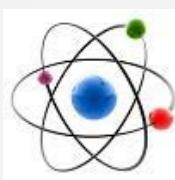


AN INTRODUCTION TO NEURAL MACHINE TRANSLATION

ATA59

Carola F. Berger



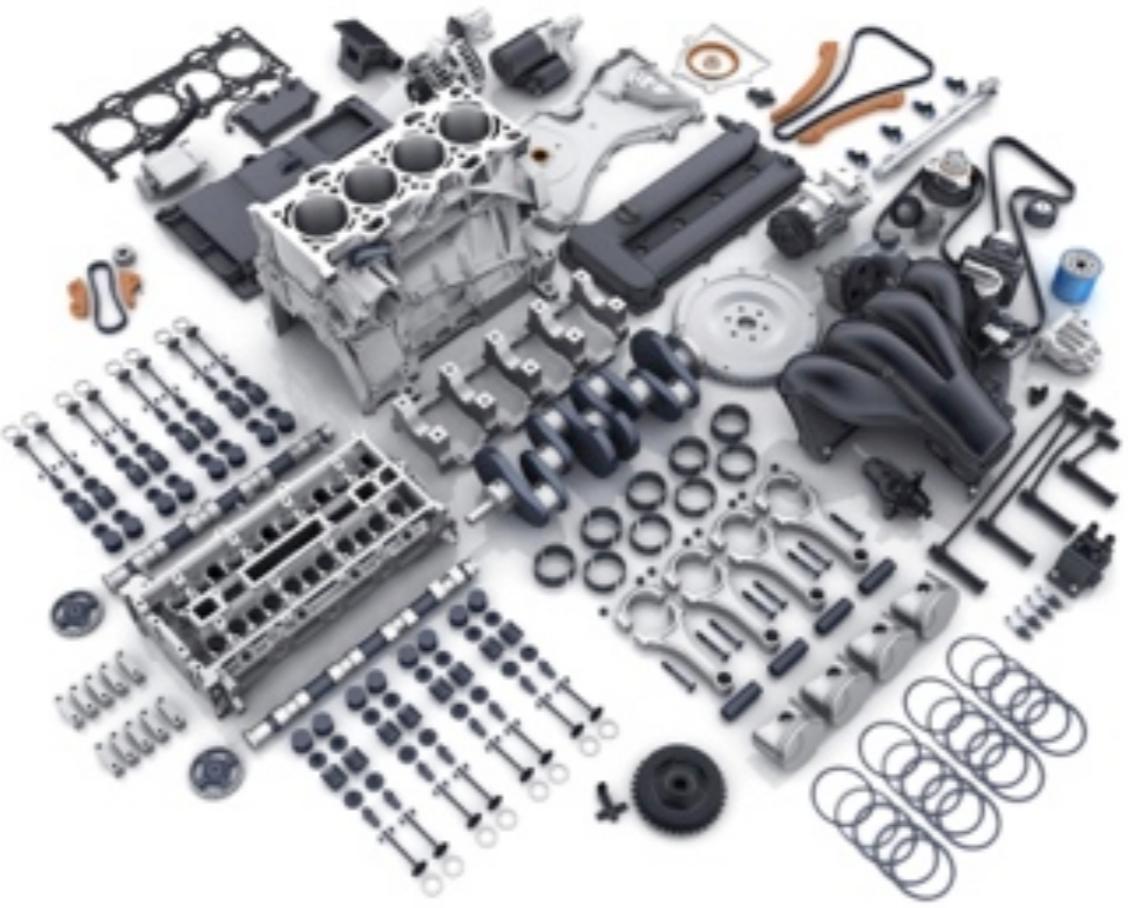
**TURN OFF
2-WAY RADIO
AND
CELL PHONE**



DISCLAIMER 1 – THIS PRESENTATION



DISCLAIMER 1 – THIS PRESENTATION



DISCLAIMER 2

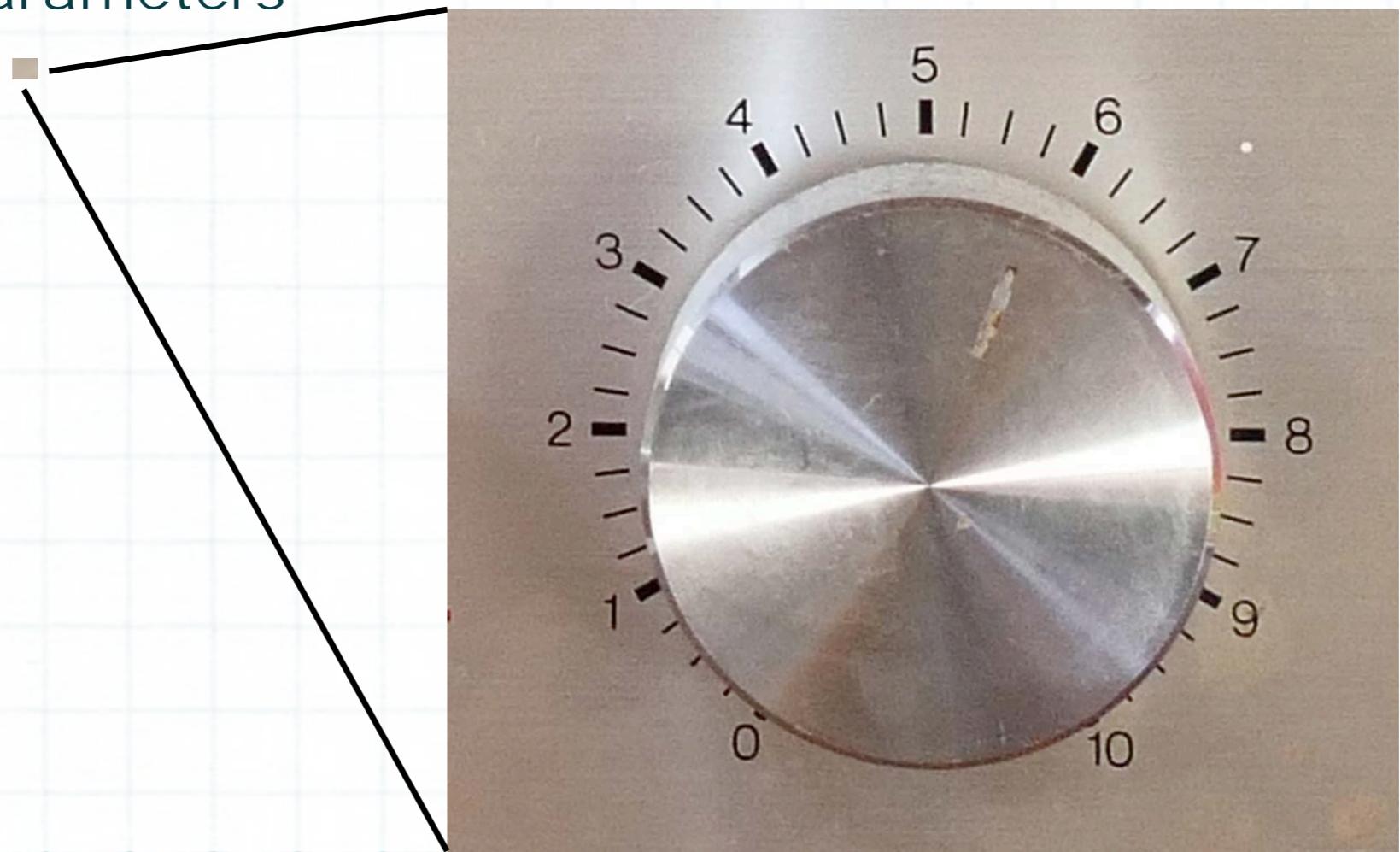
After this presentation, you will hopefully understand how neural MT works, but you will not understand what it does.



DISCLAIMER 2

After this presentation, you will hopefully understand how neural MT works, but you will not understand what it does.

~100 million parameters



OUTLINE

- Brief recap: previous MT approaches
- How do neural networks work?
- How do words get into and out of a neural network?



NMT APPROACHES

- Rules based:
Basically just grammar + dictionary
- Statistical MT:
Chop sentences up into n-grams (sequences of n words) or phrases.
Training of the engine: Calculate frequency = probability in source and target language.
Translation after training: Chop source sentences into n-grams or phrases, apply previously calculated probabilities.
1-gram SMT = ?

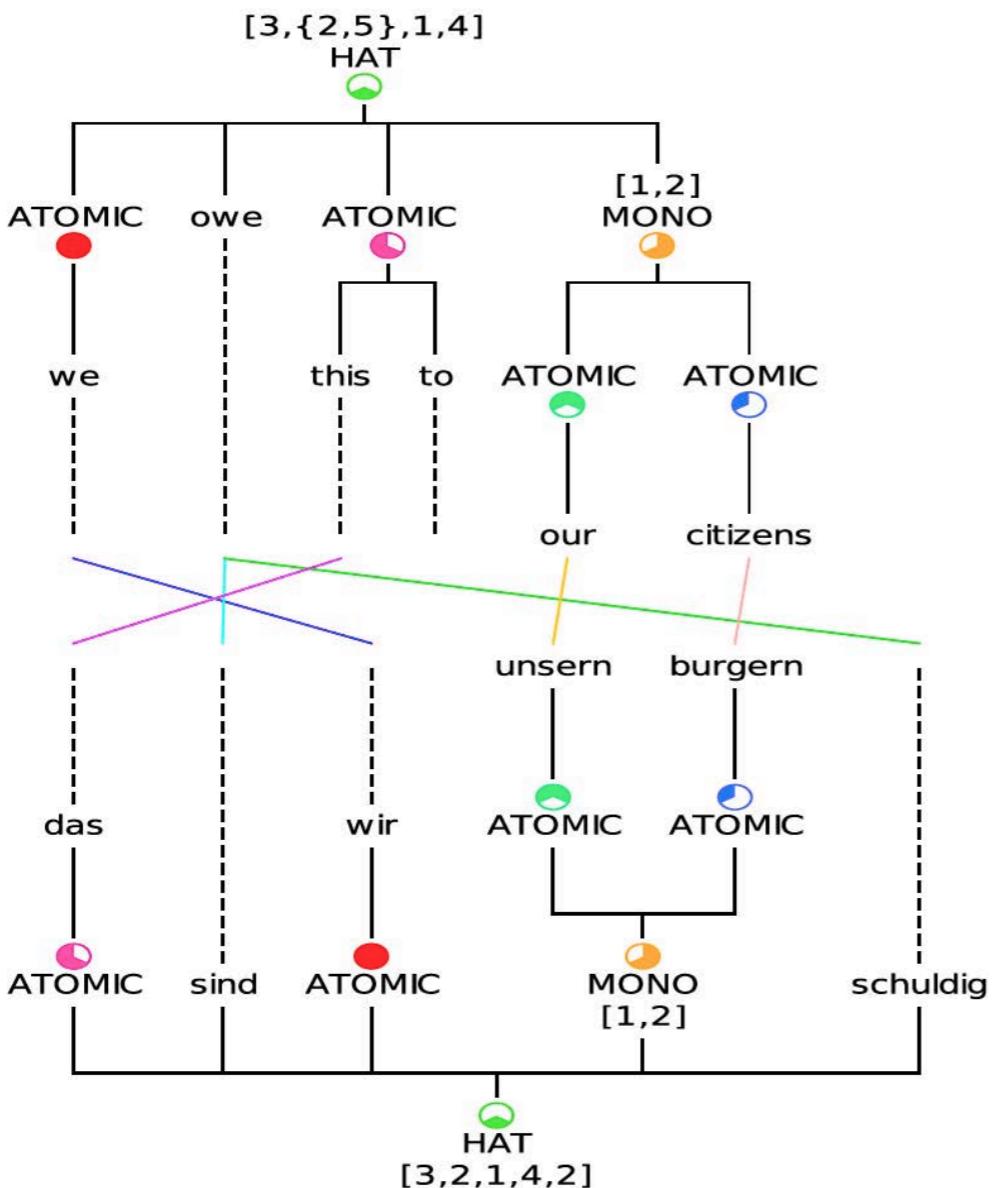


MT APPROACHES

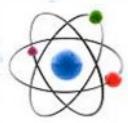
- Rules based:
Basically just grammar + dictionary
- Statistical MT:
Chop sentences up into n-grams (sequences of n words) or phrases.
Training of the engine: Calculate frequency = probability in source and target language.
Translation after training: Chop source sentences into n-grams or phrases, apply previously calculated probabilities.
1-gram SMT = dictionary
- Pros: Output is deterministic. No words missing.
Cons: Context!
- Neural MT – this presentation



STATISTICAL NMT

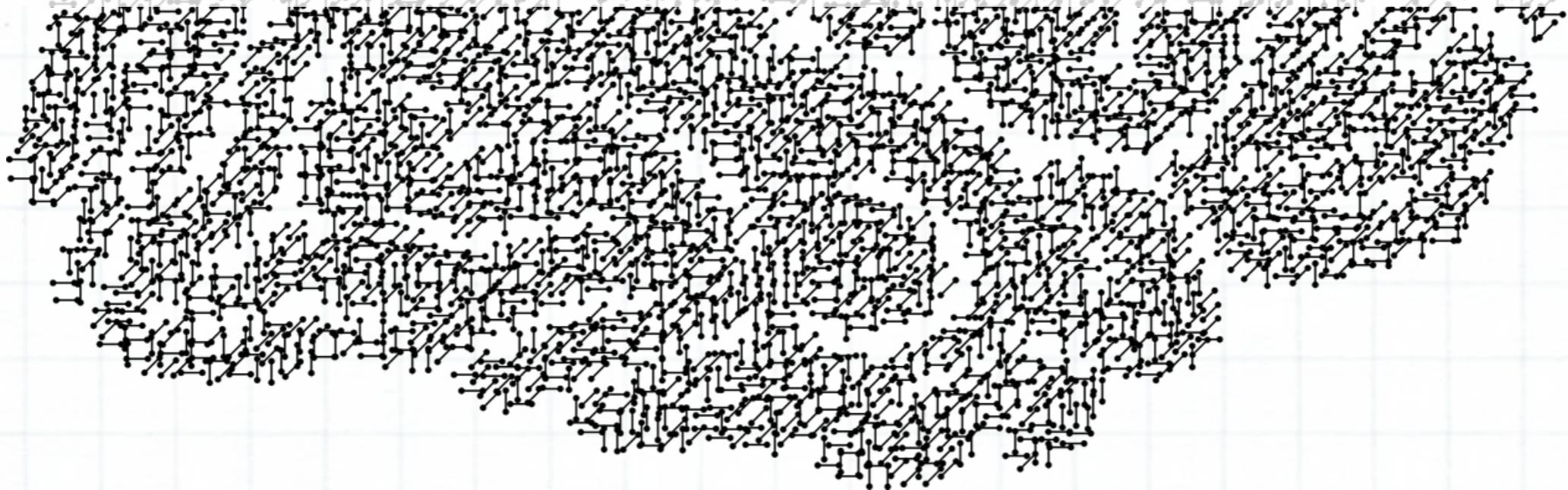


From G. M. de Buy Wenninger, K. Sima'an,
PBML No. 101, April 2014, pp. 43

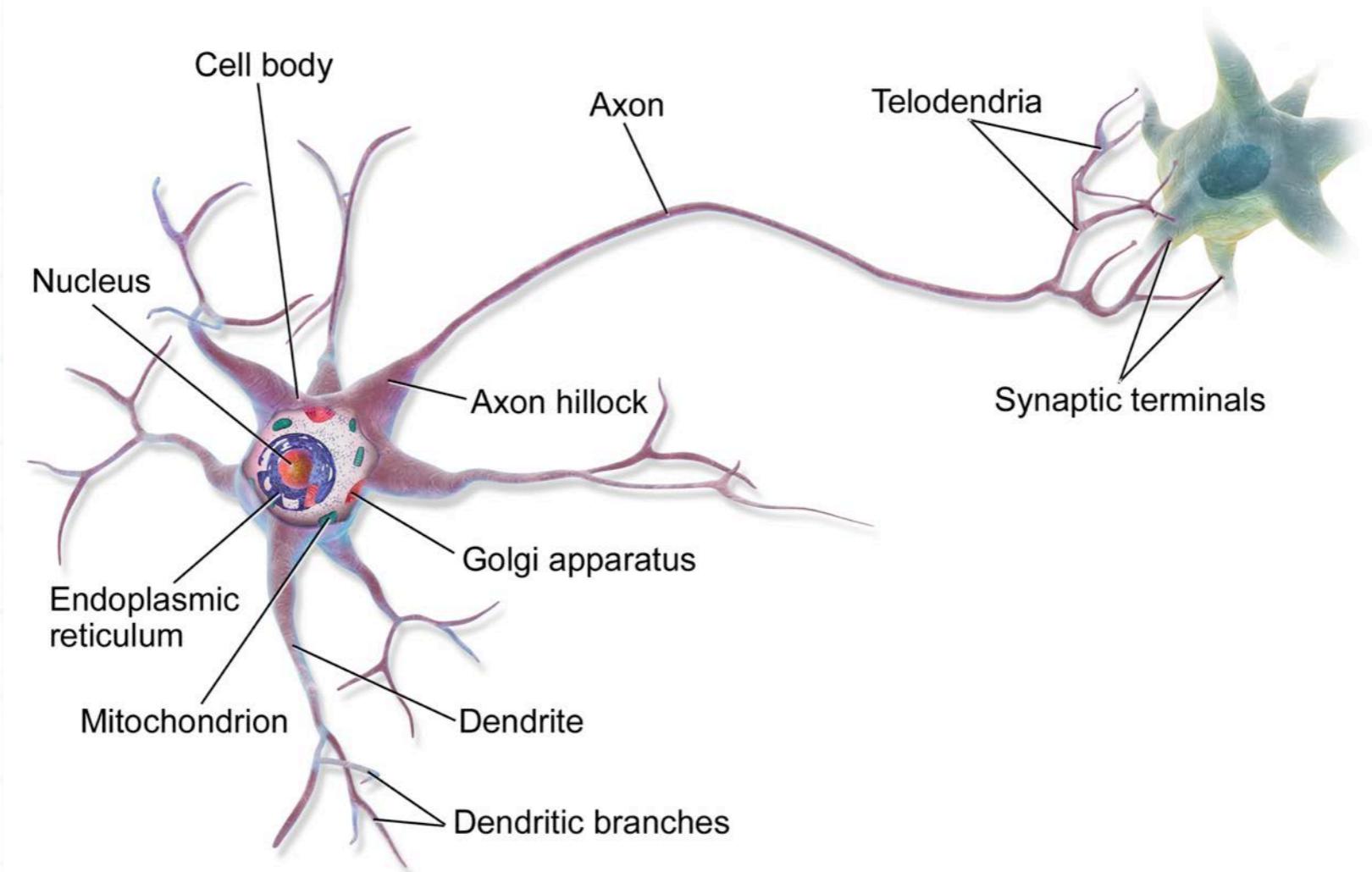


HOW DO NEURAL NETWORKS WORK?

