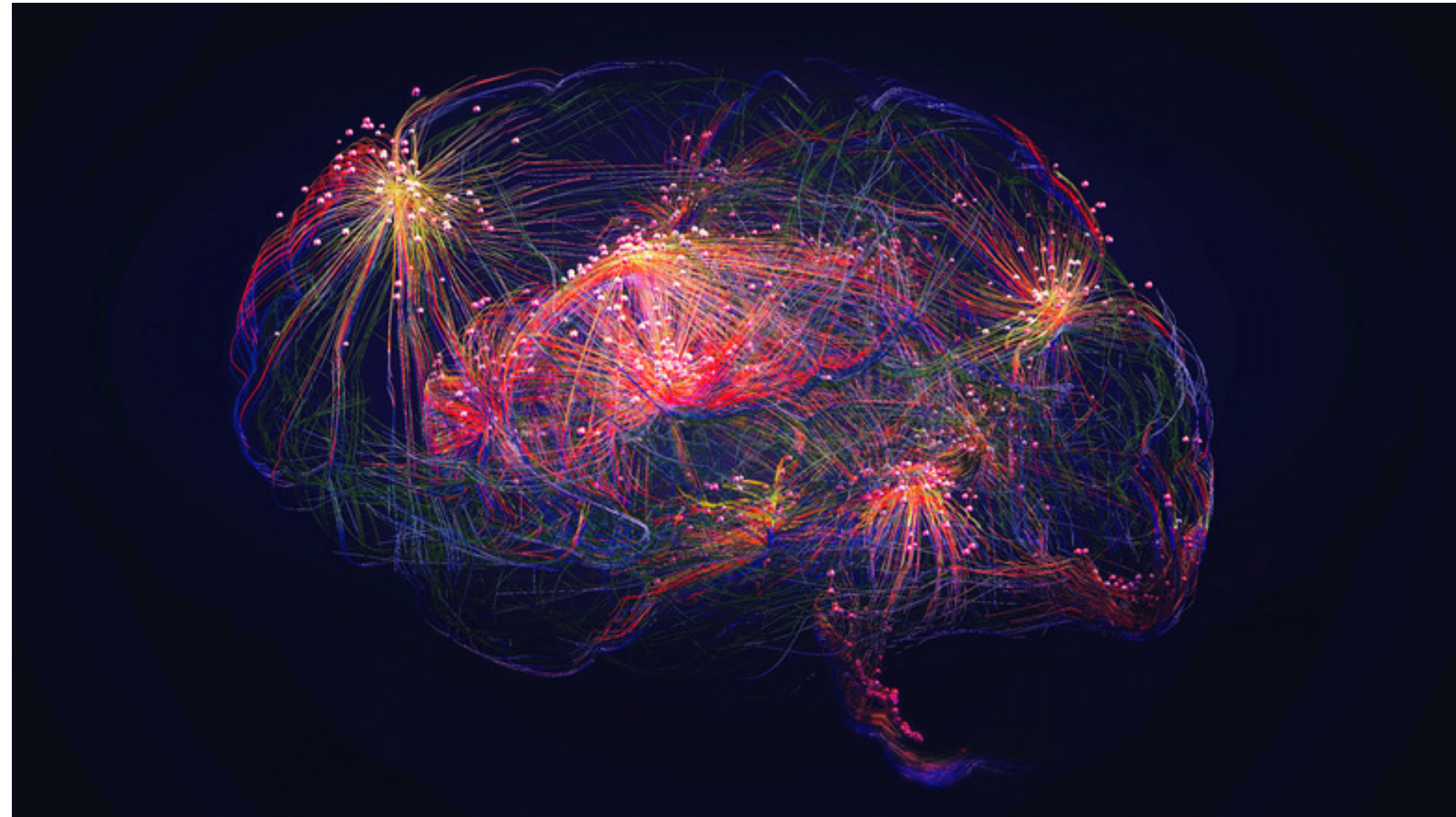# Artificial Neural Networks

Artificial Neural Networks
inspired by brain research

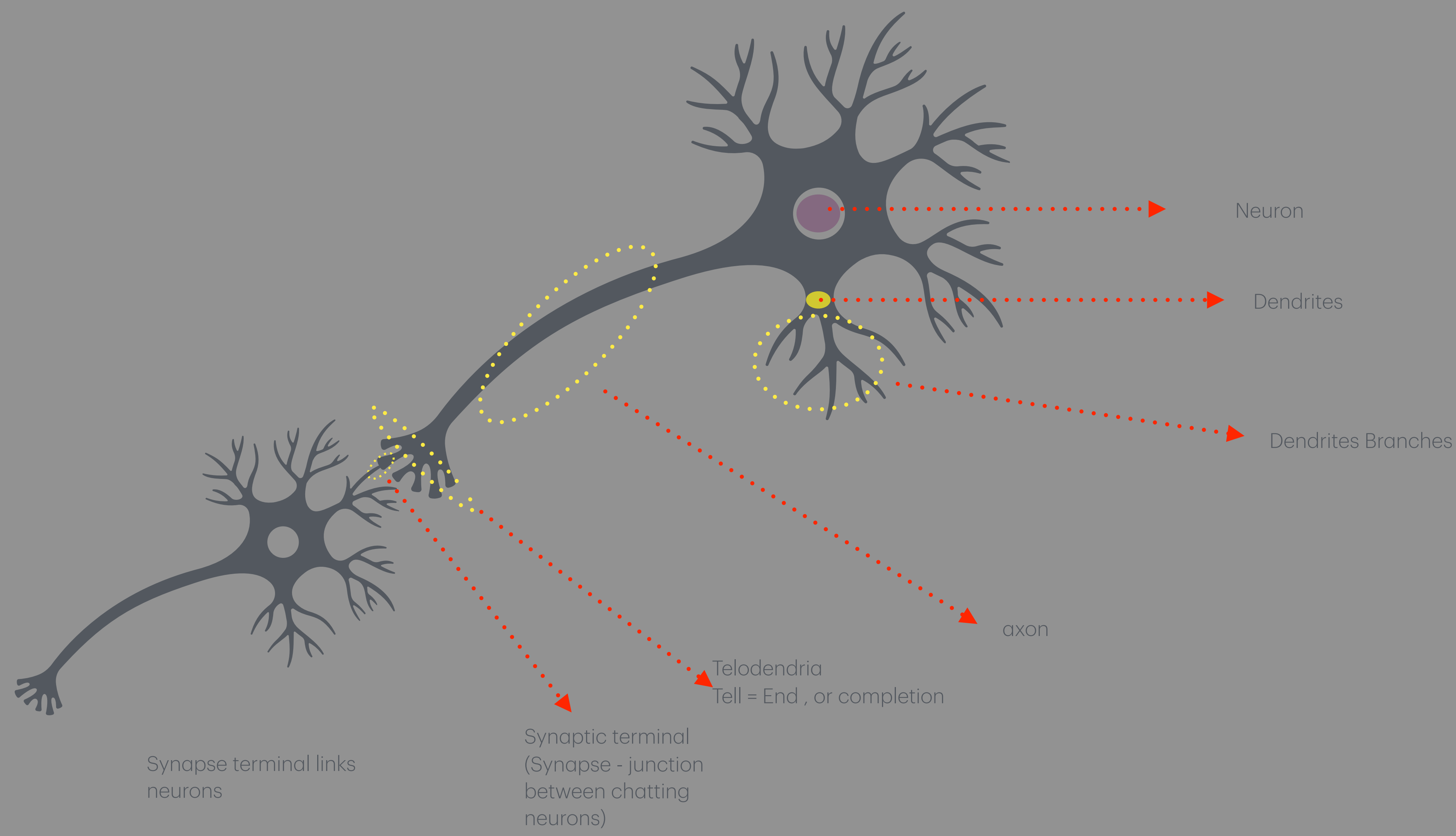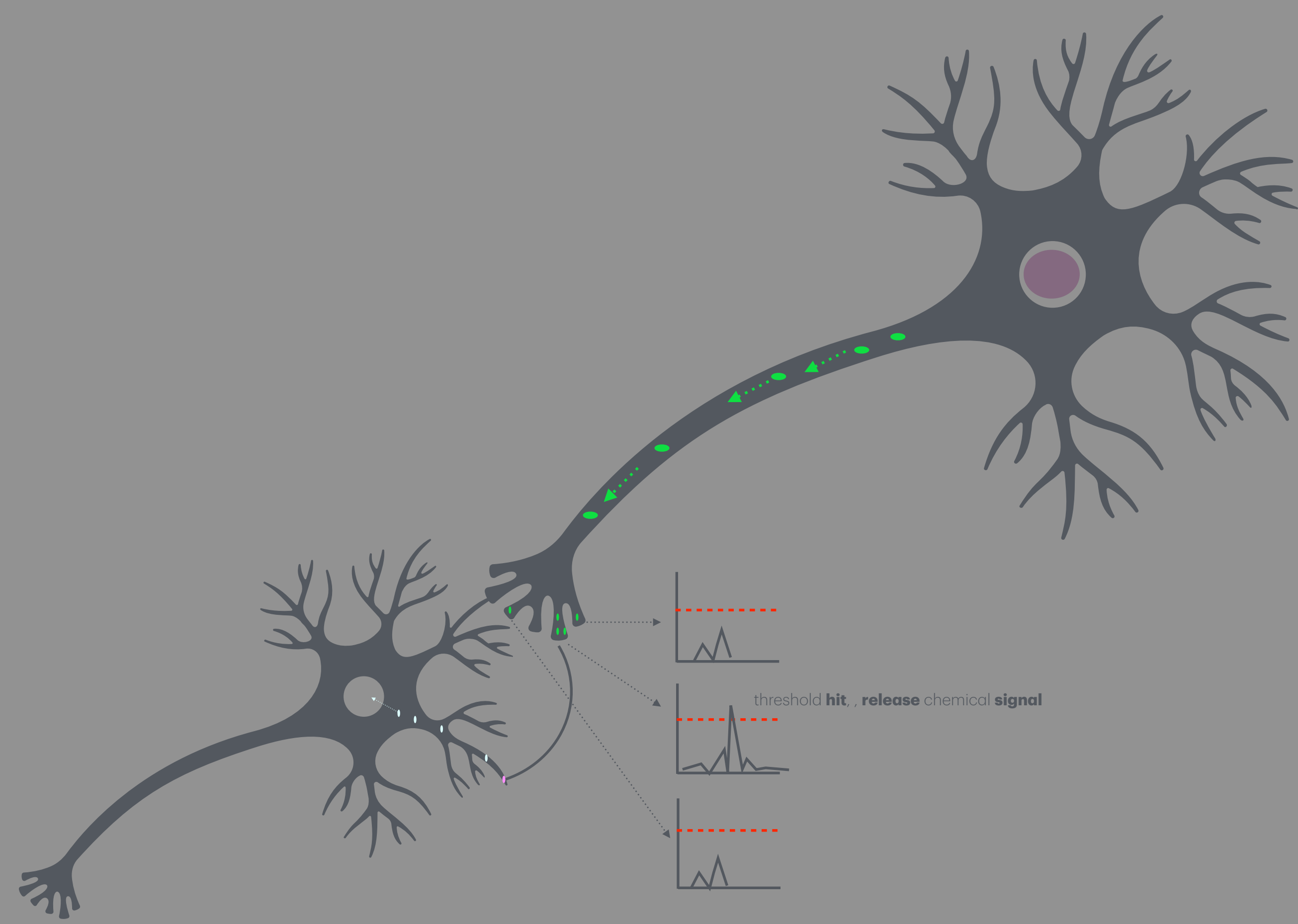| Speech Recognition | Google Images (Classification) |
| Go (DeepMind) | Recommendation Systems (Youtube) |