A (not so) brief recap of last year's presentation at ATA58. See also handouts as PDF in the app or on my website (see references).



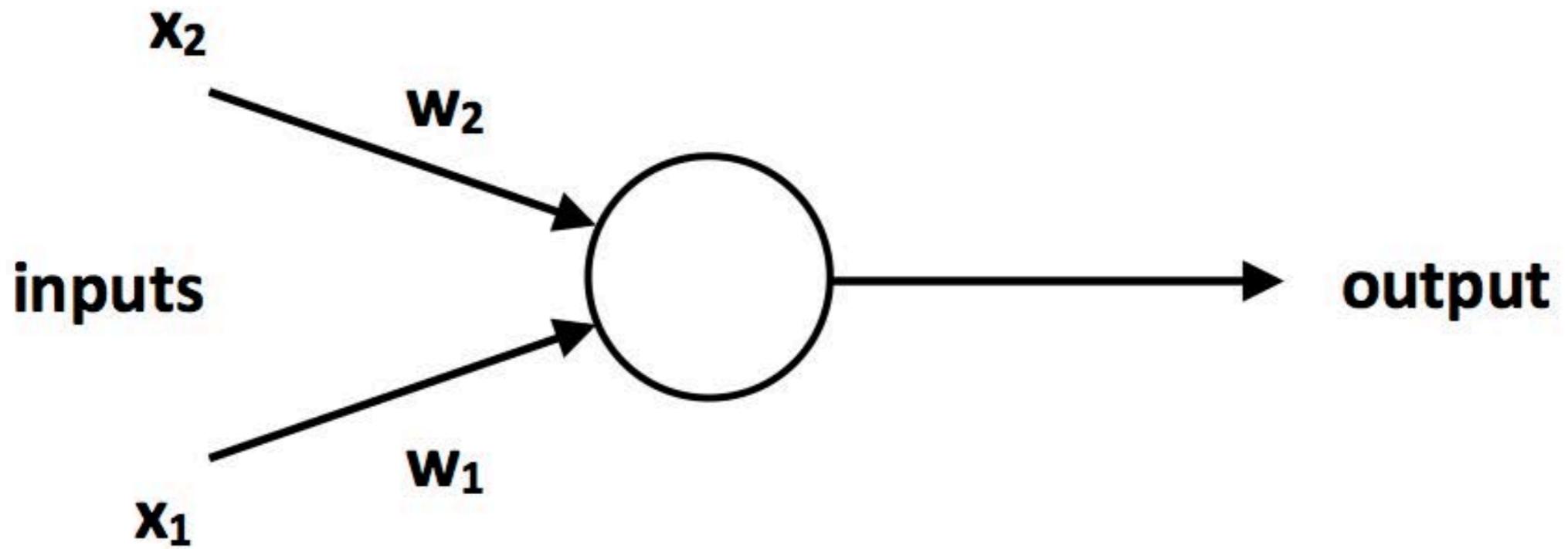
BIOLOGICAL NEURON



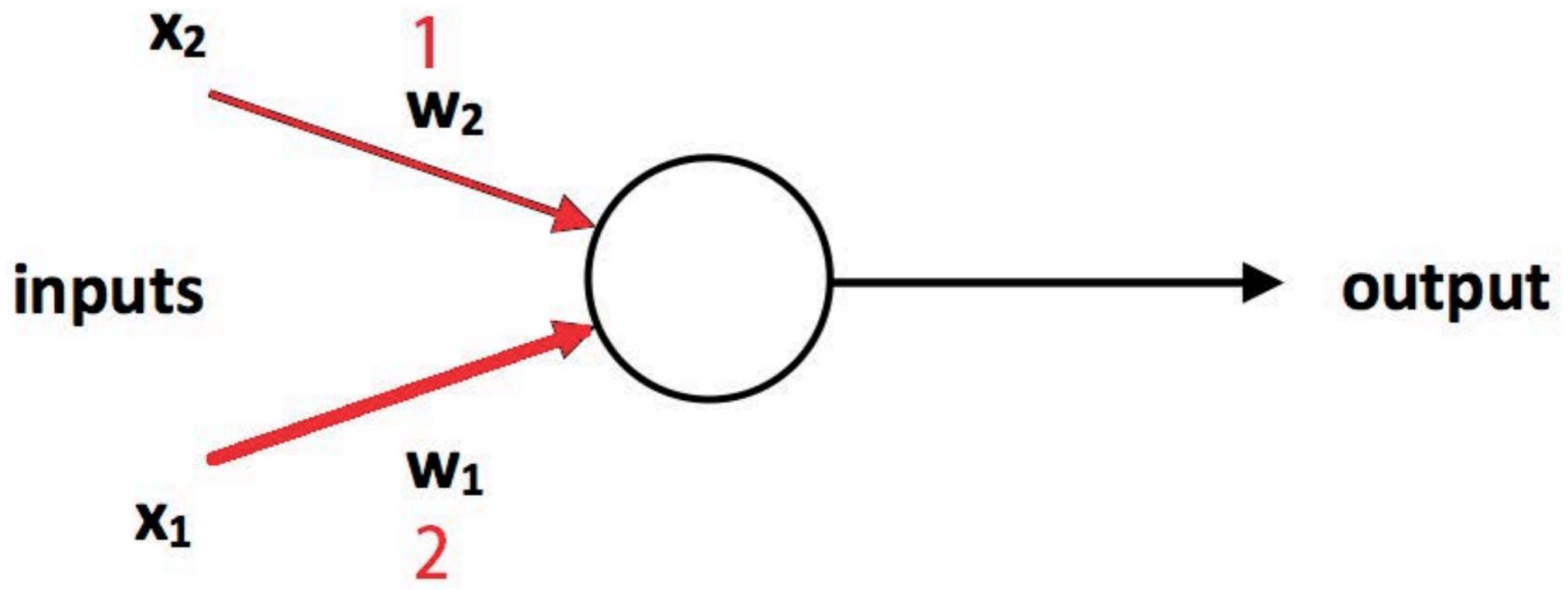
Bruce Blaus, https://commons.wikimedia.org/wiki/File:Blausen_0657_MultipolarNeuron.png



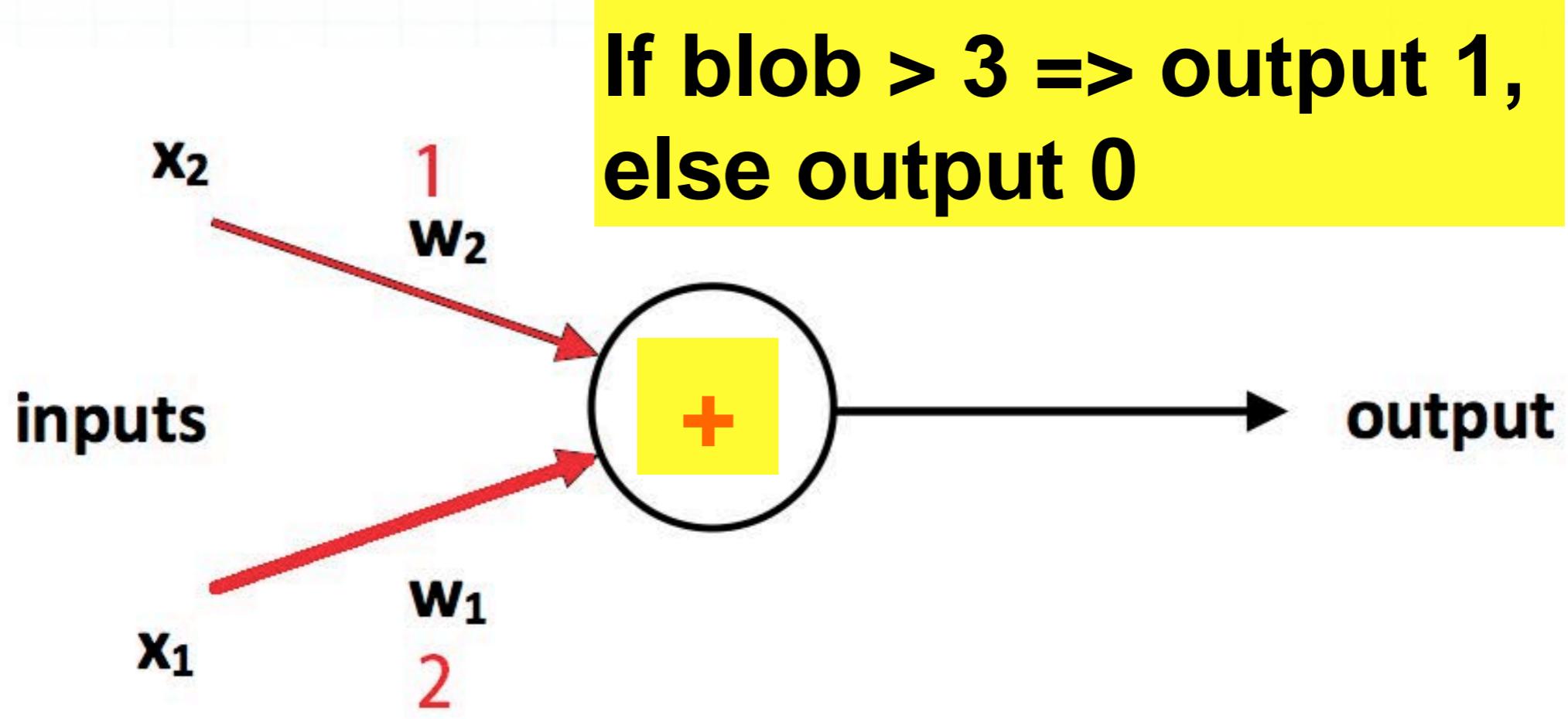
ARTIFICIAL NEURON



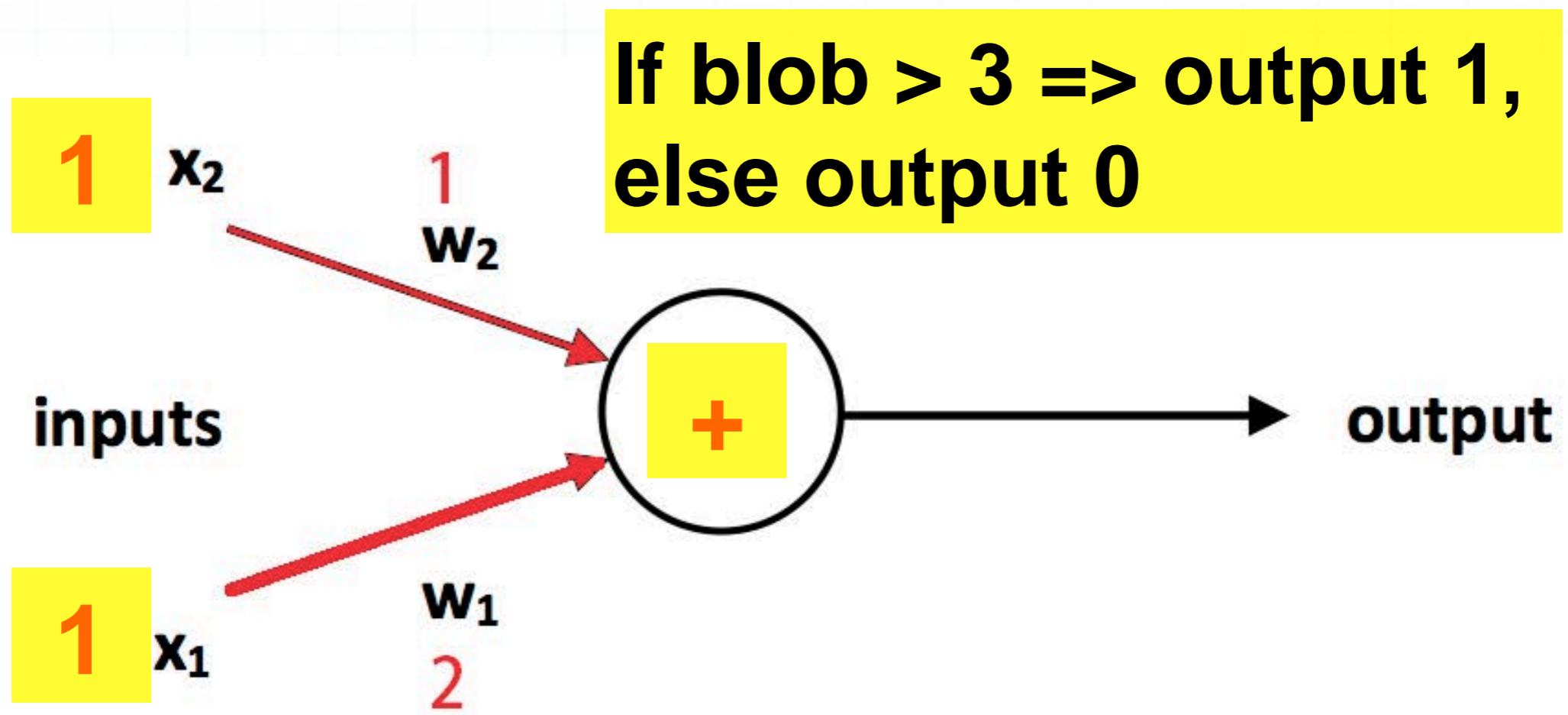
ARTIFICIAL NEURON



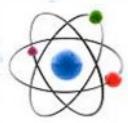
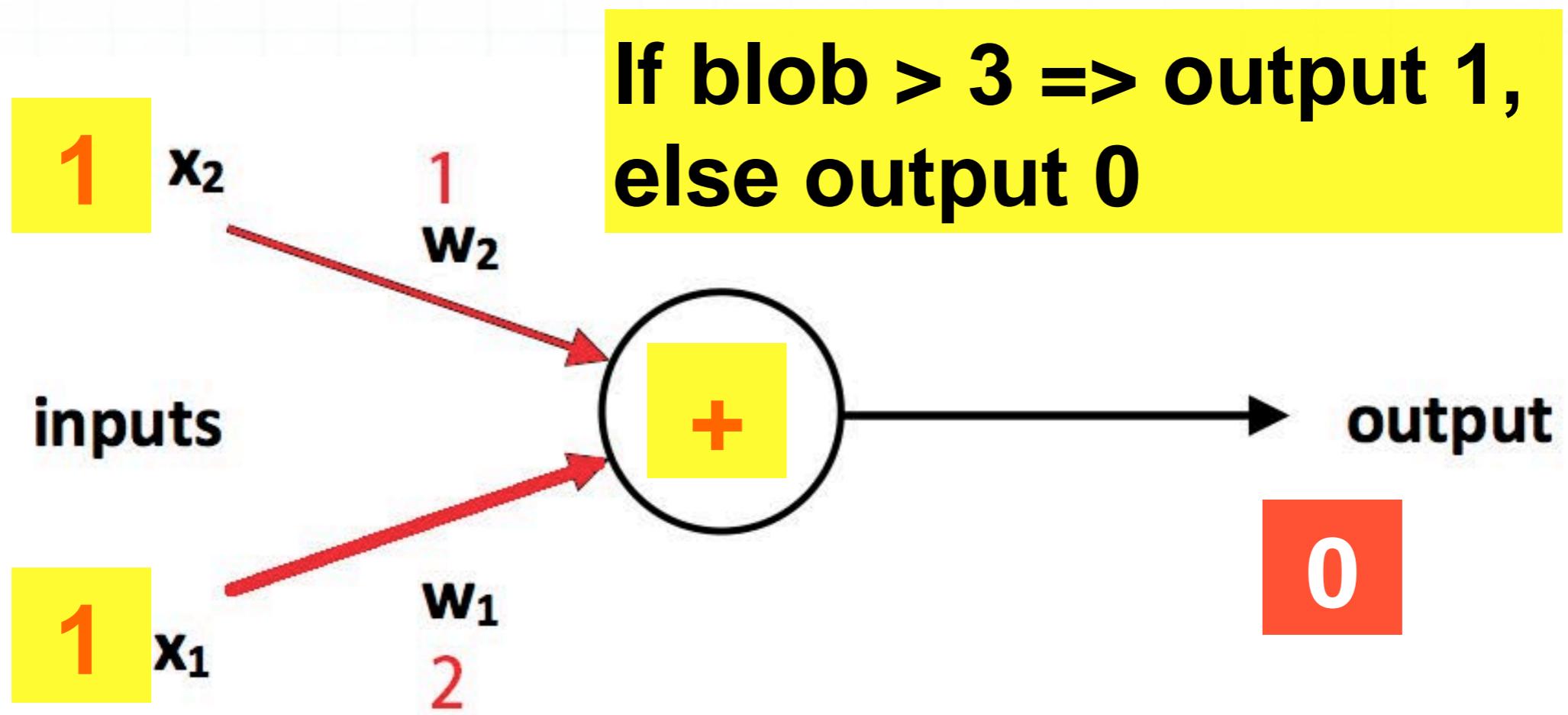
ARTIFICIAL NEURON - PERCEPTRON