# Biological Neuron



Neuron

Dendrites

Dendrites Branches

axon

Telodendria
Tell = End , or completion

Synaptic terminal
(Synapse - junction
between chatting
neurons)

Synapse terminal links
neurons

# Biological Neuron



threshold **hit**, , **release** chemical **signal**

# Logical Computation: Artificial Neuron



Neuron C = Neuron A
Equal Electrical Amplitude
Note: Off neuron cancels a single active input

Neuron A **AND** Neuron B = Neuron C

Note: Neuron activated when at least two inputs are active

Note: consumer neuron C sums all the inputs
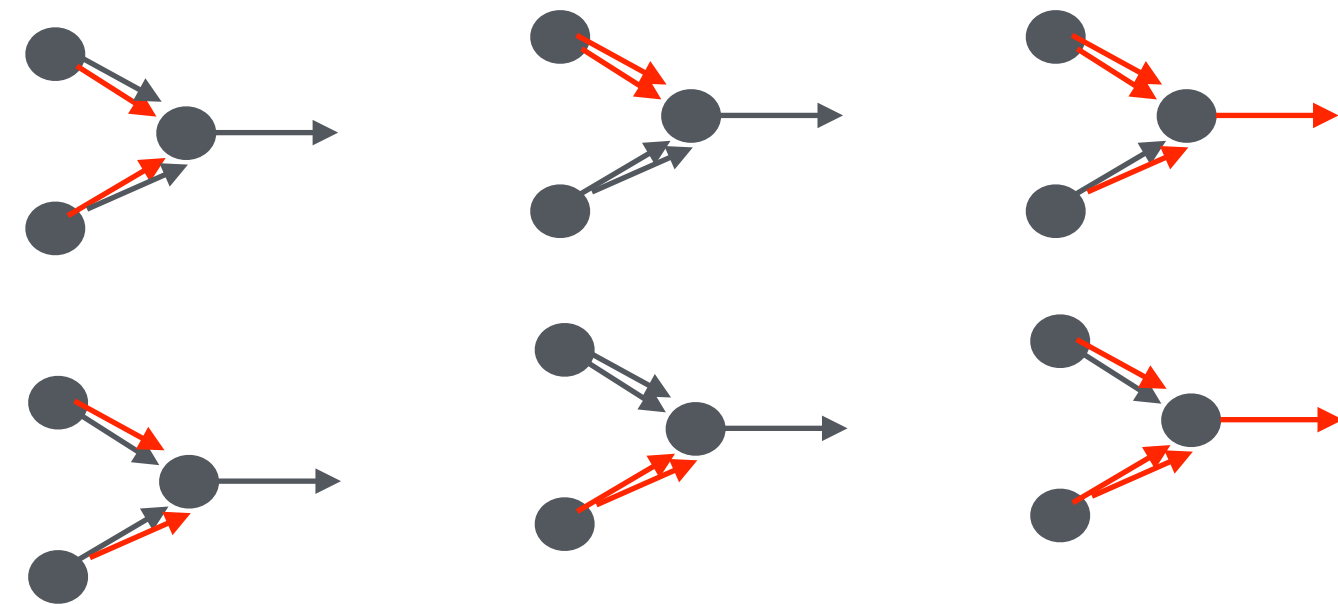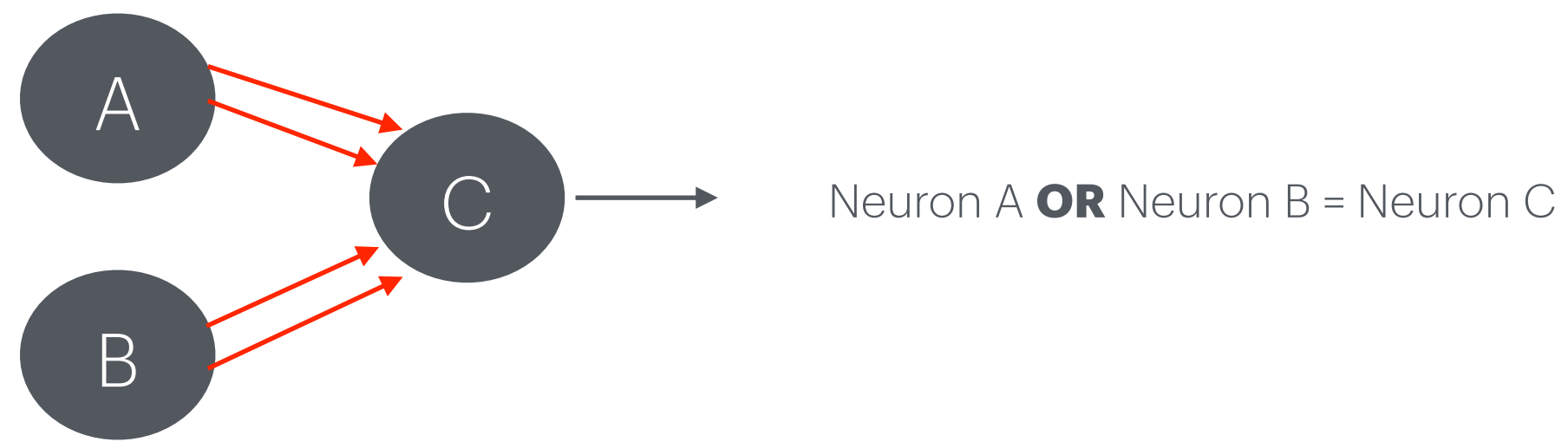
# Logical Computation: Artificial Neuron



Neuron A **OR** Neuron B = Neuron C

# Logical Computation: Artificial Neuron

# Logical Computation: Perception

**Threshold Logic Unit (TLU)**

Step(z)

$z = \mathbf{x} \cdot \mathbf{w}$

heaviside(z)

sgn(z) - sign function

$$z = \mathbf{x}^{\mathbf{T}}\mathbf{w}$$

$w_1$

$w_2$

$w_3$

$x_1$

$x_2$

$x_3$

# Logical Computation: Perception

$y$ → [ ]

$y$ **estimate**

<span style="color:red">**Threshold Logic Unit (TLU)**</span>

Step(z)

$z = \mathbf{x} \cdot \mathbf{w}$

$z = \mathbf{x^T w}$

heaviside(z)

sgn(z) - sign function

$w_0$    $w_1$    $w_2$    $w_3$

$x_0 = 1$    $x_1$    $x_2$    $x_3$

$x$    $w$    →    $x^T$    $w$    →    $z$

features

Find right values for weights that lower MSE with expectations

Used for simple linear classification (1D, 2D features )