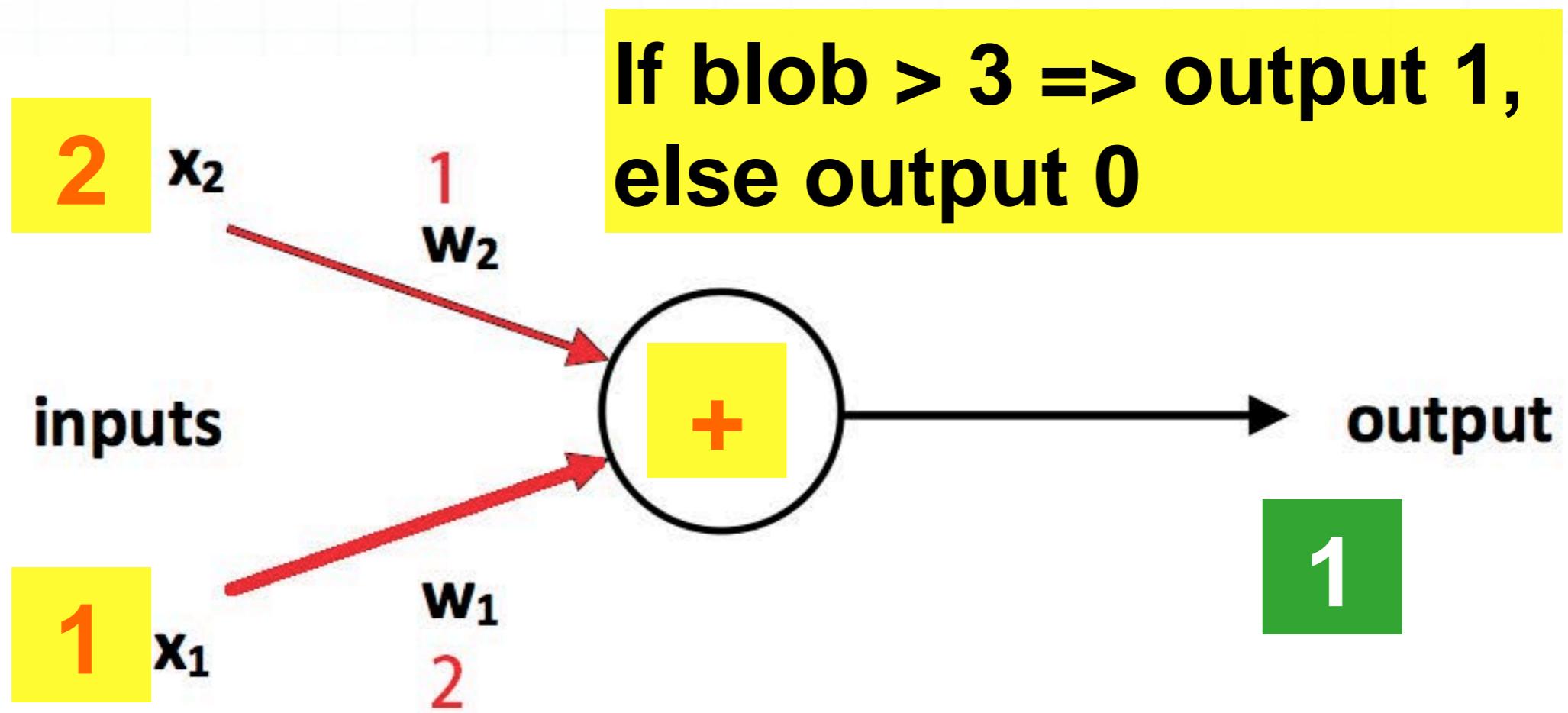
ARTIFICIAL NEURON - PERCEPTRON



ARTIFICIAL NEURON - PERCEPTRON

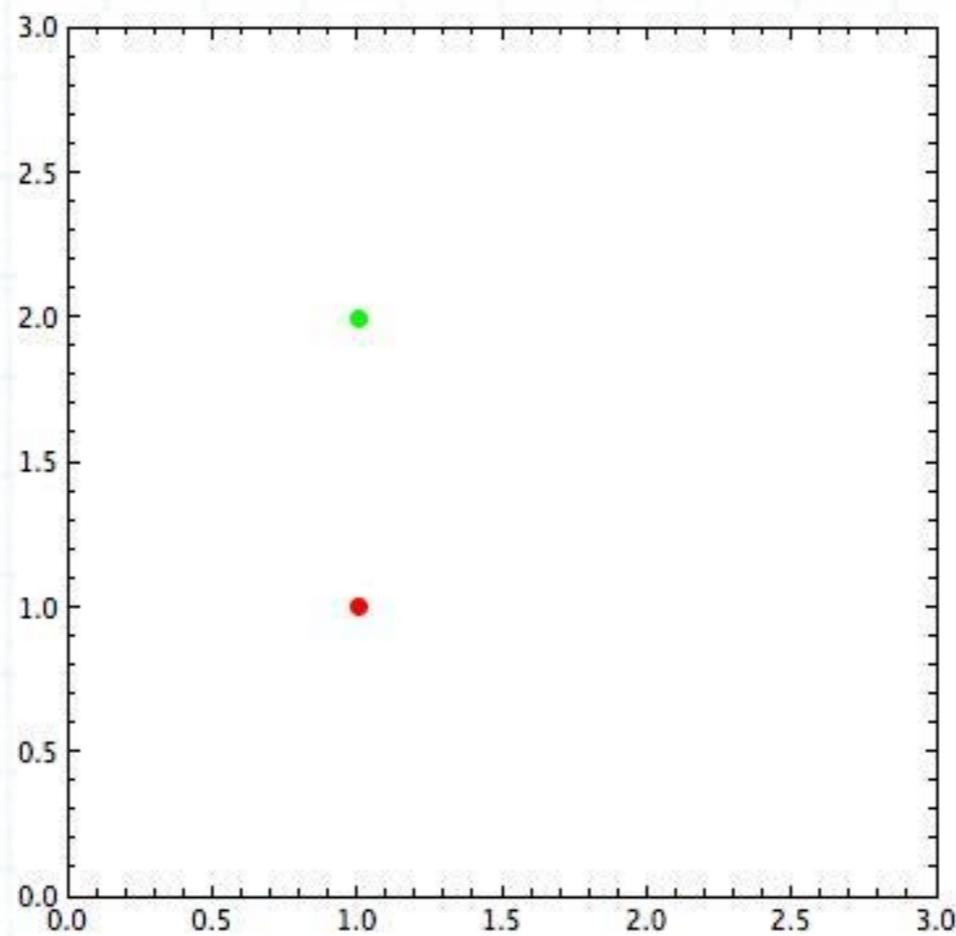


ARTIFICIAL NEURON - PERCEPTRON

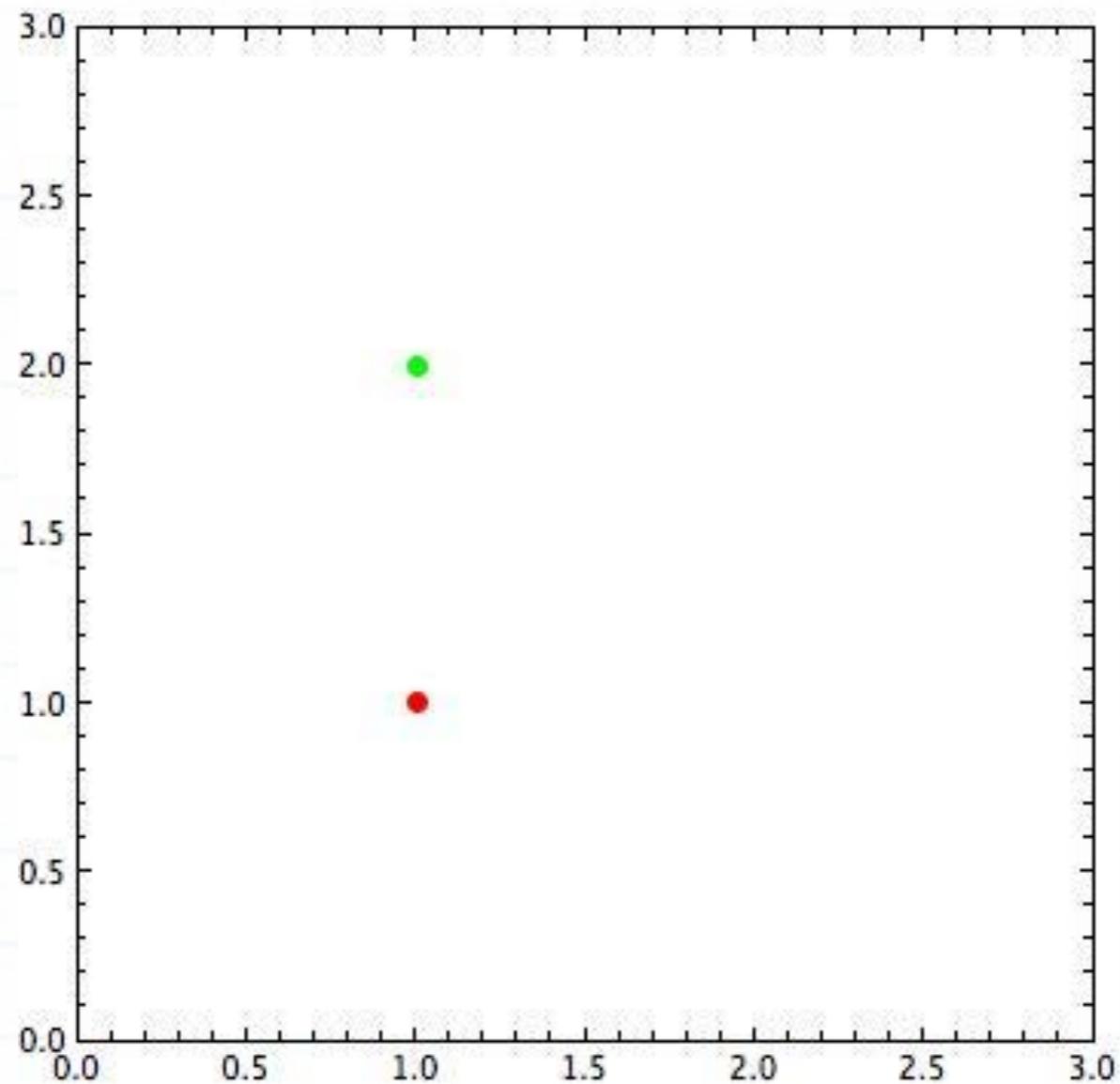


ARTIFICIAL NEURON - PERCEPTRON

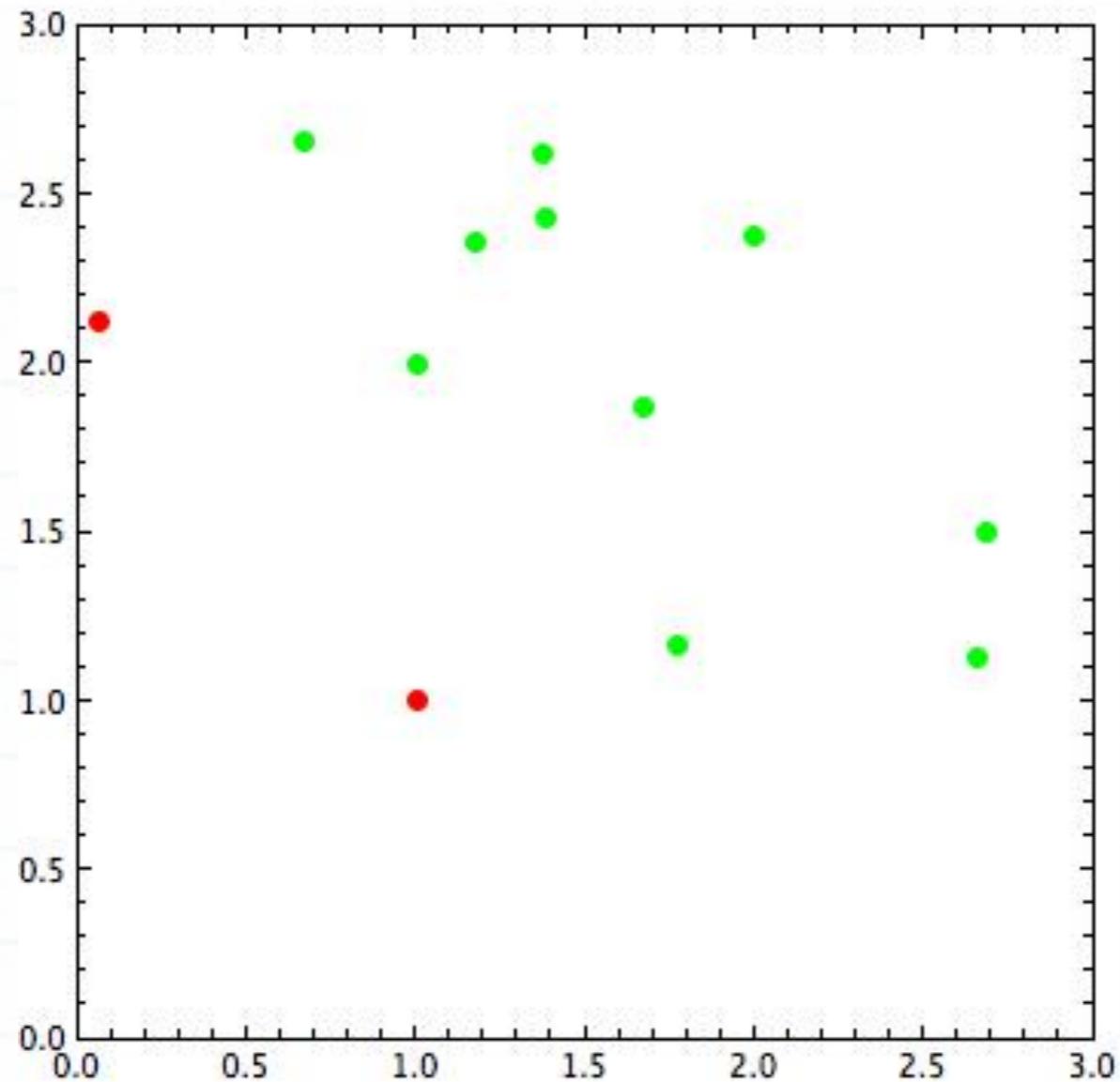
x1	x2	OUTPUT
1	1	0
1	2	1



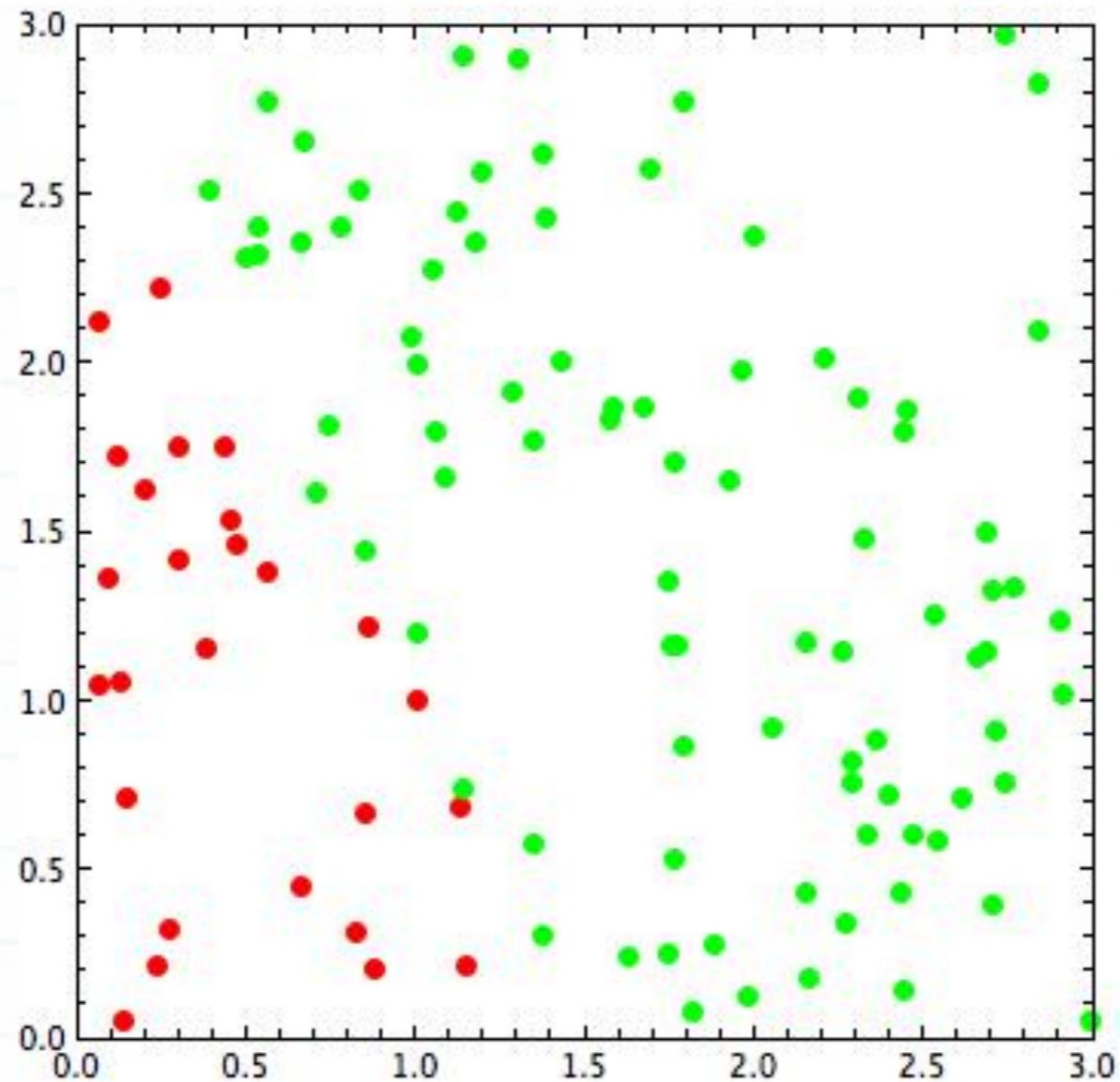
ARTIFICIAL NEURON - PERCEPTRON



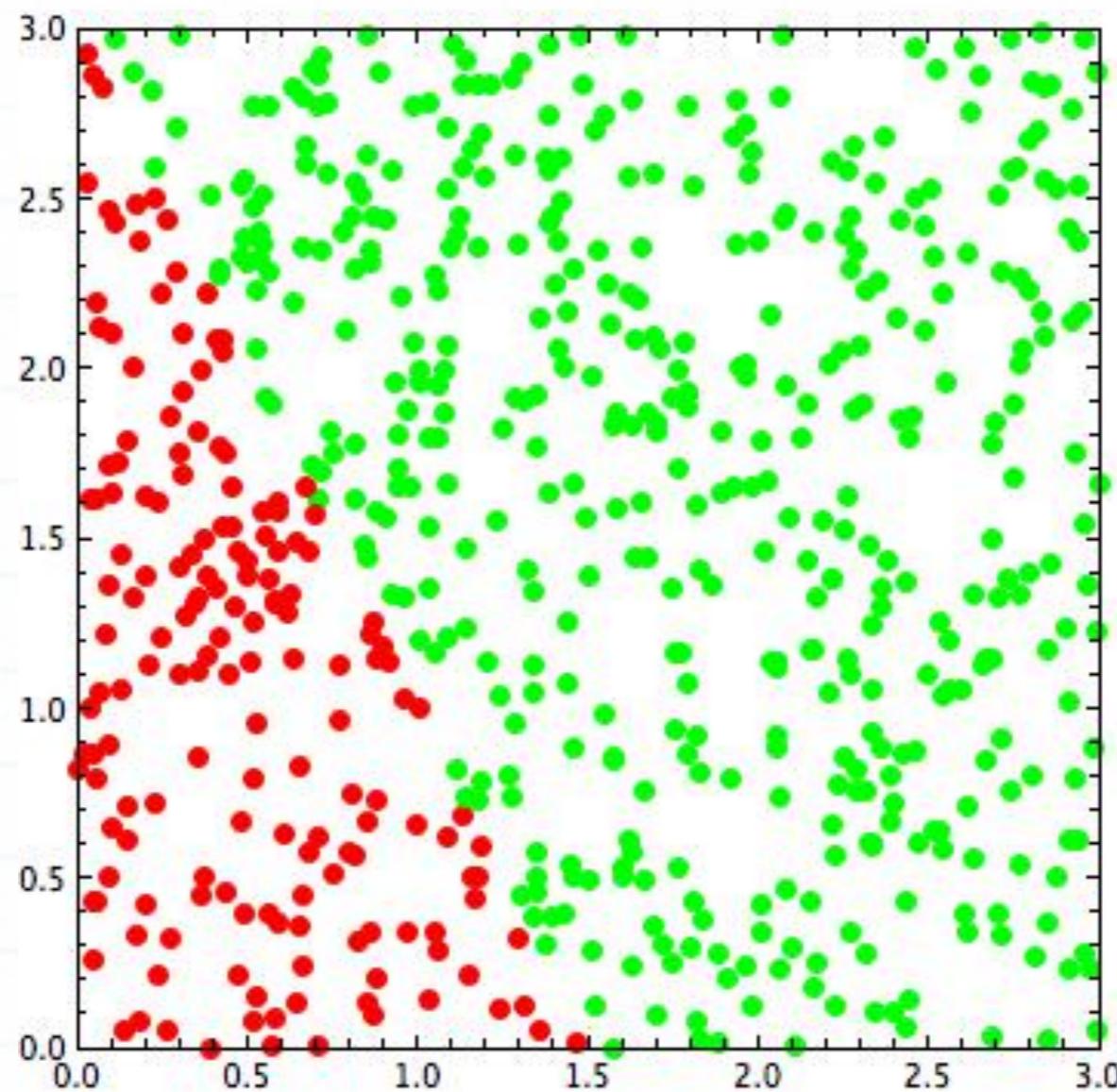
ARTIFICIAL NEURON - PERCEPTRON