# Logical Computation: Multi-Perception

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

artificial neuron

weights

$x_0 = 1$

$x_1$
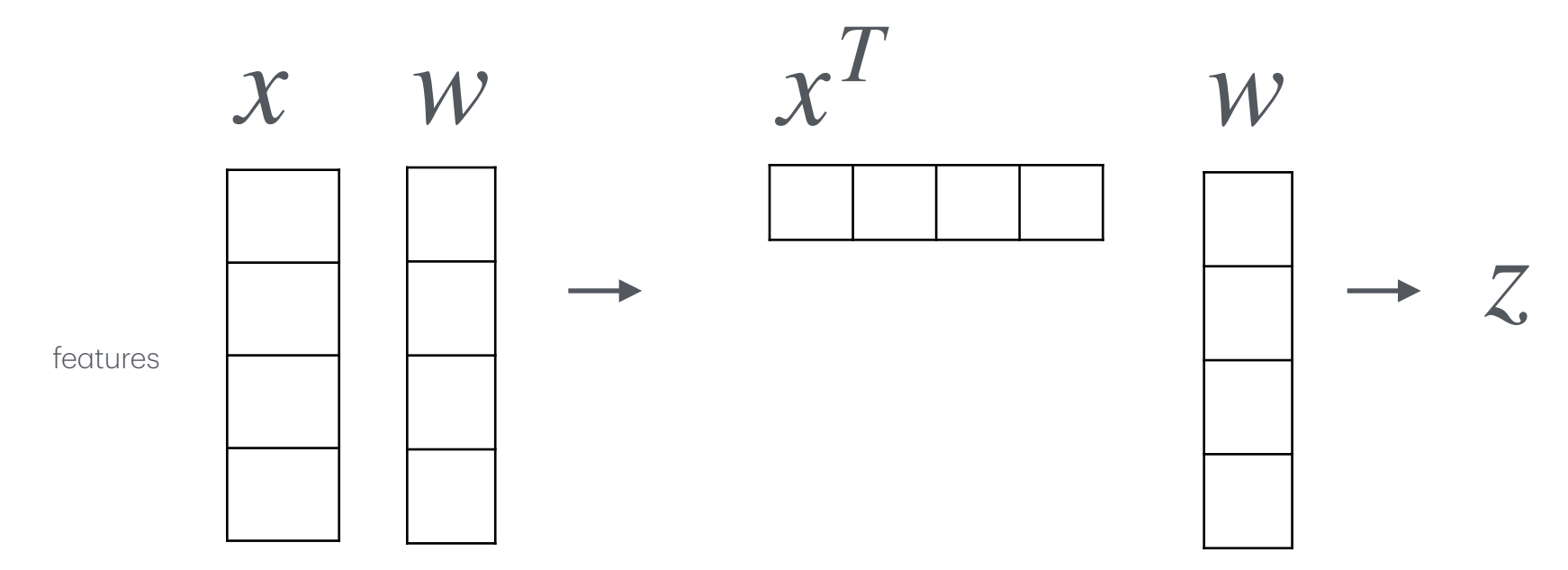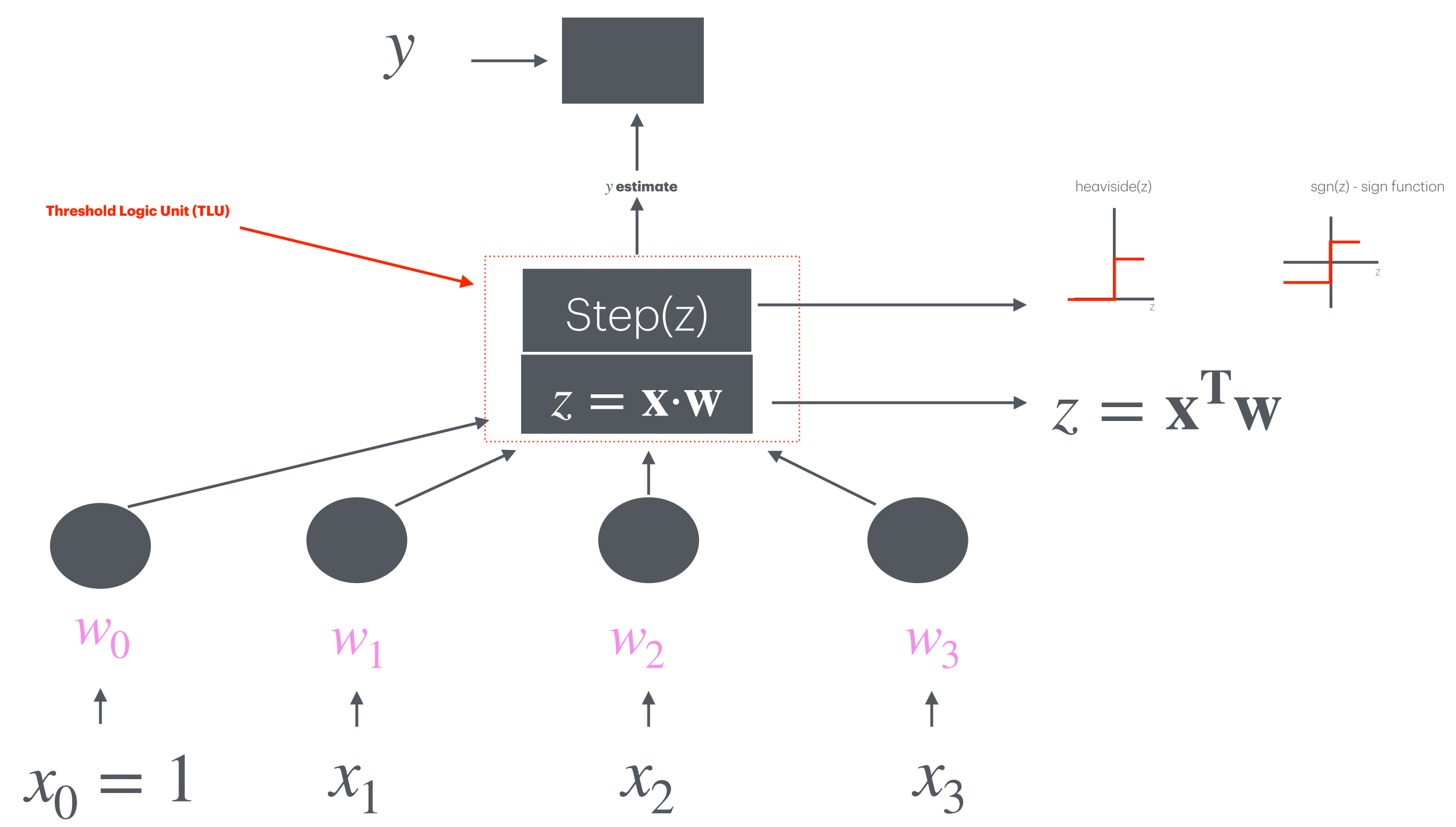
$x_2$

$x_3$

Bias neurons

Input neurons

1

1

1

1

X

W

B

# Logical Computation: Perception

# Logical Computation: Perception

# Logical Computation: Perception



Threshold Logic Unit (TLU)

$x_0 = 1$  $x_1$  $x_2$  $x_3$

X   W   B

# Logical Computation: Perception



Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

$x_0 = 1$

$x_1$

$x_2$

$x_3$

Note: Single Instance

Next layer receives single instances vector

# Logical Computation: Perception

Output Threshold Logic Unit (OTLU)

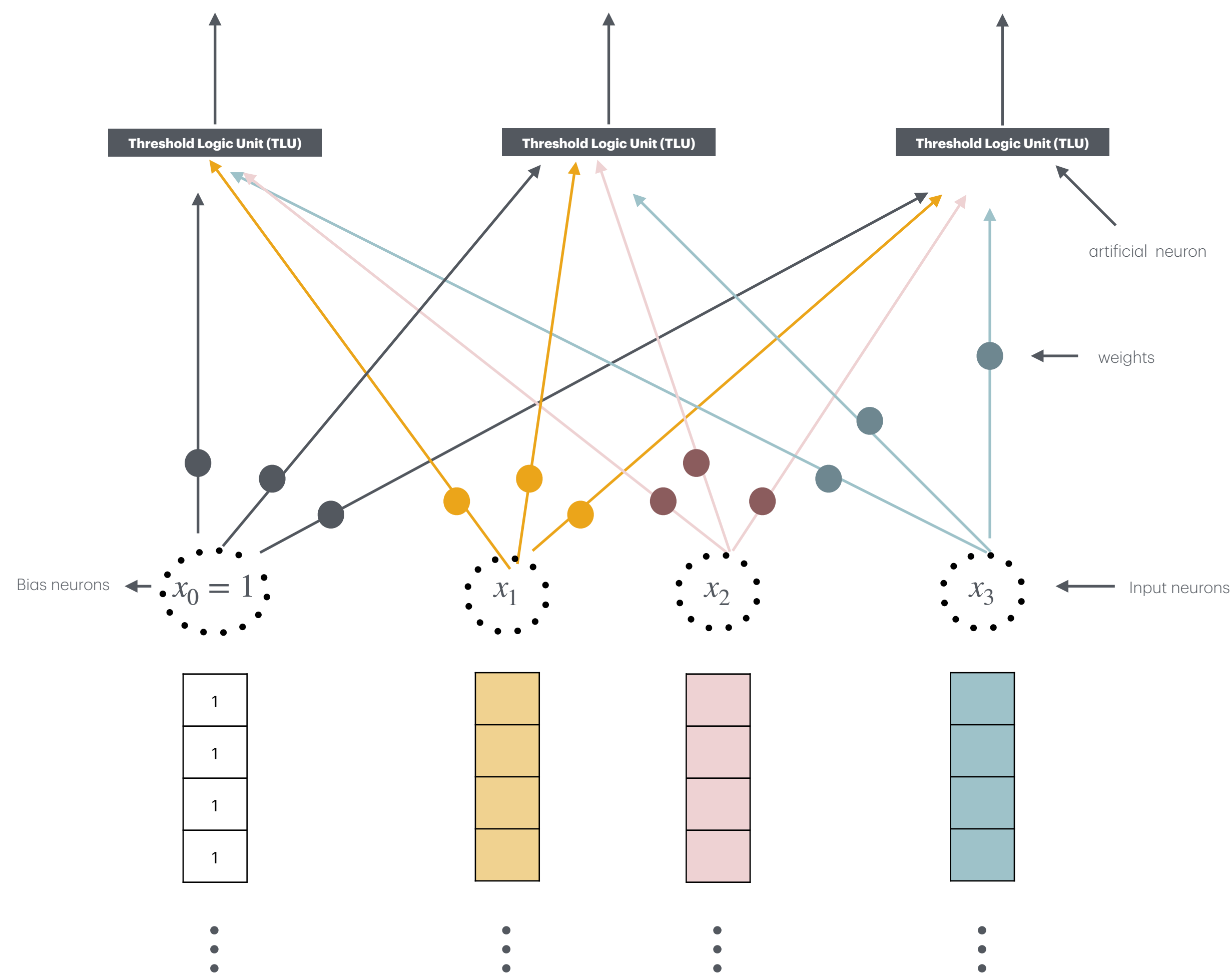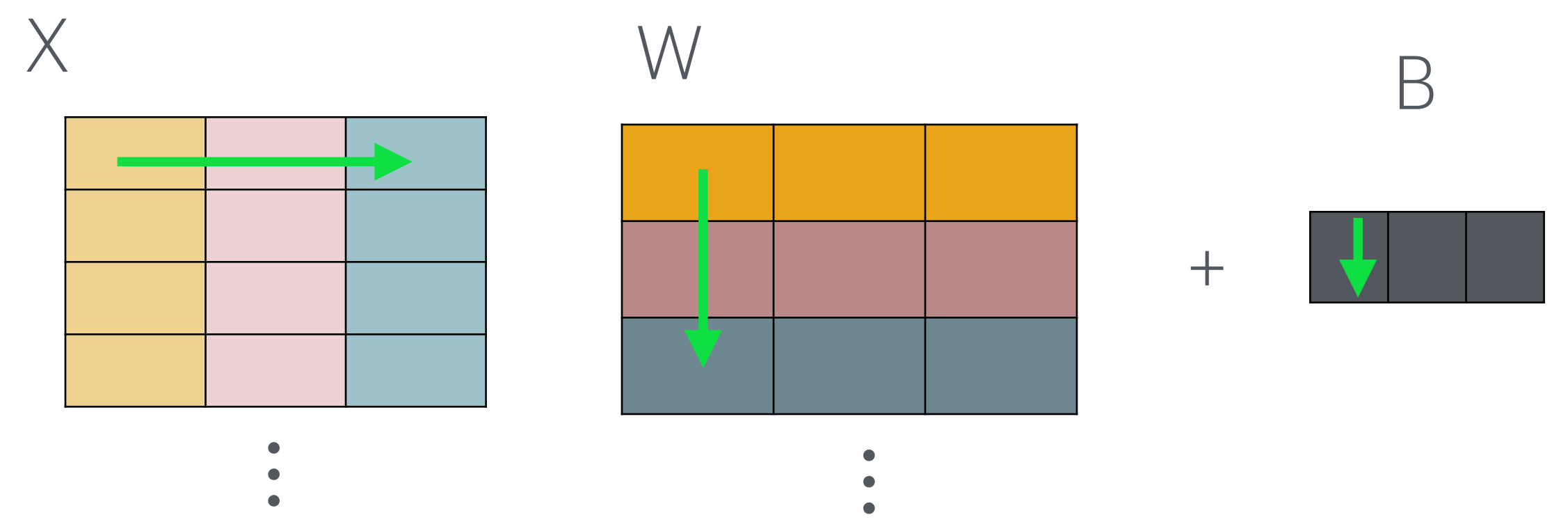OTLU usually does not have an activation unit and outputs the weighted sum

$x_0^{artifical\_neuron}$

$x_1^{artifical\_neuron}$

$x_2^{artifical\_neuron}$

$x_3^{artifical\_neuron}$

Bias =1

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

$x_0 = 1$

$x_1$

$x_2$

$x_3$

OTLU

OTLU

Output Threshold Logic Unit (OTLU)

$x_0^{artifical\_neuron}$

$x_1^{artifical\_neuron}$

$x_2^{artifical\_neuron}$

$x_3^{artifical\_neuron}$

Bias = 1

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

$x_0 = 1$

$x_1$

$x_2$

$x_3$

Networks handles mini-batch at a time

Estimation

Layer 3
Output

$y$

Cost

Output Threshold Logic Unit
(OTLU)

$\Delta$

Layer 2
Batch Process Time

$x_0^{artifical\_neuron}$

$x_1^{artifical\_neuron}$

$x_2^{artifical\_neuron}$

$x_3^{artifical\_neuron}$

Bias =1

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Forward Pass

$\Delta$

Batch Estimation Time

$\Delta$

Layer 1 (Input Layer)
Batch Process Time

$x_0 = 1$

$x_1$

$x_2$

$x_3$

$y$

Cost

Output Threshold Logic Unit
(OTLU)

All weights are
updated to varying
degrees

$x_0^{artifical\_neuron}$

$x_1^{artifical\_neuron}$

$x_2^{artifical\_neuron}$

$x_3^{artifical\_neuron}$

Bias = 1

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

$x_0 = 1$

$x_1$

$x_2$

$x_3$

Reverse Pass

How much **layer 2 connections** contributed to high cost (i.e. high error)

Cost/Error gradients are measured across connections (weights)

$$\frac{Cost}{W_{some\_connection}}$$

Gradient Descent performed on all connections (weights) using error gradients

Note: Input batch persistence is required for reverse algorithm

$y$

Cost

Output Threshold Logic Unit
(OTLU)

$x_0^{artifical\_neuron}$

$x_1^{artifical\_neuron}$

$x_2^{artifical\_neuron}$

$x_3^{artifical\_neuron}$

Bias = 1

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

All weights are
updated to varying
degrees

$x_0 = 1$

$x_1$

$x_2$

$x_3$

Reverse Pass

How much **layer 1 connections** contributed to high cost (i.e. high error)

Cost/Error gradients are measured across connections (weights)

$$\frac{Cost}{W_{some\_connection}}$$

Gradient Descent performed on all connections (weights) using error gradients

Note: Input batch persistence is required for reverse algorithm

$y$

Cost

Output Threshold Logic Unit
(OTLU)

$x_0^{artifical\_neuron}$

$x_1^{artifical\_neuron}$

$x_2^{artifical\_neuron}$

$x_3^{artifical\_neuron}$

Bias =1

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

Threshold Logic Unit (TLU)

$x_0 = 1$

$x_1$

$x_2$

$x_3$

Reverse Pass

**Stop** after performing weight update using Gradient Descent on all the connections in each layer

: Activation Function



Activation Function:

Linear Regression/Classifiers :

- **Heaviside**

- **Sign Function**

Nonlinear Regression/Classifiers :

- **Sigmoid Function**

- **Hyperbolic Tangent Function**

- **Rectified Linear Unit Function**

Non-Linear activation functions can
be used on linearly models as well

# Activation Function

**Biological neurons** have been observed to implement a roughly **sigmoid activation** function



$\sigma(z)$

# Forward-Mode AutoDifferentiation

$g(x, y) = 5 + xy$

# Forward-Mode AutoDifferentiation

$$g(x, y) = 5 + xy$$

Partial Derivative $\dfrac{g(x, y)}{\partial x}$

Symbolic Differentiation (created from AutoDiff)

$Y$

$+$

$0$

$Y$

$$\frac{\partial uv}{\partial x} = \frac{\partial v}{\partial x}u + v\frac{\partial u}{\partial x}$$

$+$

$\times$

$\times$

$0$    $X$      $Y$    $1$

$\dfrac{\partial y}{\partial x}$    $u$      $v$    $\dfrac{\partial x}{\partial x}$

AutoDiff Computation Graphs

$+$

$5$

$u = 5$

$\times$     $\cdots\cdots$    $\dfrac{\partial uv}{\partial x} = \dfrac{\partial v}{\partial x}u + v\dfrac{\partial u}{\partial x}$

$X$    $Y$

U=X     V=y

$+$    $\cdots\cdots\cdots\cdots$   $\dfrac{\partial(u + v)}{\partial x} = \dfrac{\partial u}{\partial x} + \dfrac{\partial v}{\partial x}$

$Y$

$5$

$u = 5$

$\times$

$X$    $Y$

U=X     V=y

1

Forward
AutoDiff

Forward
AutoDiff

$g(x, y) = 5 + xy$

1

Y

Y

+

5

$u = 5$

×

X

Y

If $\epsilon$ is a infinitesimal number with $\epsilon^2 = 0$, dual numbers can be used to solve forward-mode autodiff

$$\frac{\partial f(3,4)}{\partial x} = ?$$

$$f(3 + \epsilon, 4) = f(3,4) + f'(3,4)\epsilon$$

$$f(3 + \epsilon, 4) = f(3,4) + \frac{\partial f(3,4)}{dx}\epsilon$$

Rule: $h(a + \epsilon) = h(a) + h'(a)\epsilon$

value at point          derivative at point

Run forward diff to find real and $\epsilon$ components

$$f(x, y) = x^2y + y + 2$$

42+24$\epsilon$

$$f(x, y) = 42 \quad \text{real component}$$

$$f'(x, y) = 24 \quad \epsilon \text{ component}$$

Note: $\epsilon^2 = 0$



Partial derivative with respect to y requires the same process

3 + $\epsilon$          3 + $\epsilon$          4

# Reverse-Mode AutoDifferentiation

$$f(x, y) = x^2 y + y + 2$$

$$\frac{\partial f(3,4)}{\partial x} = \,?$$

42

$+$

36

$\times$

6

$+$

9

$\times$

| X |
|---|

3

| Y |
|---|

4

| 2 |
|---|

$n_7$ change caused by $n_5$

$f(x, y)$ change caused by $n_7$

Output $n_x$ at each node. $\dfrac{\partial f}{\partial n_7} = 1$

$$\frac{\partial f}{\partial n_5} = \frac{\partial f}{\partial n_7} \times \frac{\partial n_7}{\partial n_5}$$

$\dfrac{\partial f}{\partial n_7}$

$n_7$

$+$

$$\frac{\partial f}{\partial n_6} = \frac{\partial f}{\partial n_7} \times \frac{\partial n_7}{\partial n_6}$$