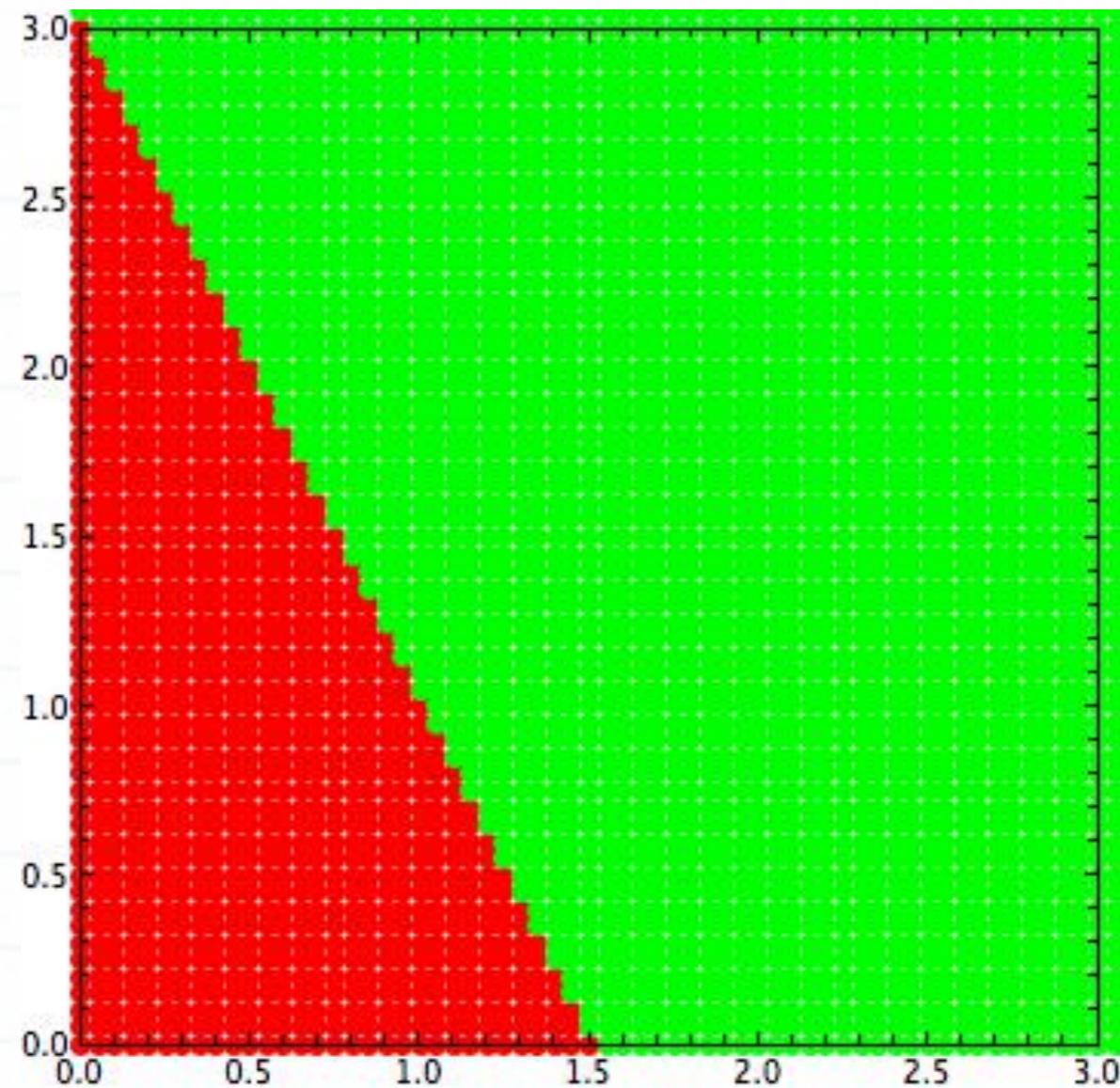
ARTIFICIAL NEURON - PERCEPTRON



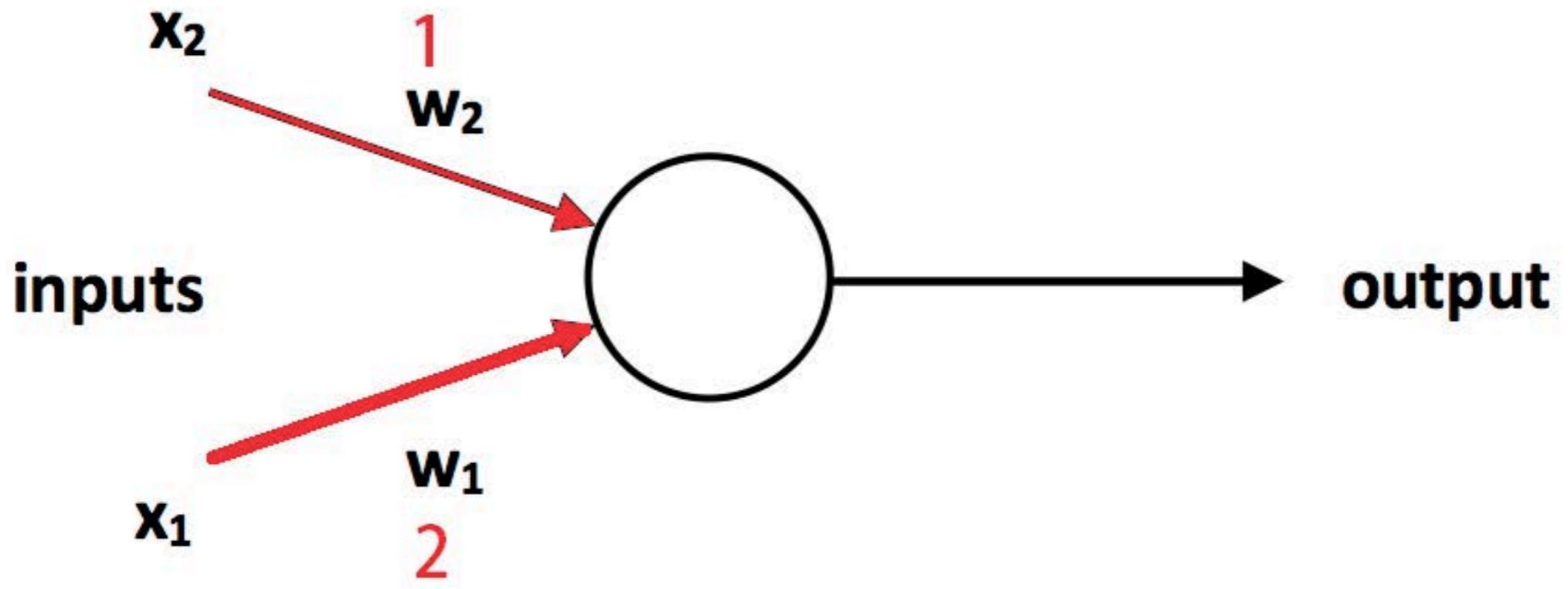
ARTIFICIAL NEURON - PERCEPTRON



ARTIFICIAL NEURON - PERCEPTRON

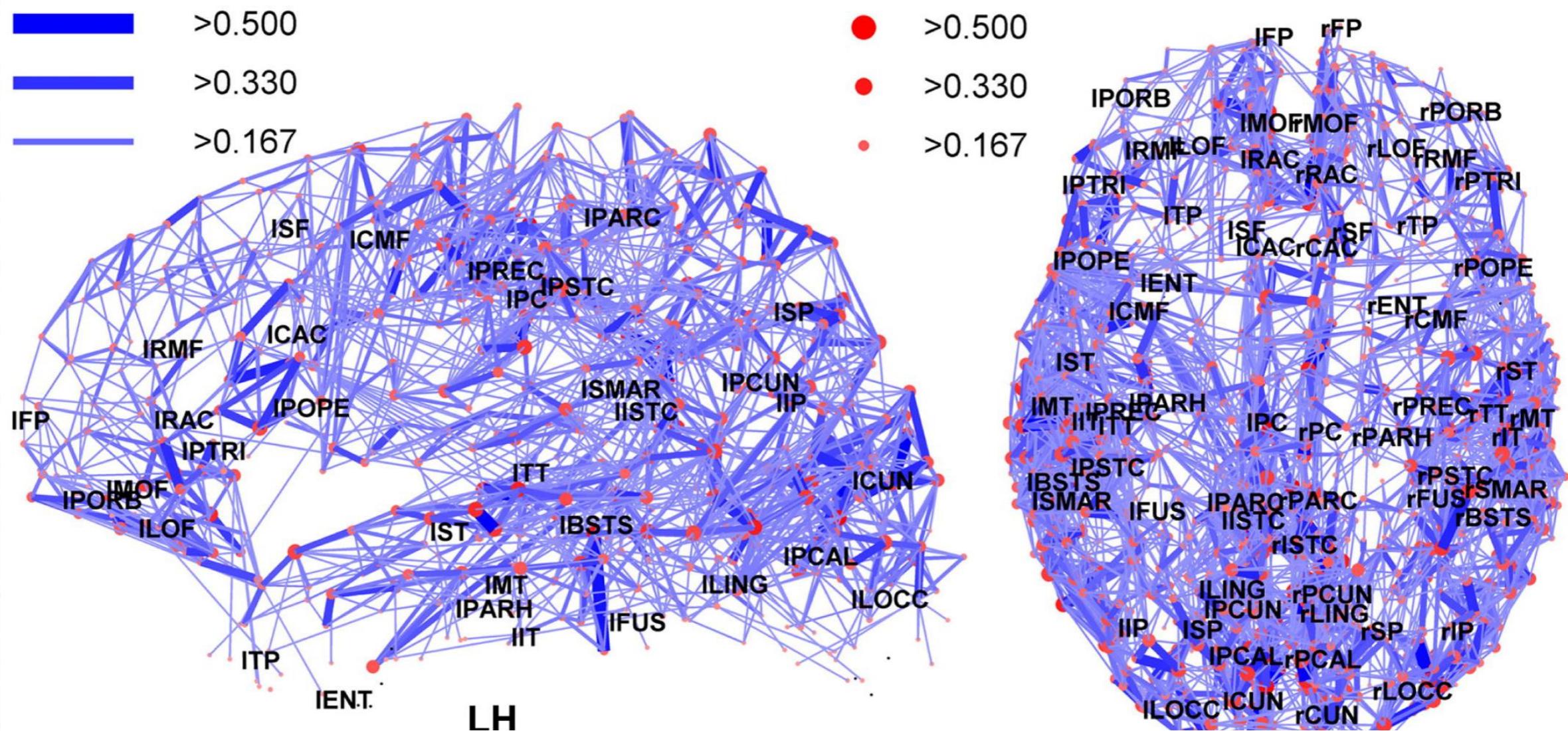


ARTIFICIAL NEURON



NEURAL NETWORK

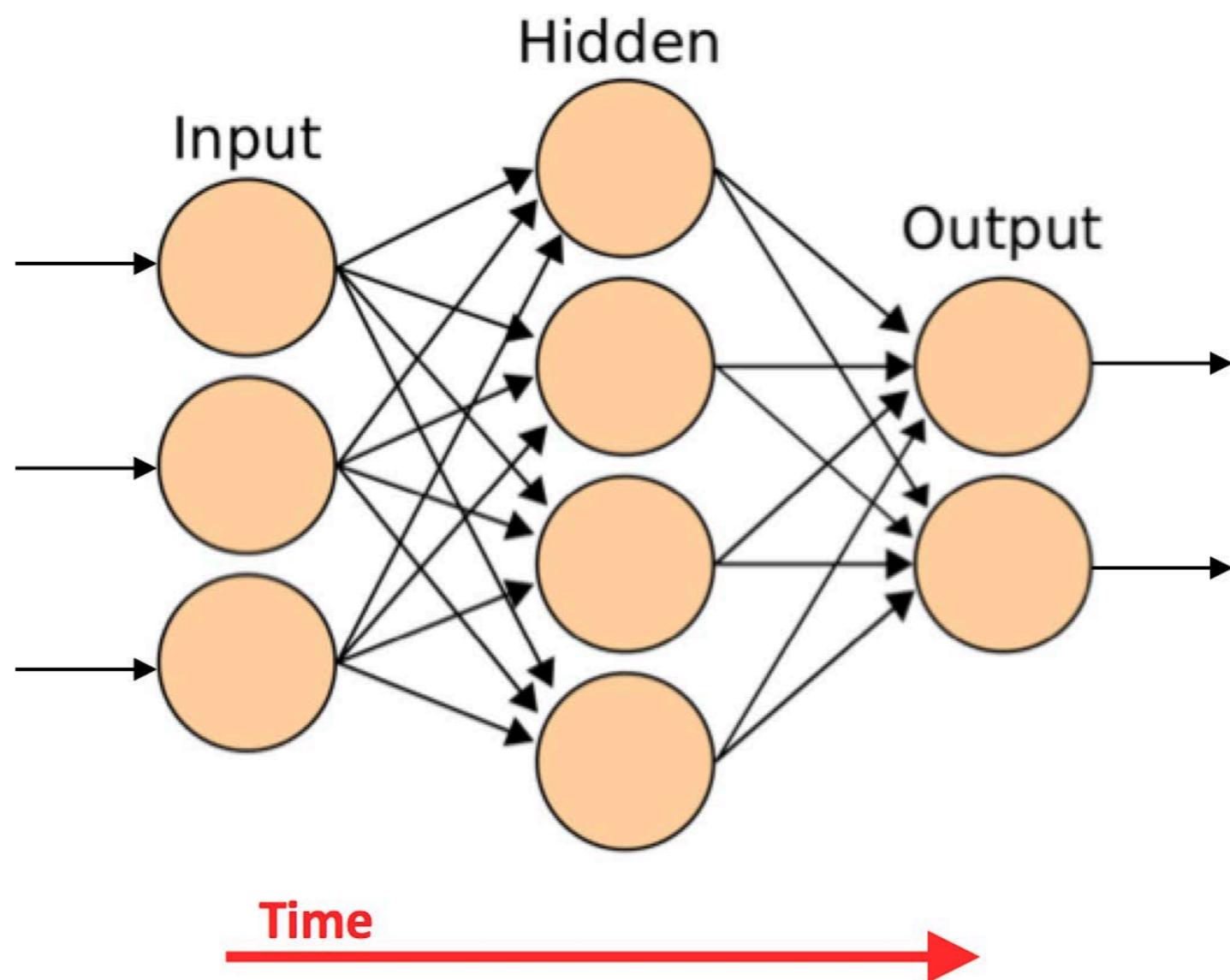
→ Human brain:



Hagmann P, Cammoun L, Gigandet X, Meuli R, Honey CJ, Wedeen VJ, Sporns O (2008)
Mapping the structural core of human cerebral cortex. PLoS Biology Vol. 6, No. 7, e159.



ARTIFICIAL NEURAL NETWORK