$f(x, y)$ change caused by $n_5$

36   $n_5$

$\times$

$n_6$   6

$+$

$\dfrac{\partial f}{\partial n_4}$   $n_4$

$\dfrac{\partial f}{\partial n_3}$

$\dfrac{\partial f}{\partial n_3}$

$\dfrac{\partial f}{\partial n_2}$

$\dfrac{\partial f}{\partial n_1}$

$n_1$

$\times$

$n_3$

| Y |
|---|

$n_2$

| 2 |
|---|

| X |
|---|

3

4

# Contrived Example

$$\frac{\partial cost(y_1)}{\partial x_3} = \partial cost(y_1 + \epsilon)$$

$$\frac{\partial cost(y_1)}{\partial x_3} = \partial cost(y_1 + \epsilon) = cost(y_1) + cost'(y_1)\epsilon$$

$$\frac{\partial cost(y_1)}{\partial x_3} = \partial cost(y_1 + \epsilon) = cost(y_1) + \frac{\partial cost(y_1)}{\partial x_3}\epsilon$$

Contrived Example...



$$\frac{\partial cost}{\partial y_1}$$

$y_1$

How?

$$\frac{\partial cost}{\partial w_{31}} = ?$$

$w_{b,1}$    $w_{1,1}$    $w_{2,1}$    $w_{3,1}$

$x_{bias} = 1$    $x_1$    $x_2$    $x_3$

# Contrived Example More...



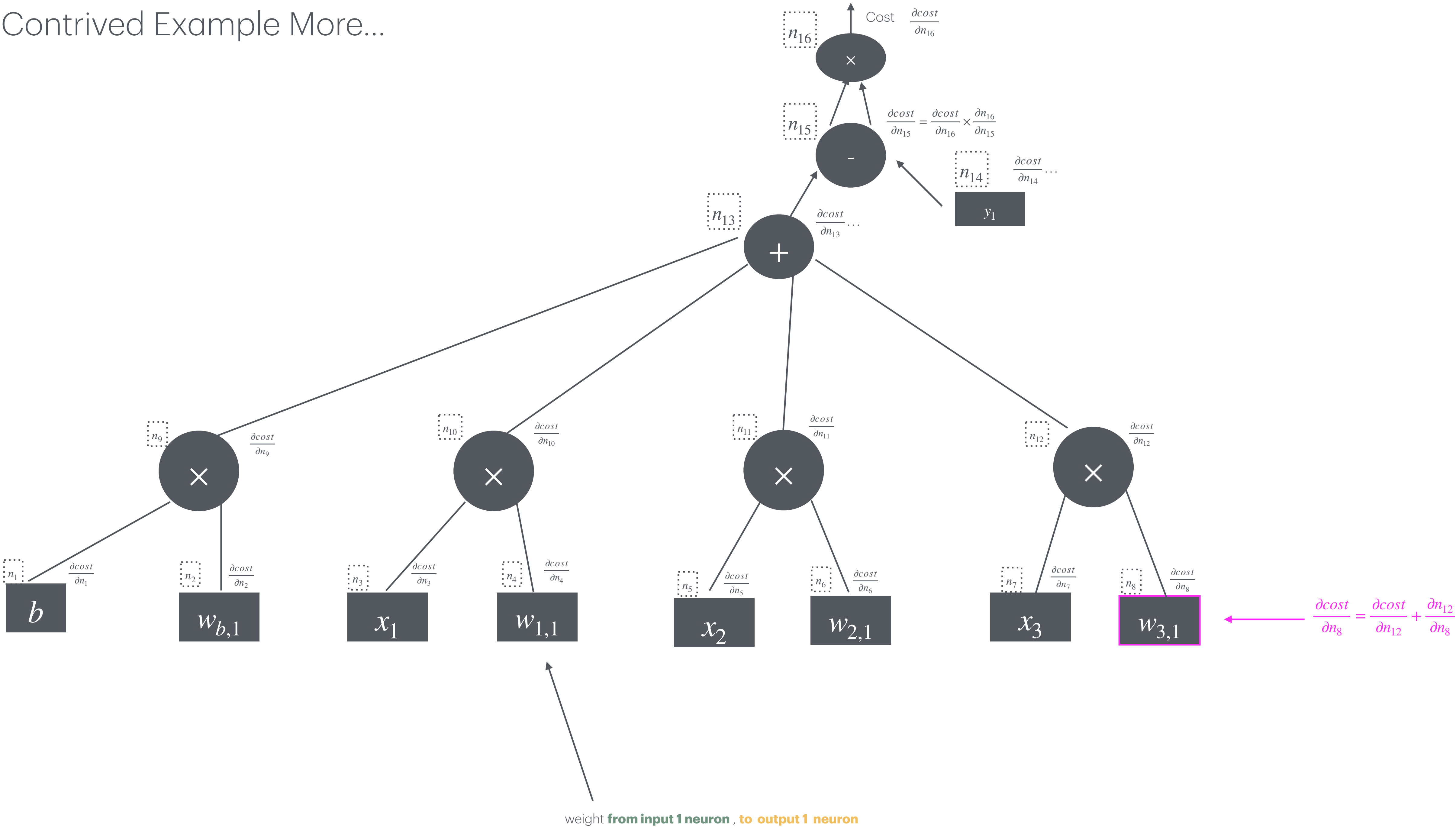Cost $\frac{\partial cost}{\partial n_{16}}$

$n_{16}$   $\times$

$n_{15}$   -   $\frac{\partial cost}{\partial n_{15}} = \frac{\partial cost}{\partial n_{16}} \times \frac{\partial n_{16}}{\partial n_{15}}$

$n_{14}$   $\frac{\partial cost}{\partial n_{14}} \cdots$

$y_1$

$n_{13}$   +   $\frac{\partial cost}{\partial n_{13}} \cdots$

$n_9$   $\times$   $\frac{\partial cost}{\partial n_9}$

$n_{10}$   $\times$   $\frac{\partial cost}{\partial n_{10}}$

$n_{11}$   $\times$   $\frac{\partial cost}{\partial n_{11}}$

$n_{12}$   $\times$   $\frac{\partial cost}{\partial n_{12}}$

$n_1$   $\frac{\partial cost}{\partial n_1}$

$b$

$n_2$   $\frac{\partial cost}{\partial n_2}$

$w_{b,1}$

$n_3$   $\frac{\partial cost}{\partial n_3}$

$x_1$

$n_4$   $\frac{\partial cost}{\partial n_4}$

$w_{1,1}$

$n_5$   $\frac{\partial cost}{\partial n_5}$

$x_2$

$n_6$   $\frac{\partial cost}{\partial n_6}$

$w_{2,1}$

$n_7$   $\frac{\partial cost}{\partial n_7}$

$x_3$

$n_8$   $\frac{\partial cost}{\partial n_8}$

$w_{3,1}$

$\frac{\partial cost}{\partial n_8} = \frac{\partial cost}{\partial n_{12}} + \frac{\partial n_{12}}{\partial n_8}$

weight **from input 1 neuron** , **to output 1 neuron**

# Tensorflow

## Keras(API)

Keras(implementation )

## Libraries containing Keras

- ★ Tensorflow
- Microsoft Cognitive Toolkit
- Theano
- ★ Keras(API)
- ★ PyTorch

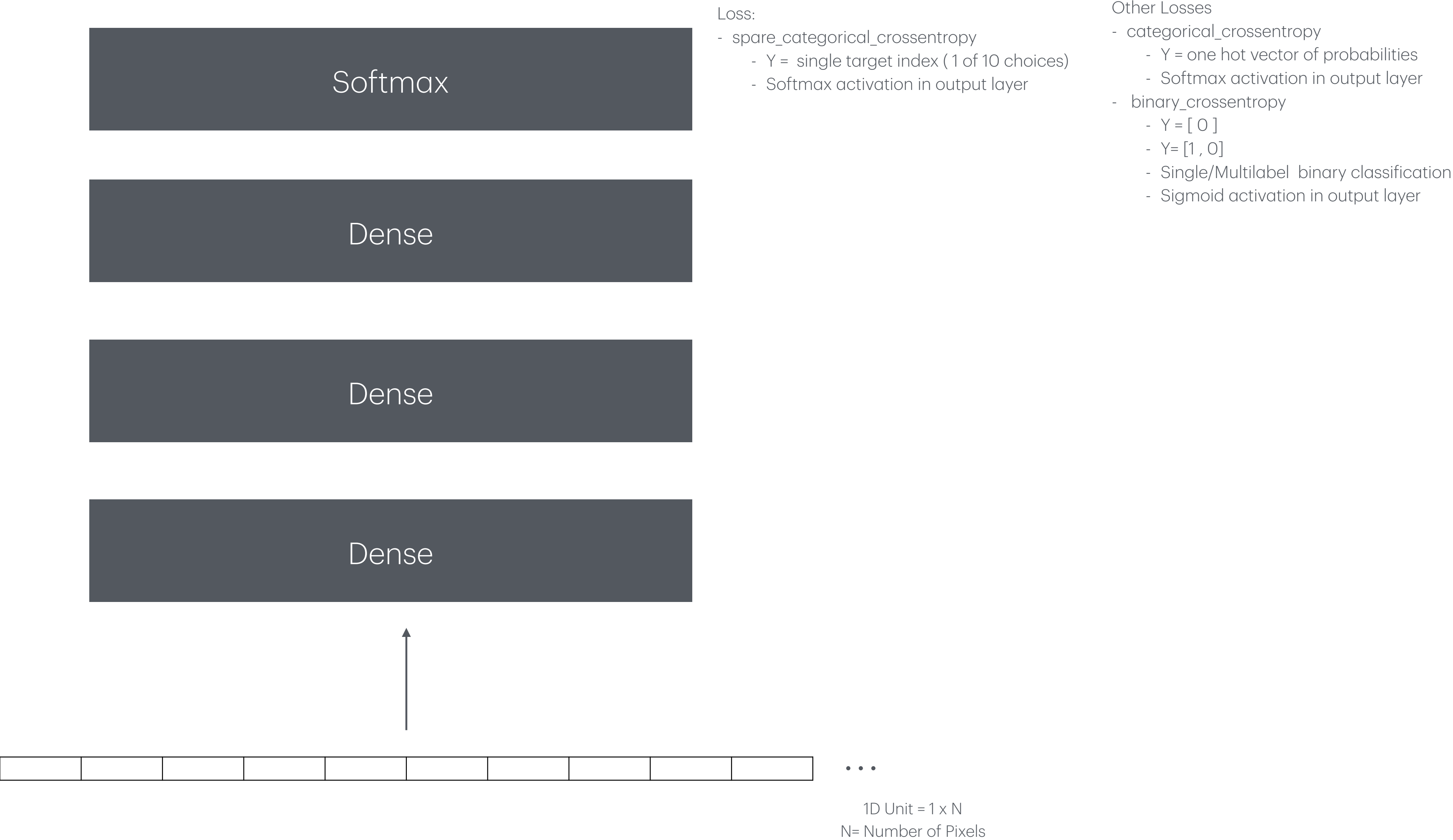$\star - popular$

Tensorflow

keras          Tensor flow Features

## Others containing Keras

- Javascript/Typescript
- PlaidML
- Apple's Core ML
- Apache MXNet

# Sequential Model: Classify Fashion MNIST

Softmax

Dense

Dense

Dense

Loss:
- spare_categorical_crossentropy
    - Y = single target index ( 1 of 10 choices)
    - Softmax activation in output layer

Other Losses
- categorical_crossentropy
    - Y = one hot vector of probabilities
    - Softmax activation in output layer
- binary_crossentropy
    - Y = [ 0 ]
    - Y= [1 , 0]
    - Single/Multilabel  binary classification
    - Sigmoid activation in output layer

• • •

1D Unit = 1 x N
N= Number of Pixels

# Dealing with skewed data

| Class A | Class B | Class C | Class D | |
|---------|---------|---------|---------| |
| | | | | Class weights |

class weight example:
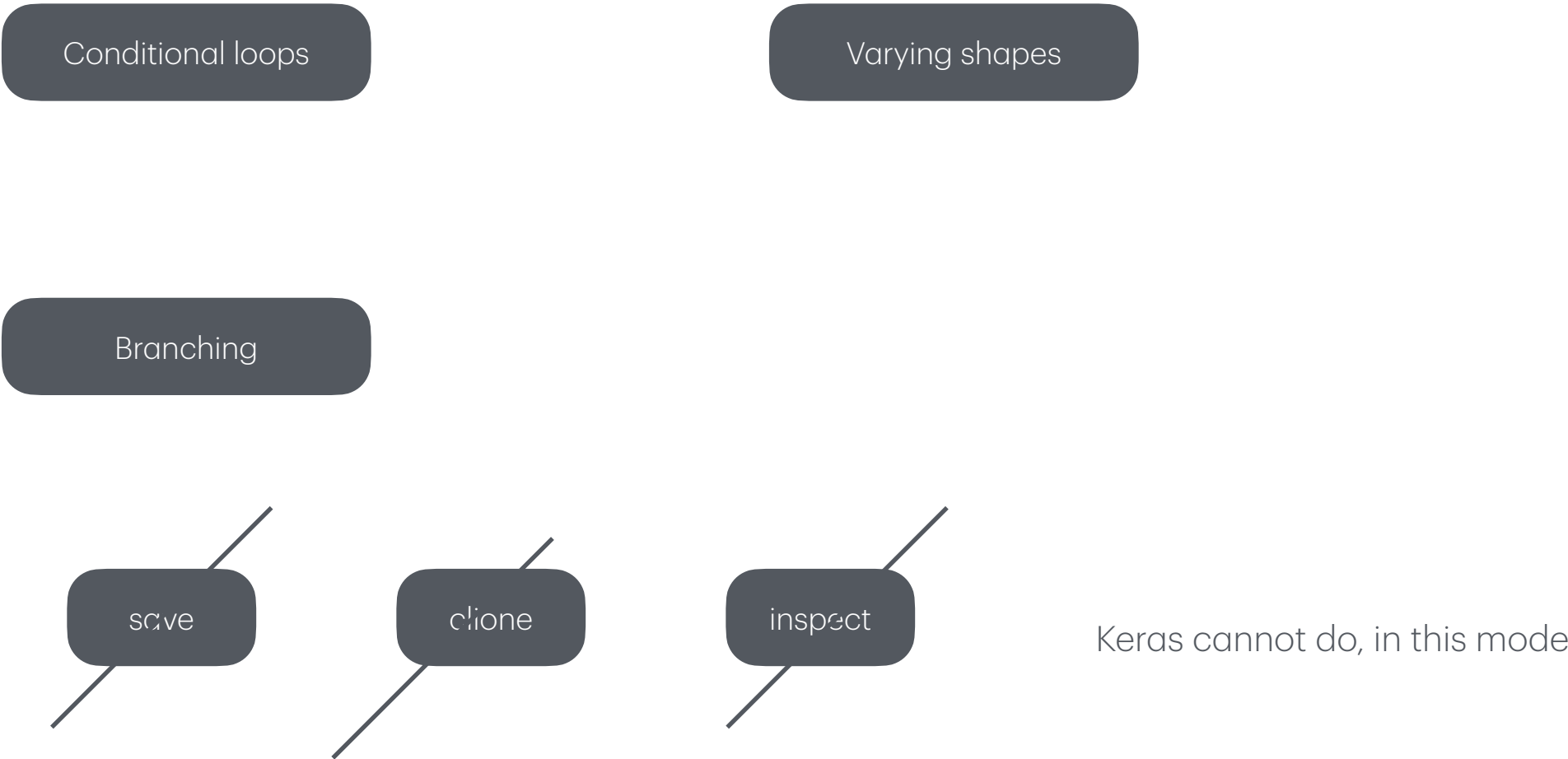Class A **overrepresented** in dataset. Give **more weight** to Class B,C, and D

| |
|--|
| Instance 1 |
| Instance 2 |
| Instance 3 |
| Instance 4 |

⋮

Samples weights

sample weight example:
Instances labeled by **expert** and **crowdsourcing**
More weight towards the **expert** instances

# Saving

- Functional

- Sequential

- Subclassing

Declarative    Static Graph

Conditional loops    Varying shapes

Branching

save    clione    inspect

save    clione    inspect    Keras cannot do, in this mode

*Save and load model weights yourself*

# Saving

| | | | |
|---|---|---|---|
| **Architecture** | **Compile**<br>- loss<br>- search/optimizer method (i.e. SGD) | **Fit** | **Save** |

Export Artifacts- TF Serving            Save SOME_NAME.keras file

# Tensorboard

**Program**

**Root Directory**

Model → Train →

**Logs(events) File**



↔ Tensorboard Server

# Fine-Tuning Neural Networks

keras_Reg

Predict          Score          Higher the better

Note: Scikit-Learn wants scores not losses

Fit          Passed to keras model -->Add training data, validation, callbacks, etc ...

Keras Regressor          Wraps model into scikit-learn regressor (e.g. scikit-learn classifier used to train or fit clustered data

Build Model          Build DNN model and compile          build_model (Parameters )

# Fine-Tuning Neural Networks

Predict

Higher the better

Score

Note: Scikit-Learn wants scores not losses

Fit

Explore **hyper-parameters** and find optimal *build_model*
*- Search uses K-fold cross-validation*

Randomized Search

Hyperparameters

GridSearch could be used instead.

Keras Regressor

Wraps model into scikit-learn regressor (e.g. scikit-learn classifier used to train or fit clustered data

Build Model

build_model (Parameters )

**Networks learns in hierarchical way**

faces

squares, circles

Line segments, shapes, orientations

Complex problems deep networks
Have higher parameter efficiency.
Fewer neurons needed per layer