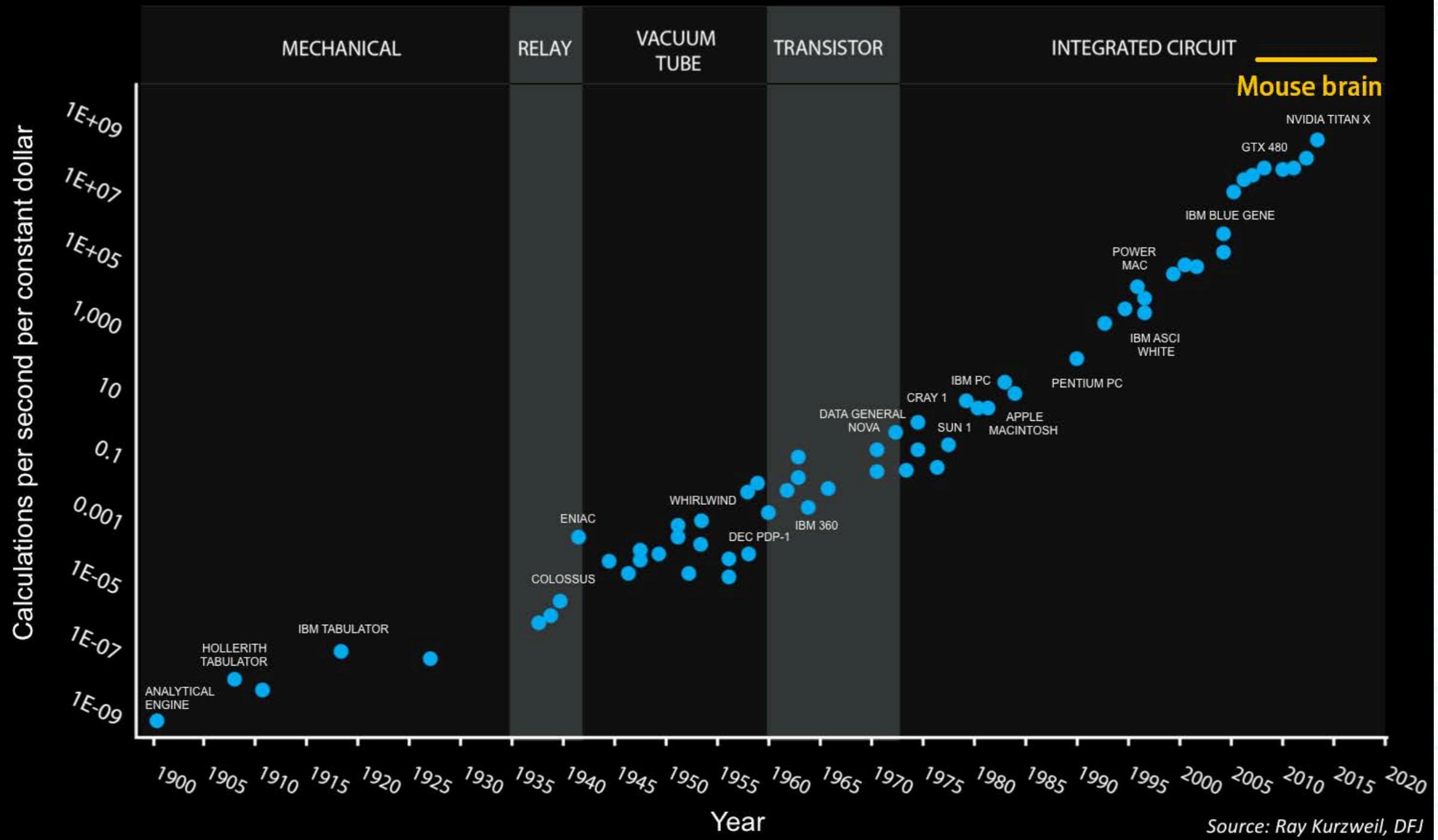


Adapted from: Cburnett, https://commons.wikimedia.org/wiki/File:Artificial_neural_network.svg



120 Years of Moore's Law

Human brain

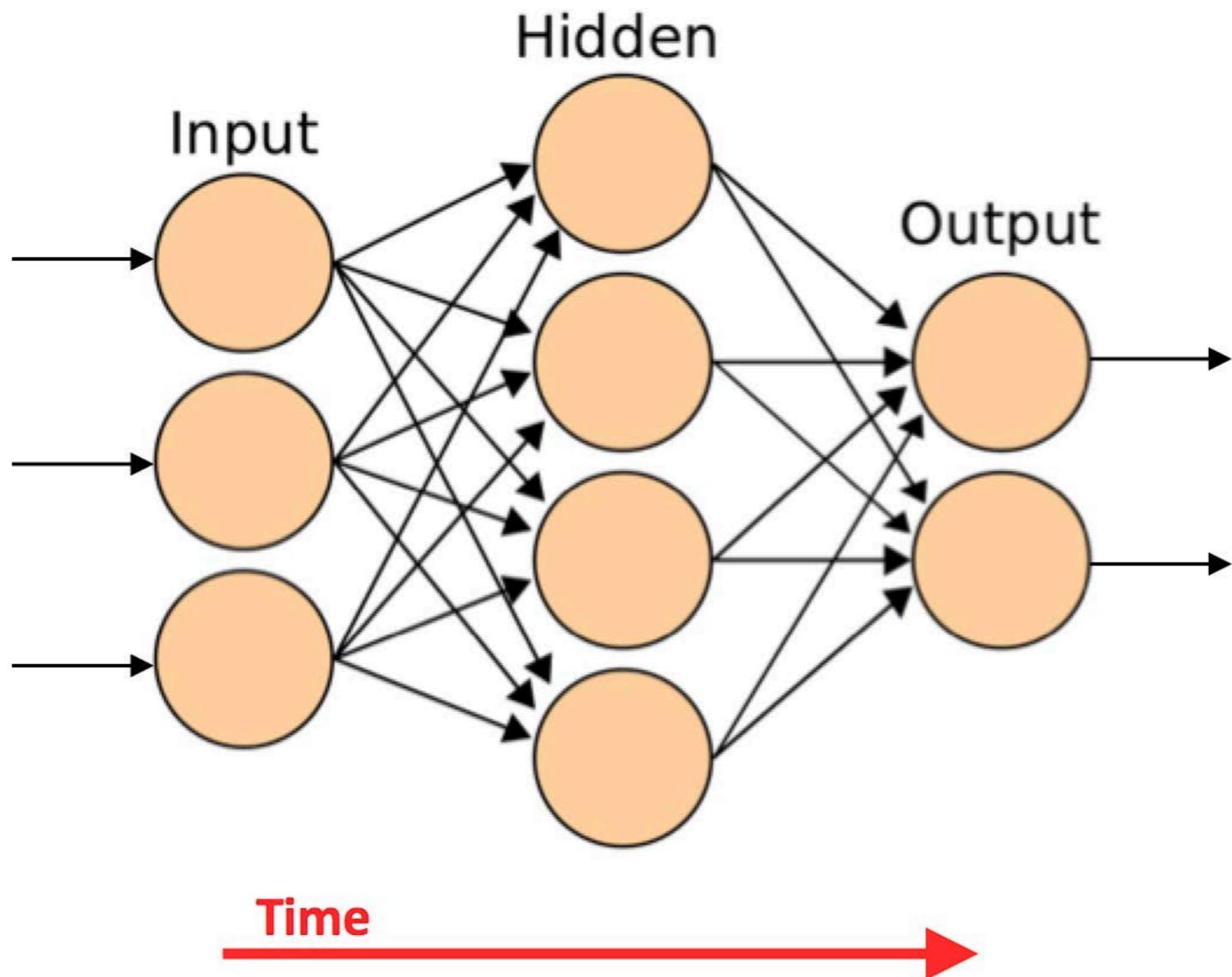


Adapted from: S. Jurvetson, <https://www.flickr.com/photos/jurvetson/31409423572>



ARTIFICIAL NEURAL NETWORK

→ Training:



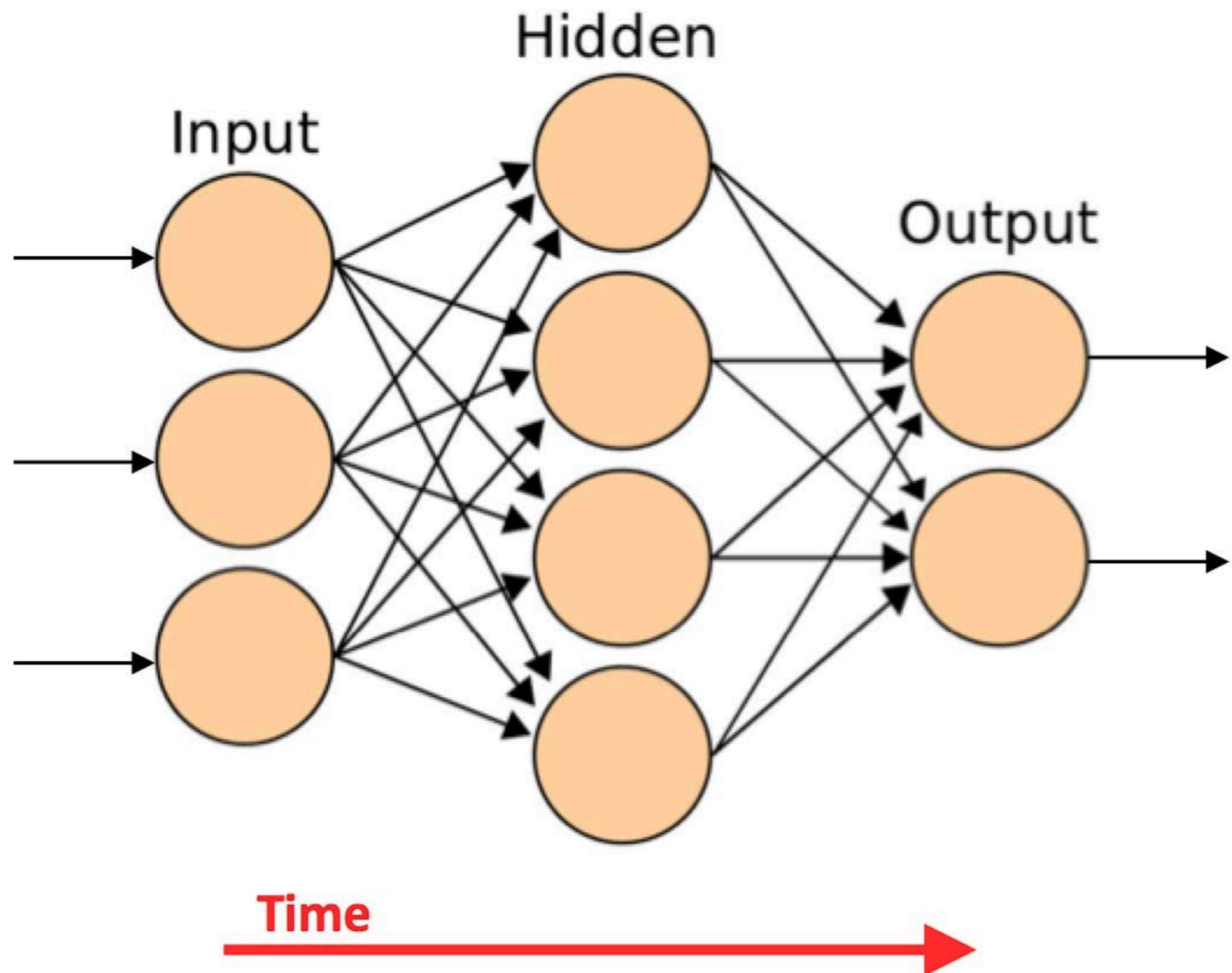
Adapted from: Cburnett, https://commons.wikimedia.org/wiki/File:Artificial_neural_network.svg



ARTIFICIAL NEURAL NETWORK

→ Training:

Feed in
training data



Adapt weights (“arrows”) according to difference between desired output and actual output, e.g. by backpropagation

Adapted from: Cburnett, https://commons.wikimedia.org/wiki/File:Artificial_neural_network.svg



NEURAL NET EXAMPLE – DIGIT RECOGNITION

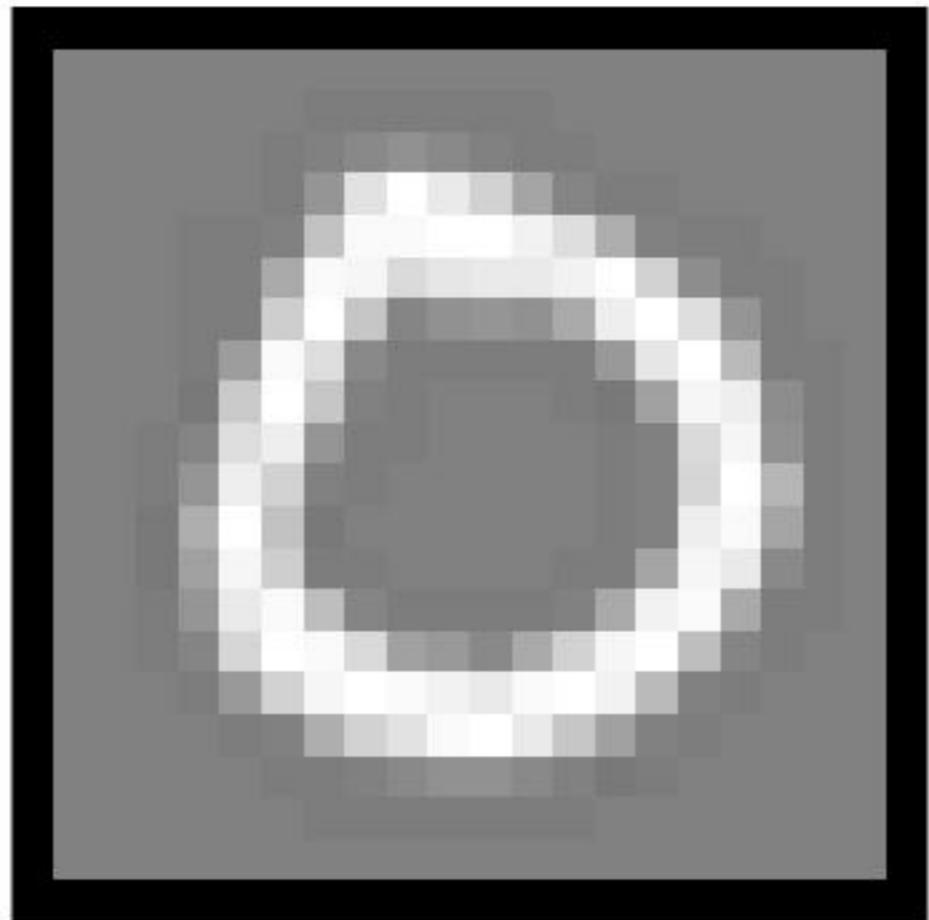
Sample input data

8	9	3	1	4	5	9	0	3	3
5	3	7	6	7	5	8	8	5	3
8	9	8	5	7	2	0	9	8	7
4	6	6	6	0	3	9	6	8	9
8	1	8	3	5	9	3	3	2	7
8	5	1	3	9	8	2	0	8	7
9	8	8	1	5	6	5	9	4	9
6	5	0	0	2	7	4	8	3	1
4	5	2	2	2	1	2	4	8	1
4	6	9	2	2	7	6	0	8	5



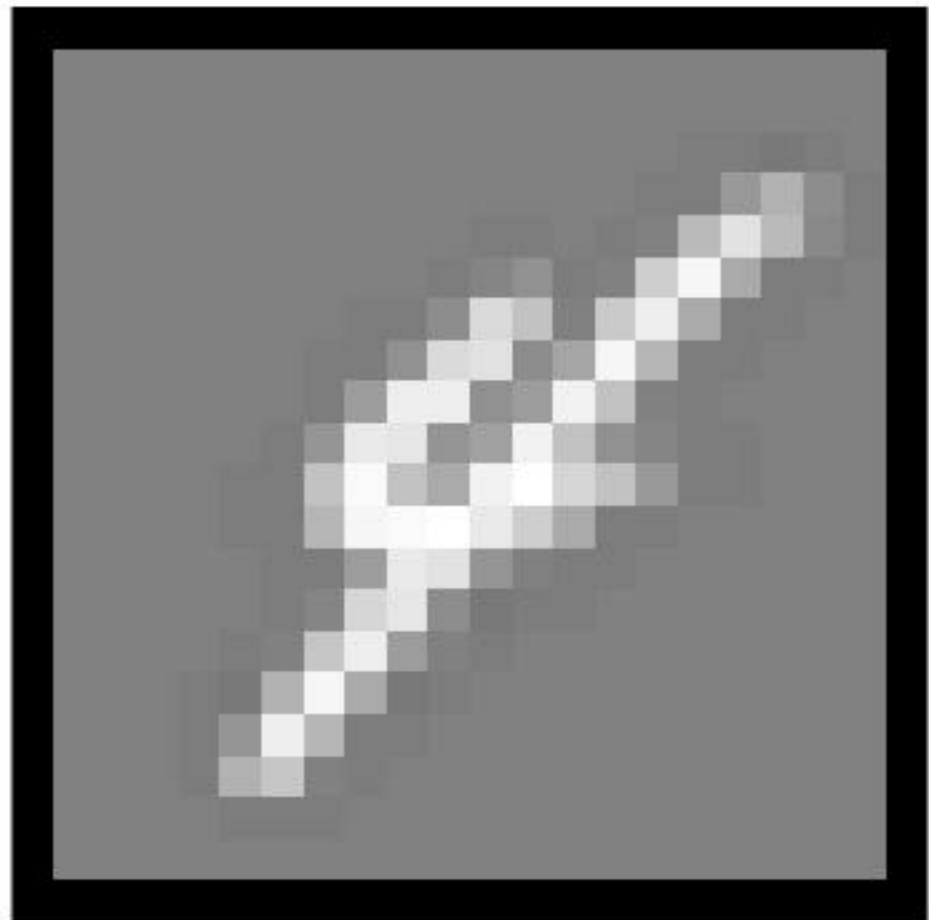
NEURAL.NET EXAMPLE – DIGIT RECOGNITION

Sample input (20x20 pixels)



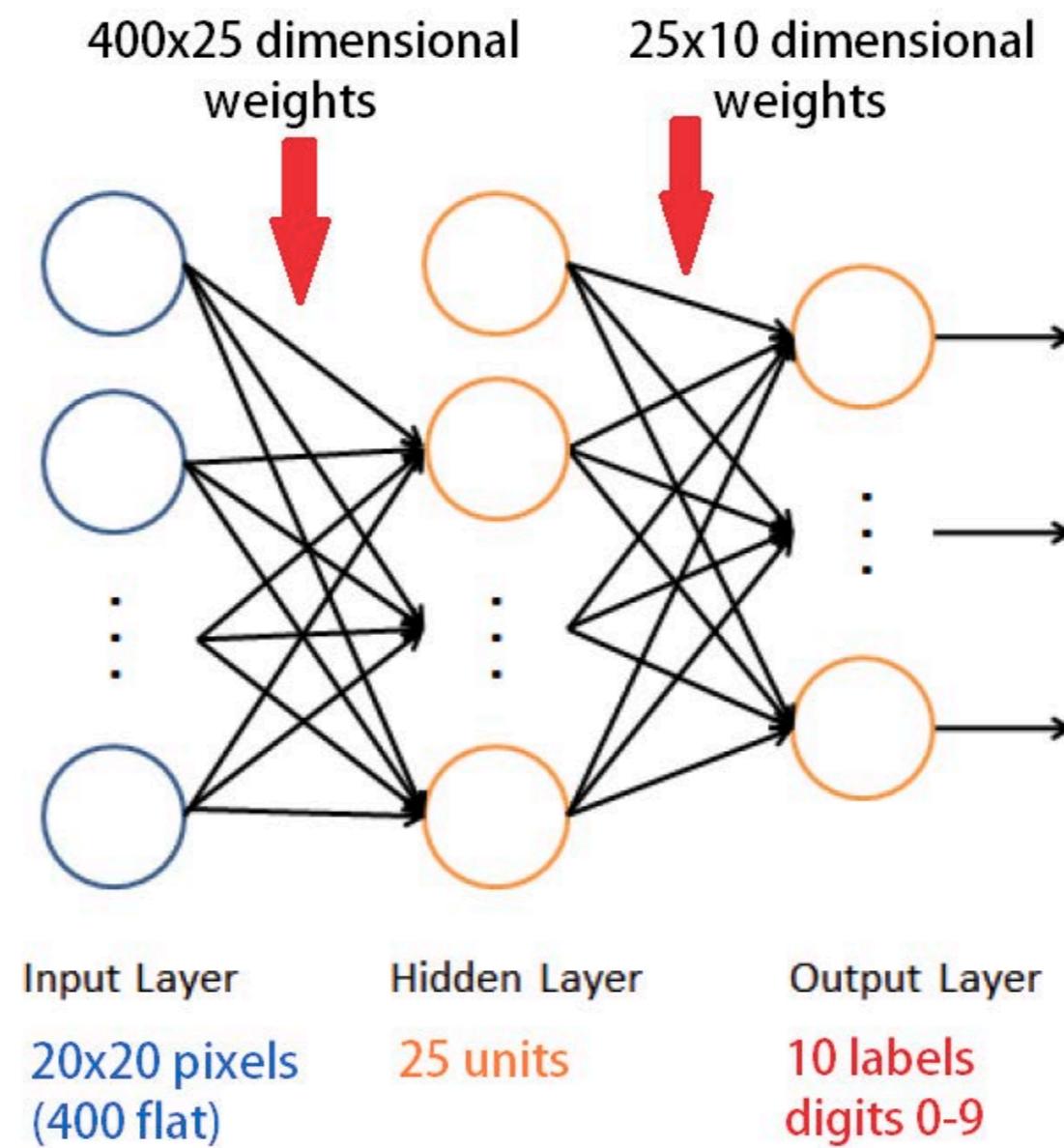
NEURAL.NET EXAMPLE – DIGIT RECOGNITION

Sample input (20x20 pixels)



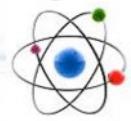
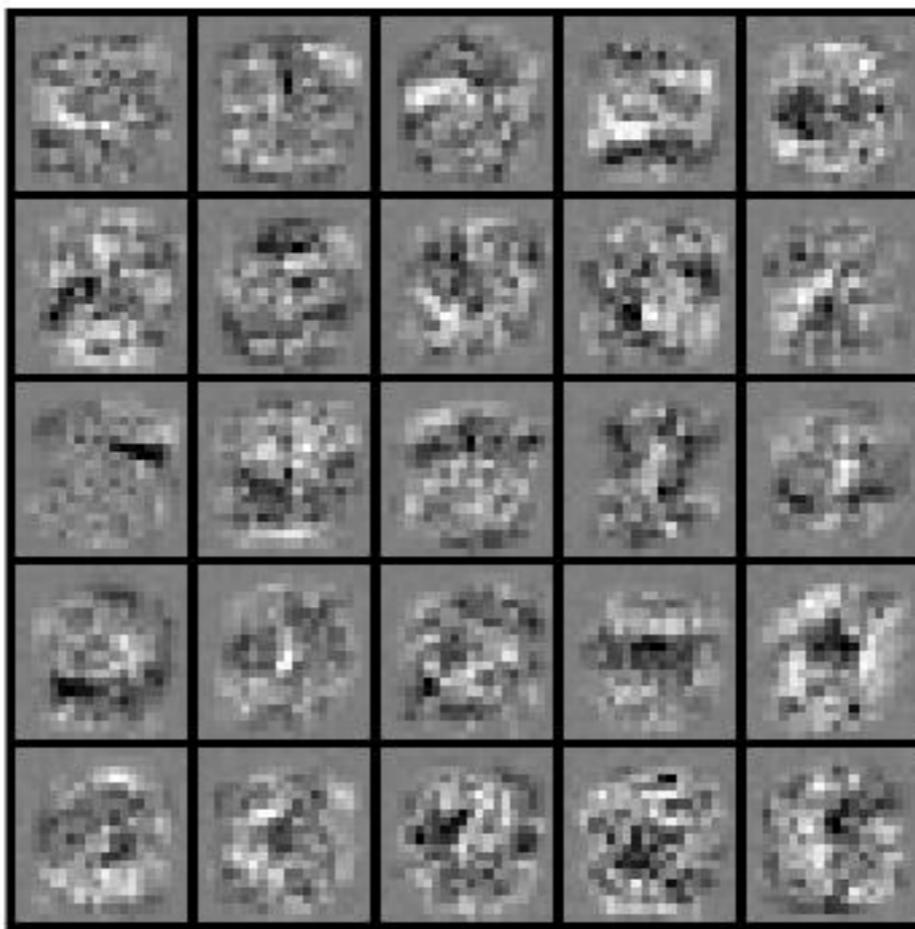
NEURAL NET EXAMPLE – DIGIT RECOGNITION