Shallow network can solve many
problems with enough neurons
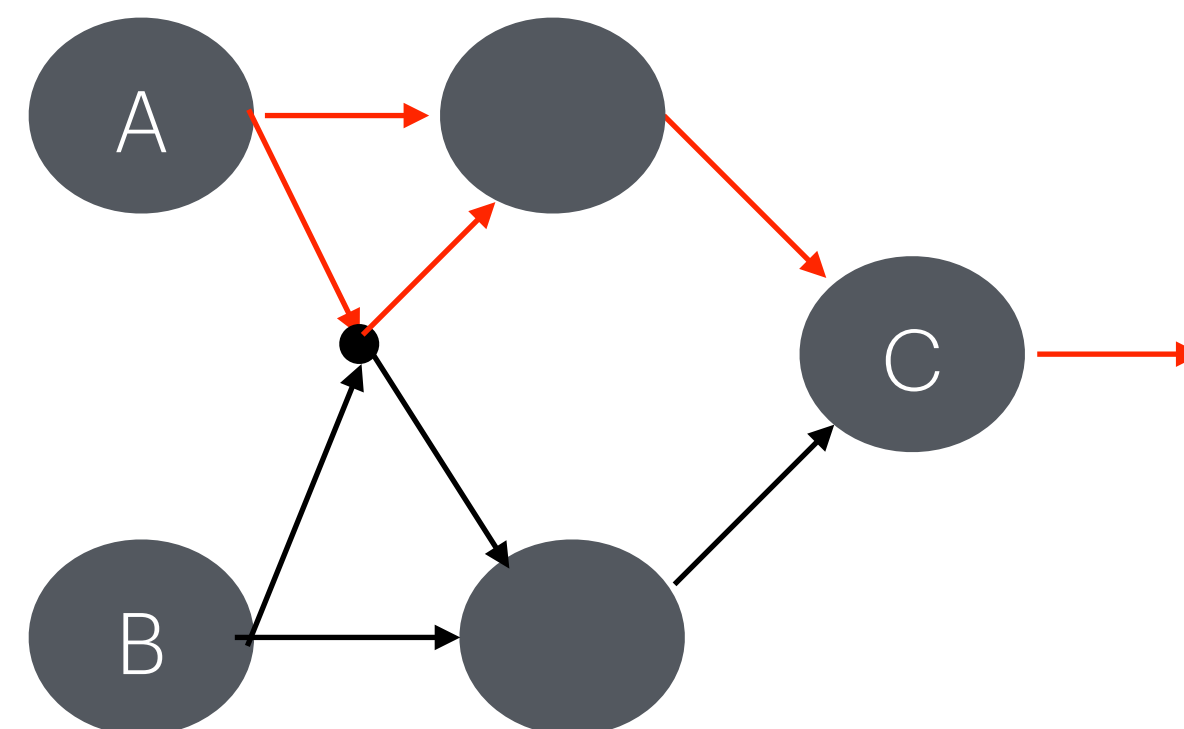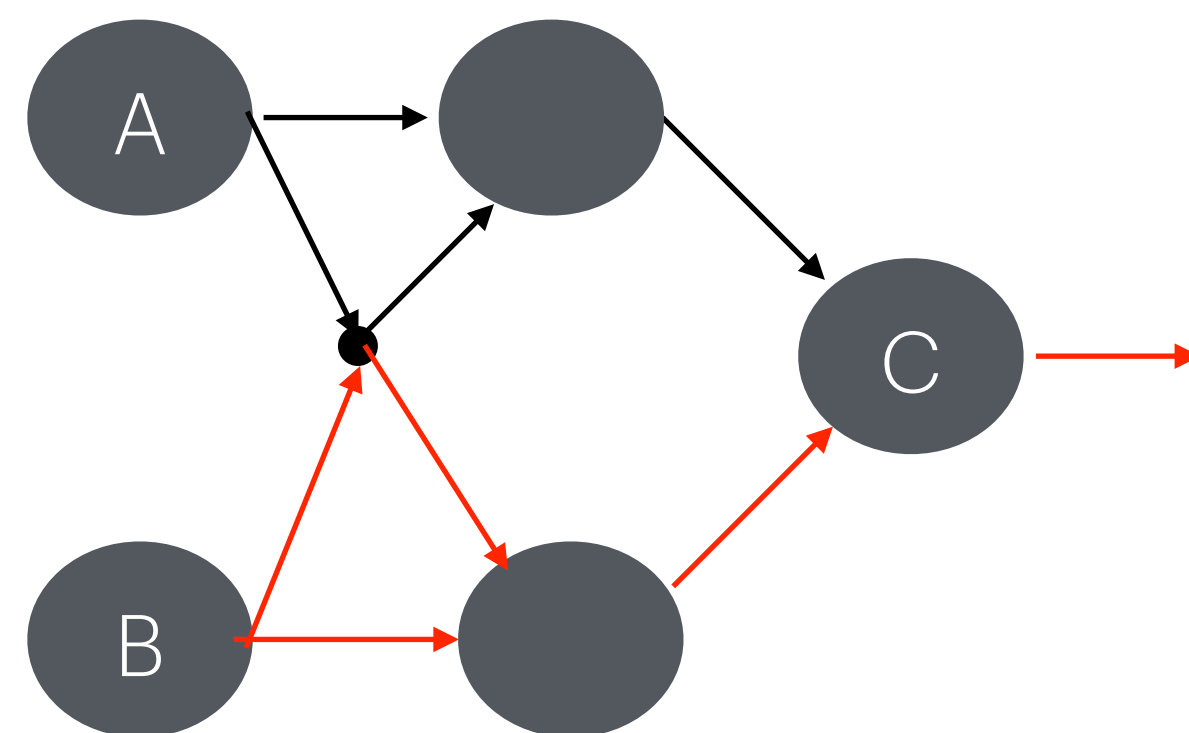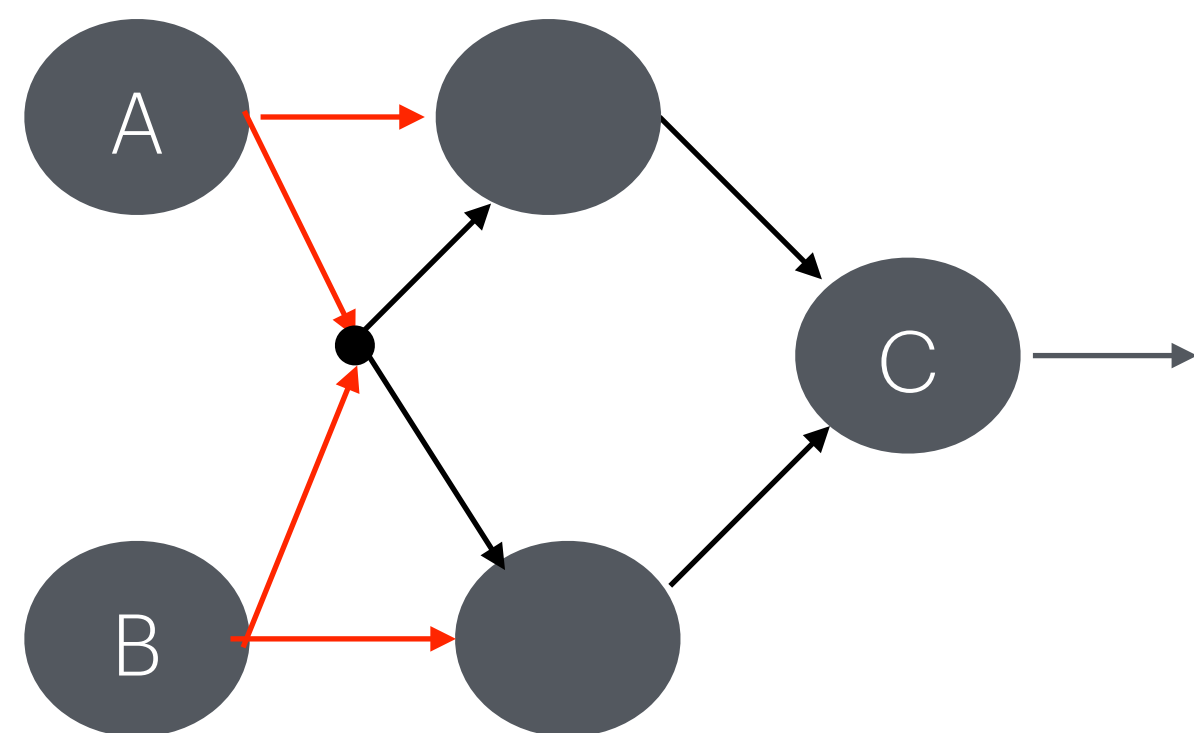
**RELU faster, notice linear boundaries**

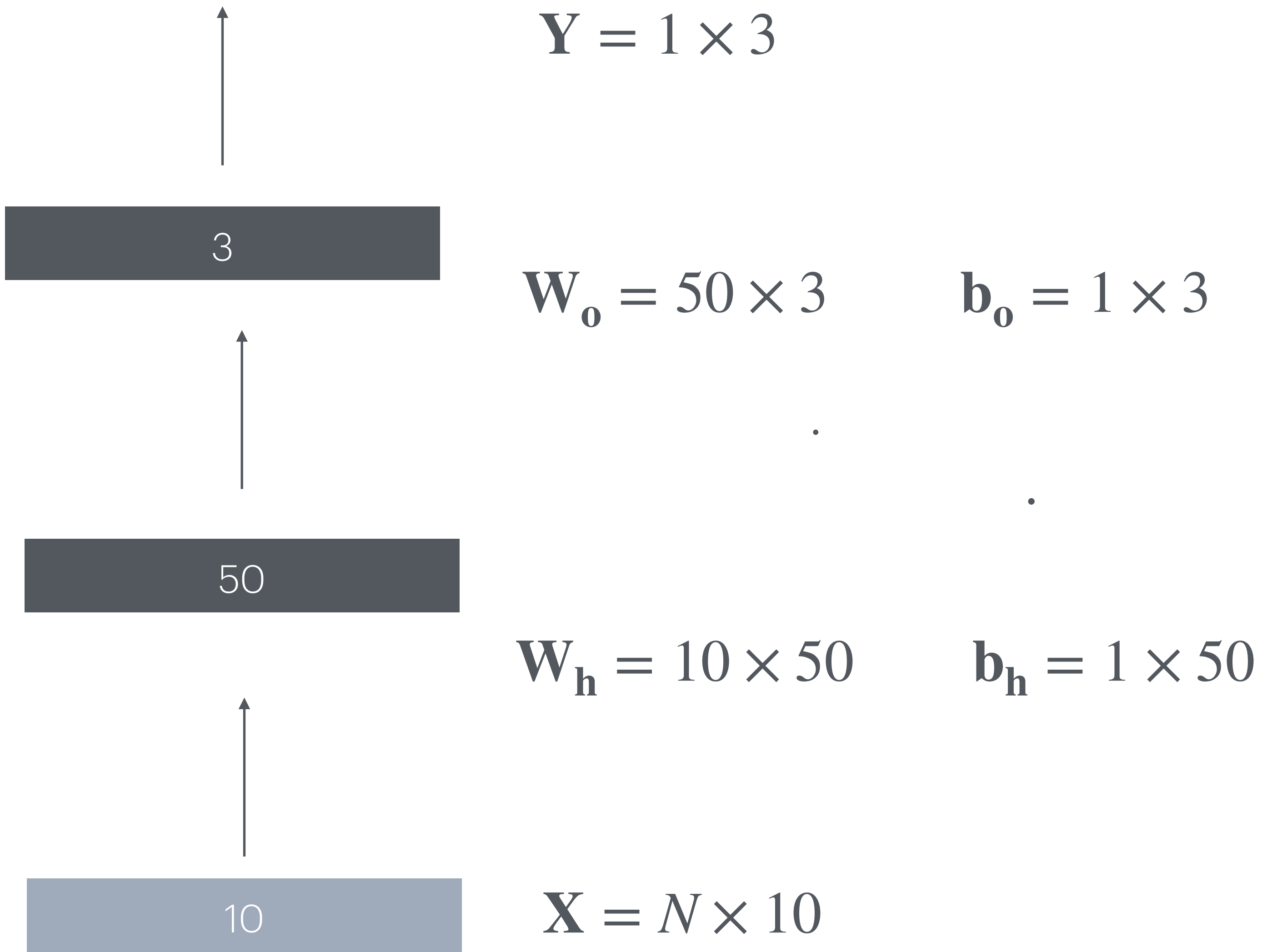**TANU takes time to converge on a solution**