Weights



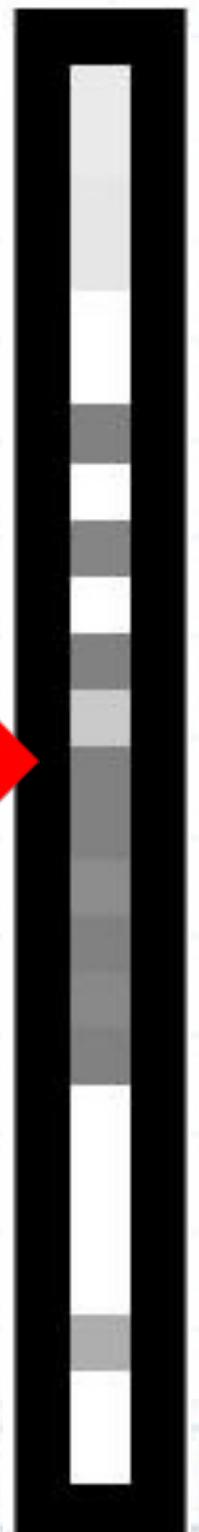
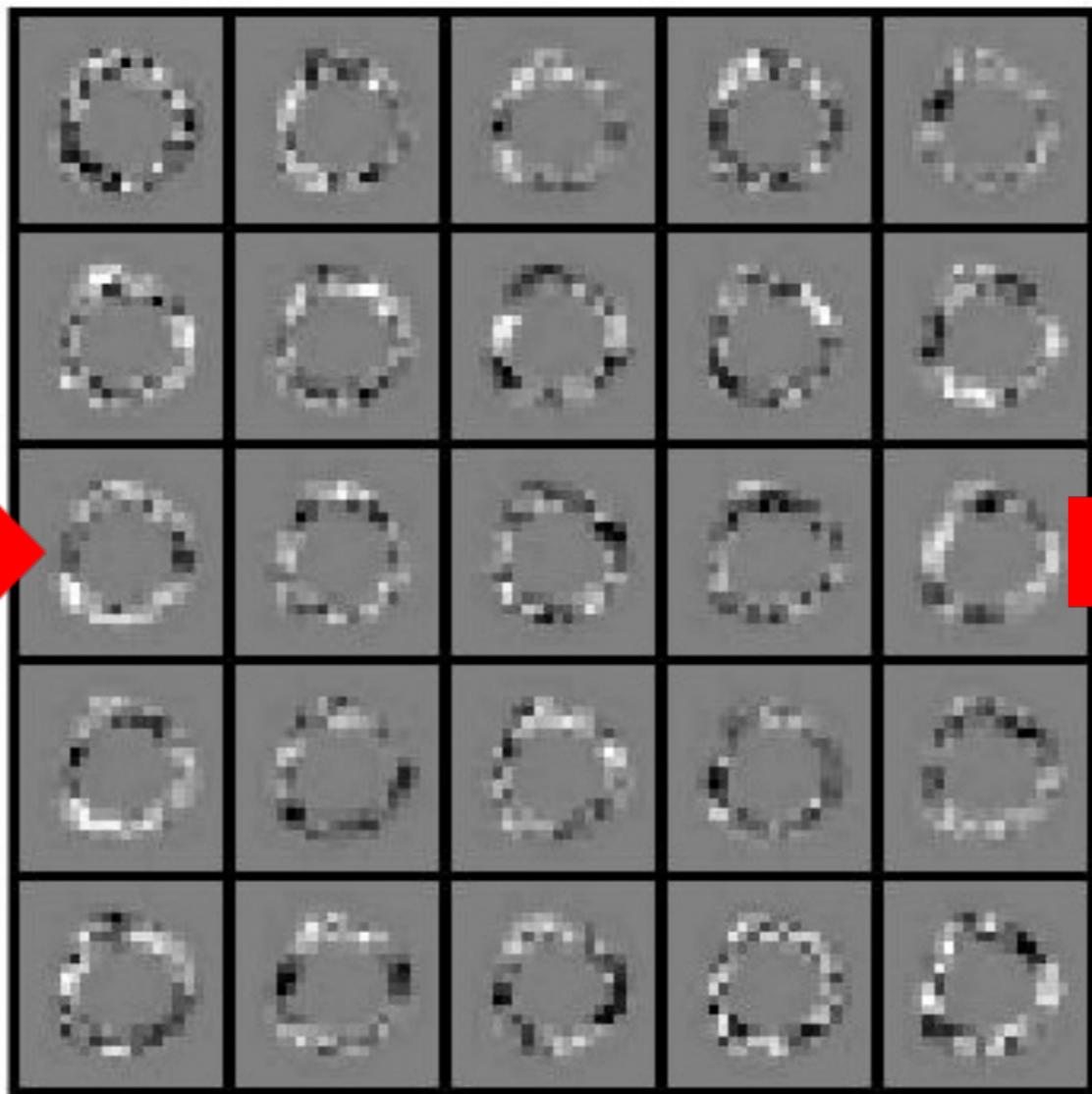
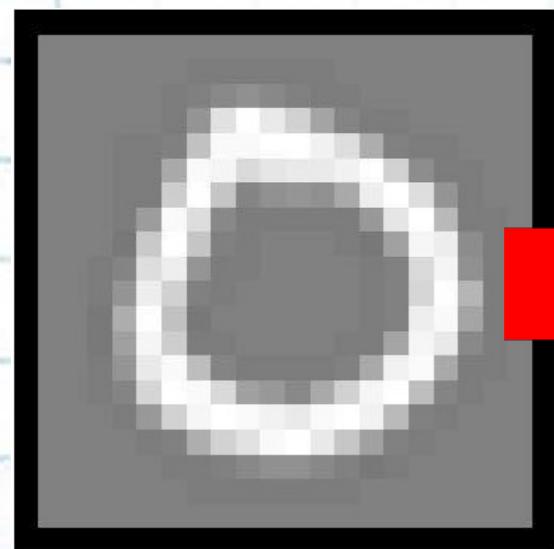
NEURAL NET EXAMPLE – DIGIT RECOGNITION

Weights to hidden units –
“feature” extraction



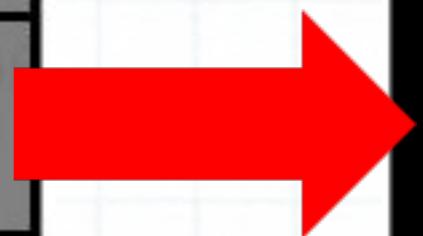
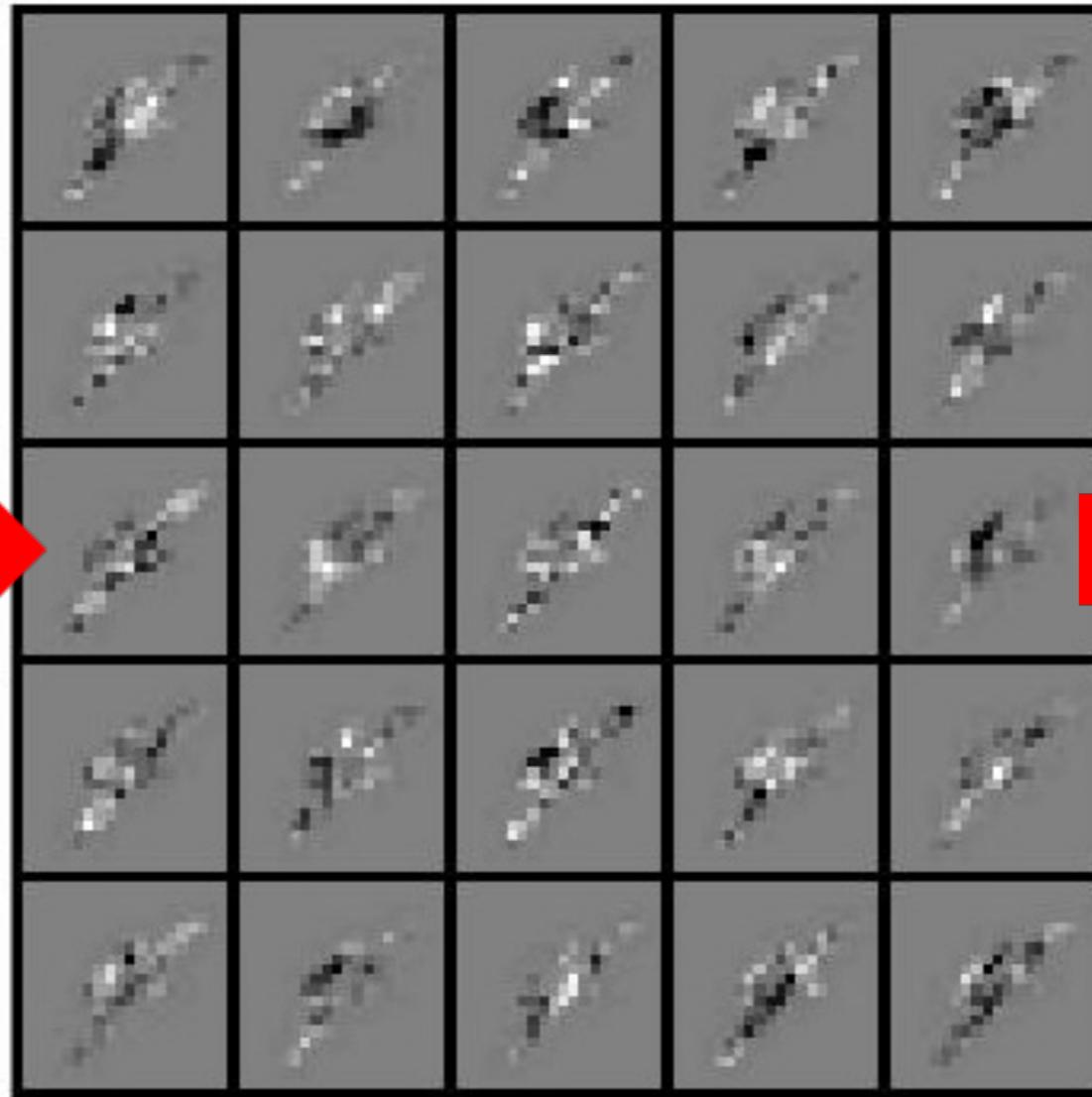
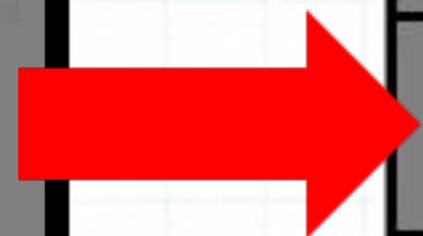
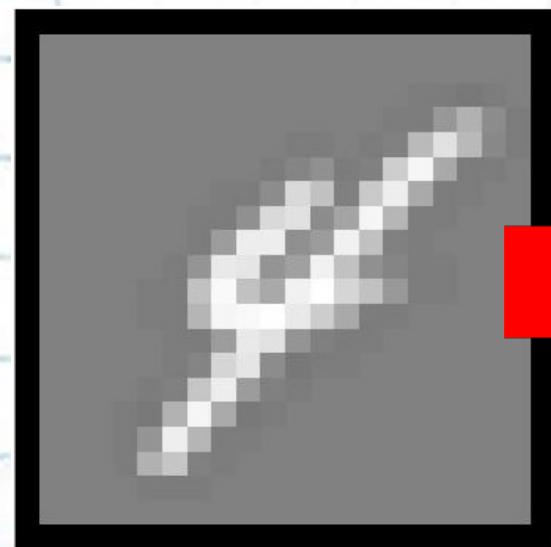
NEURAL NET EXAMPLE – DIGIT RECOGNITION