# Exercise 2

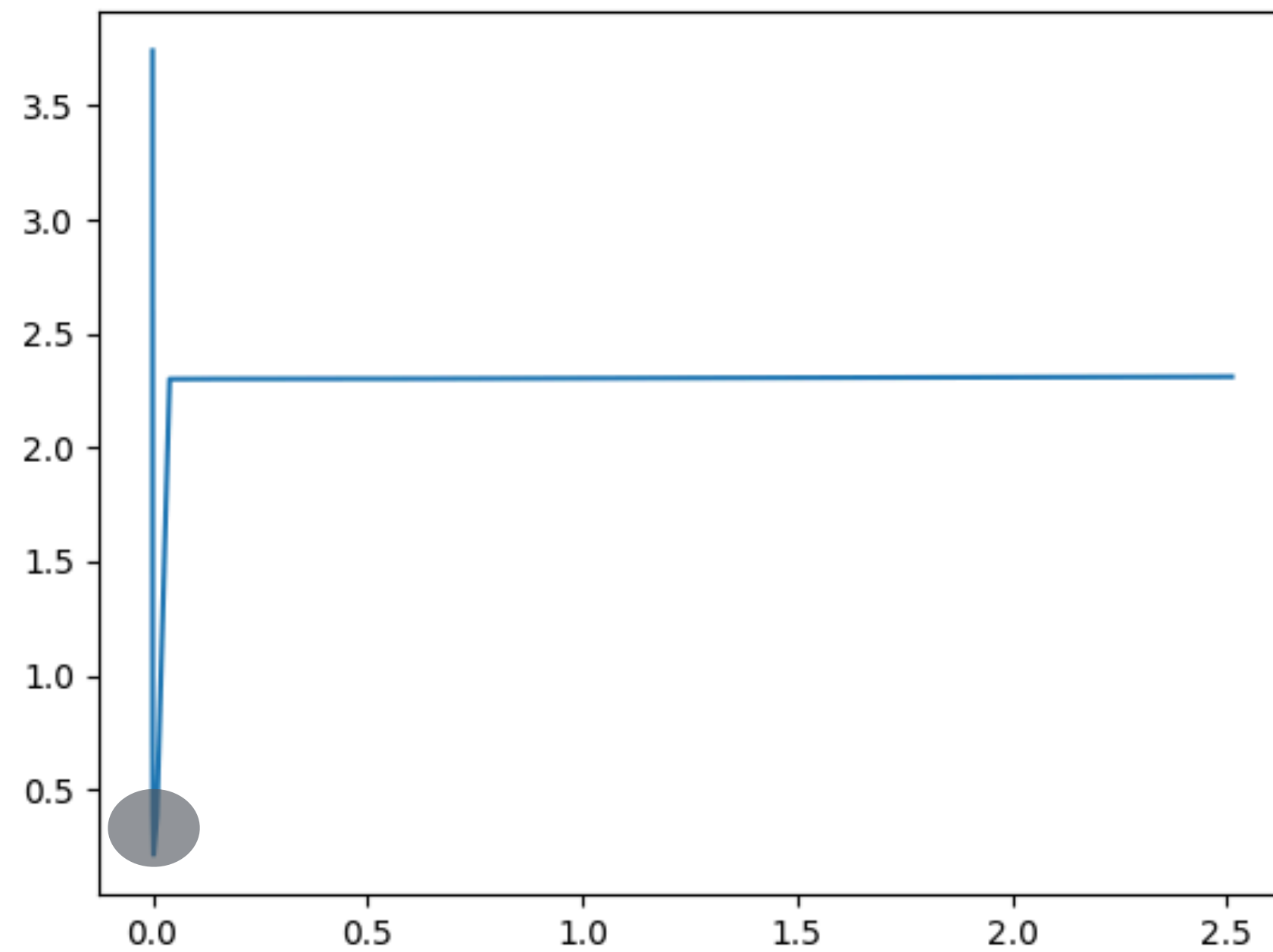Draw an ANN using the original artificial neurons that computes $A \bigoplus B$

$$(A \bigcup \neg B) \bigcup (\neg A \bigcup B)$$

Exercise 3

$$\mathbf{Y} = 1 \times 3$$

3

$$\mathbf{W_o} = 50 \times 3 \qquad \mathbf{b_o} = 1 \times 3$$

.

.

50

$$\mathbf{W_h} = 10 \times 50 \qquad \mathbf{b_h} = 1 \times 50$$
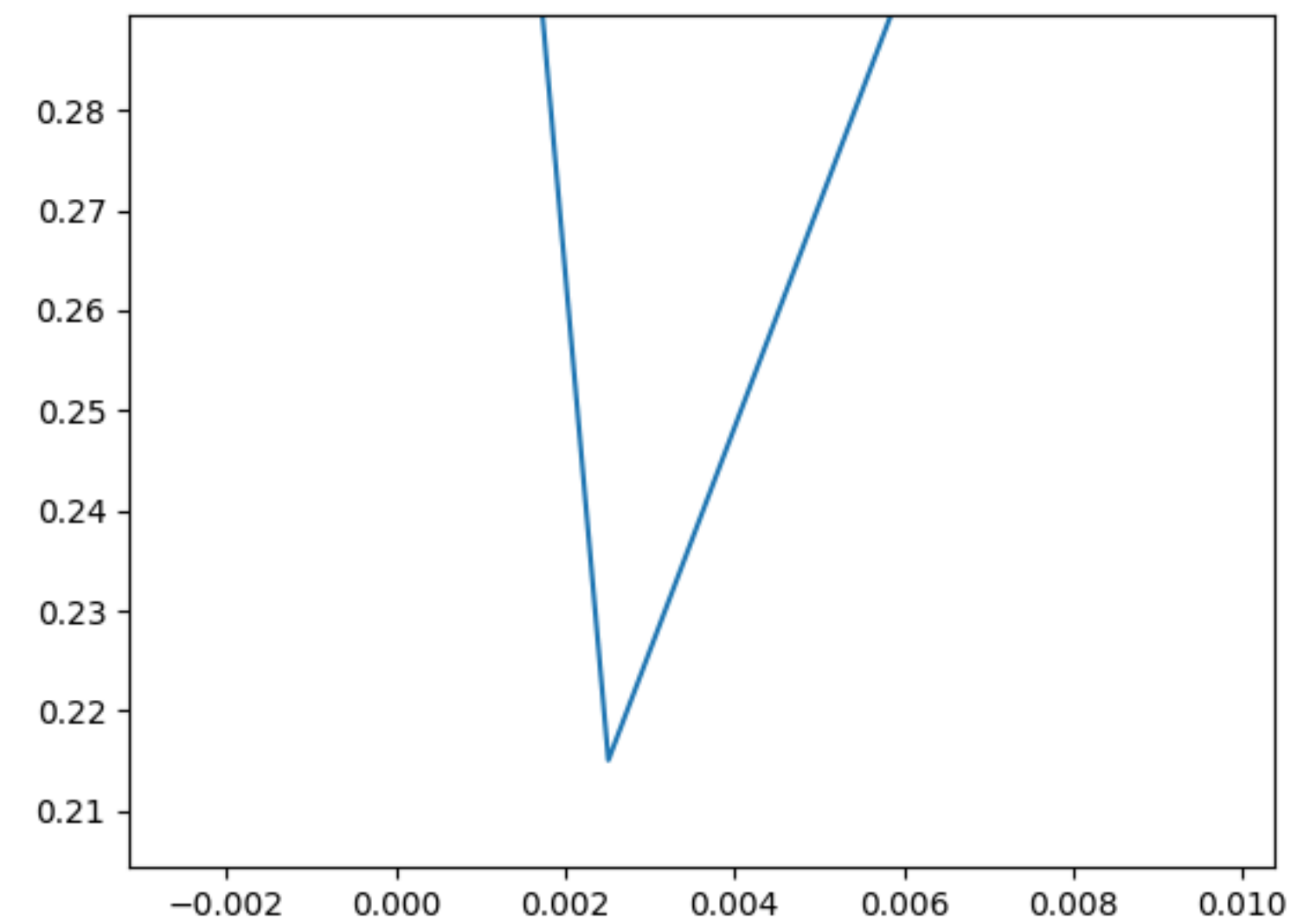
10

$$\mathbf{X} = N \times 10$$

# Find Learning Rate



Gradually increase learning run and run short epoch training session



Region where loss decreases and immediately increases is an optimal learning rate