Weights to hidden units



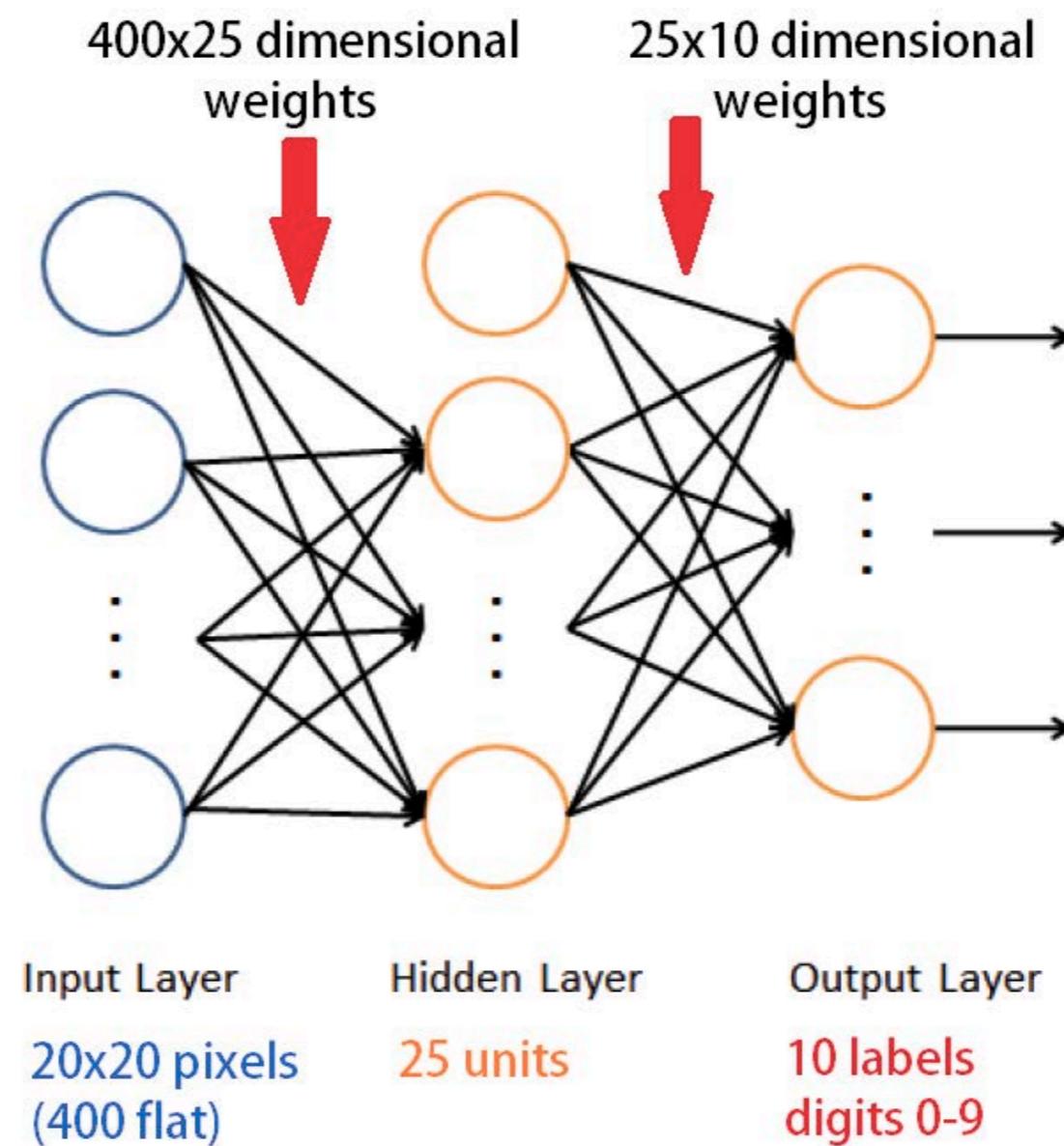
NEURAL NET EXAMPLE – DIGIT RECOGNITION

Weights to hidden units



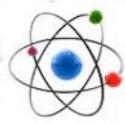
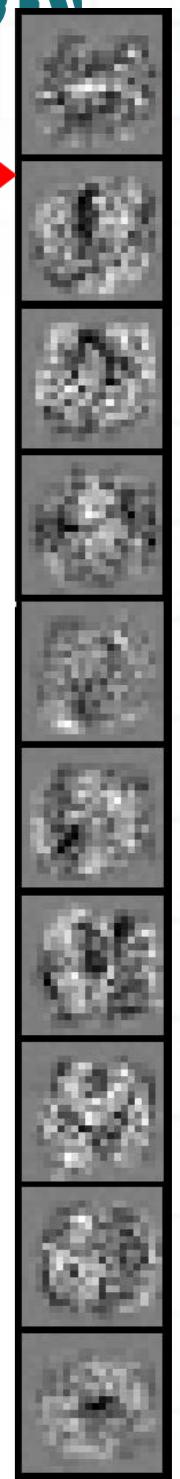
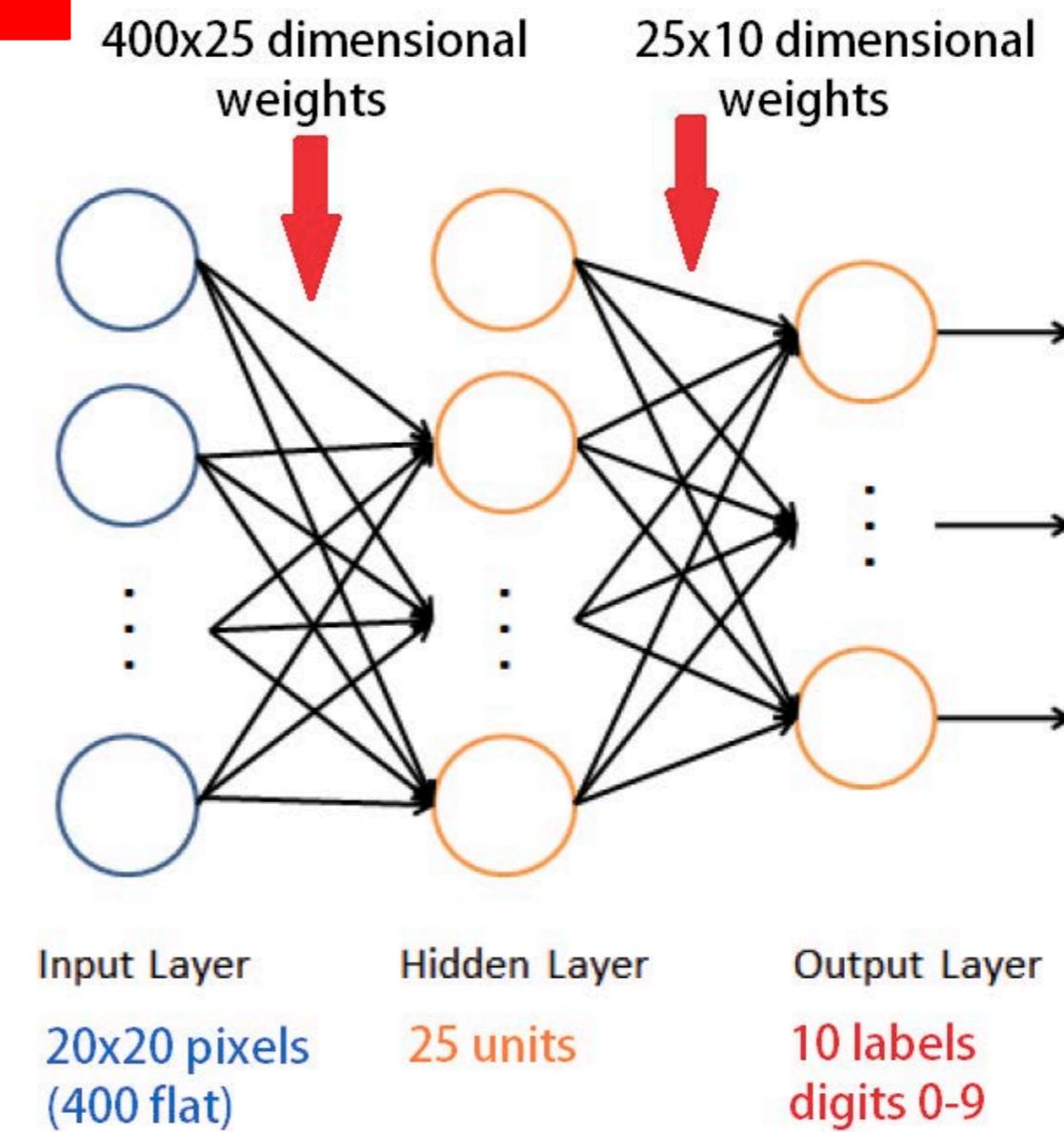
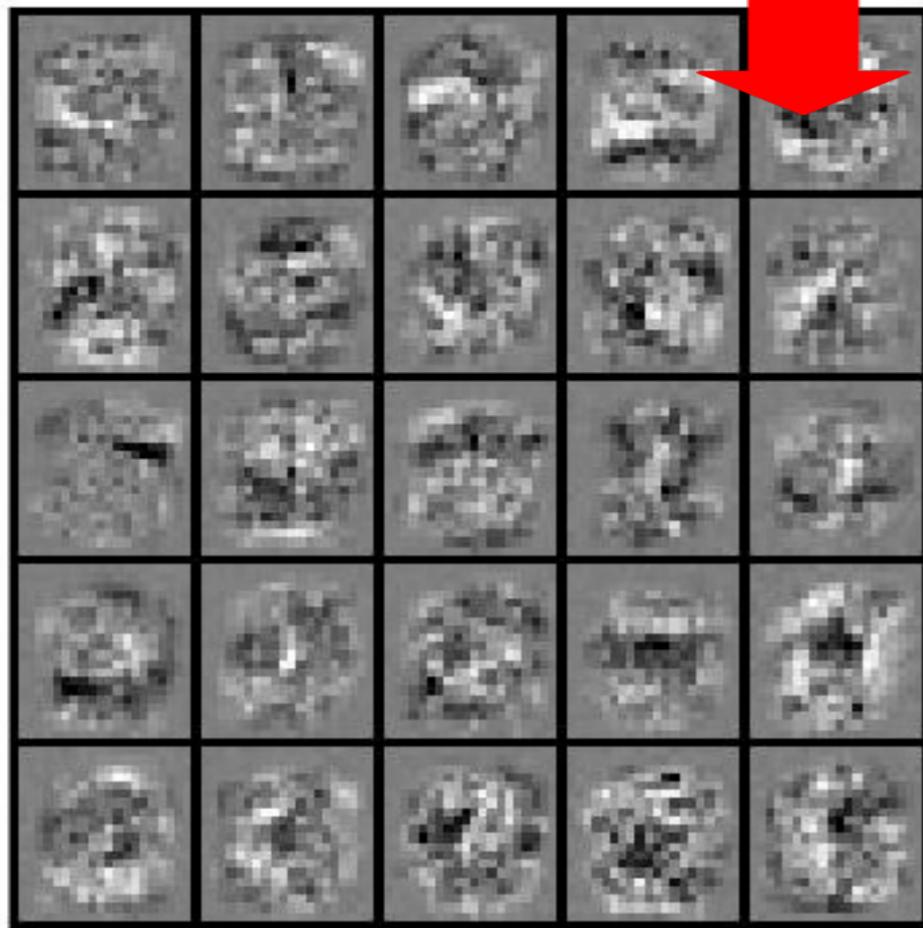
NEURAL NET EXAMPLE – DIGIT RECOGNITION

Hidden units to output



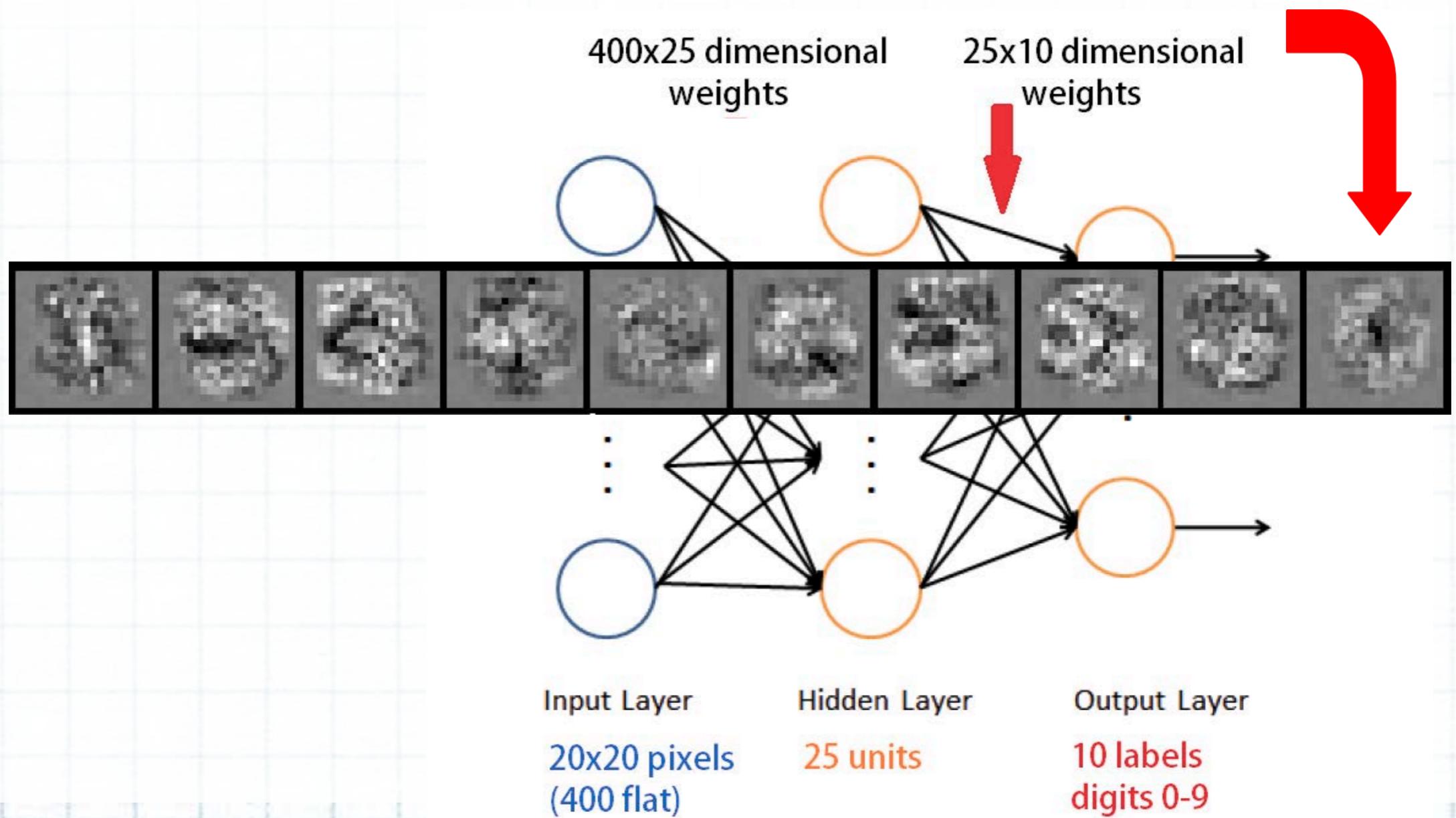
NEURAL NET EXAMPLE – DIGIT RECOGNITION

Hidden units to output

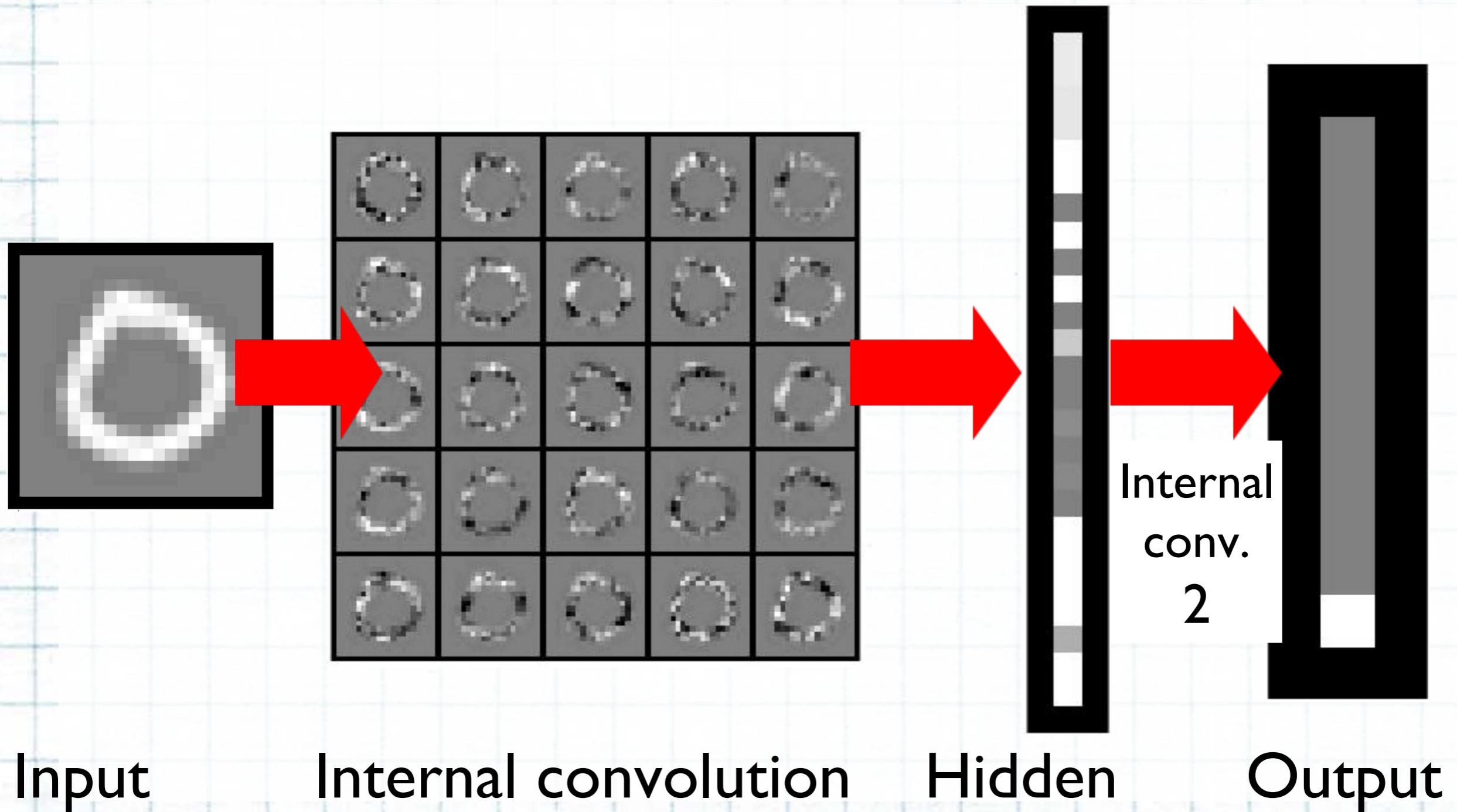


NEURAL NET EXAMPLE – DIGIT RECOGNITION

Hidden units to output

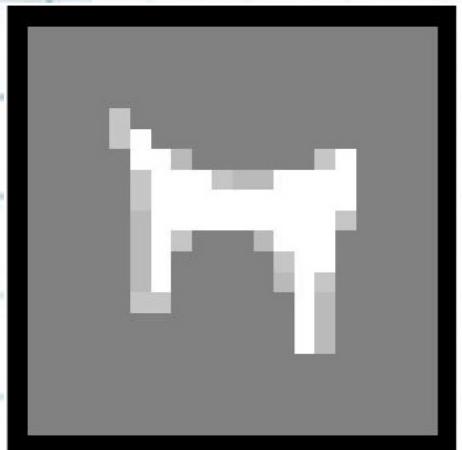


NEURAL NET EXAMPLE – DIGIT RECOGNITION



NEURAL.NET EXAMPLE – DIGIT RECOGNITION

What happens with unknowns?
Klingon 6 [jav]



Input

Internal convolution

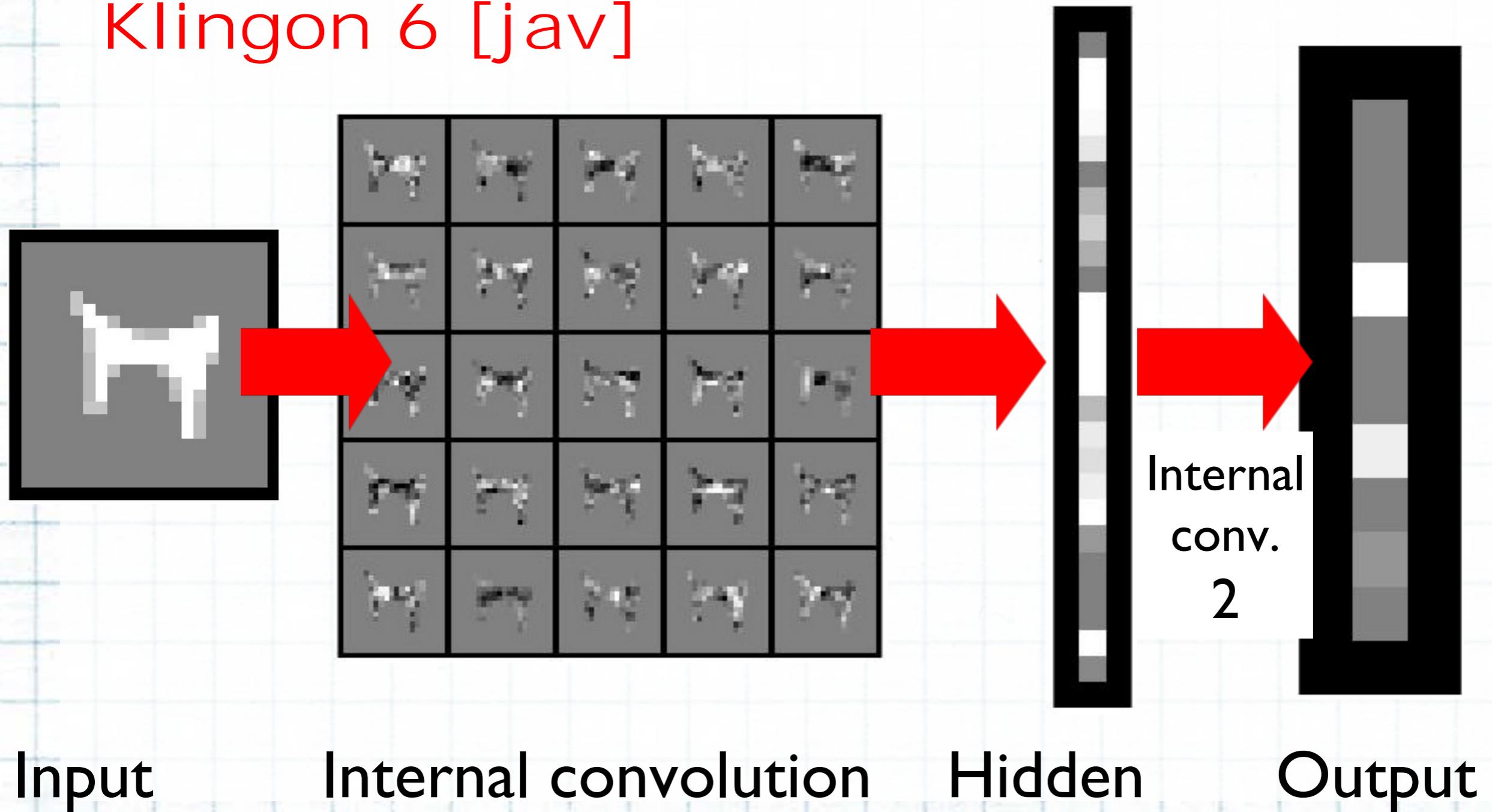
Hidden

Output



NEURAL NET EXAMPLE – DIGIT RECOGNITION

Klingon 6 [jav]



Input

Internal convolution

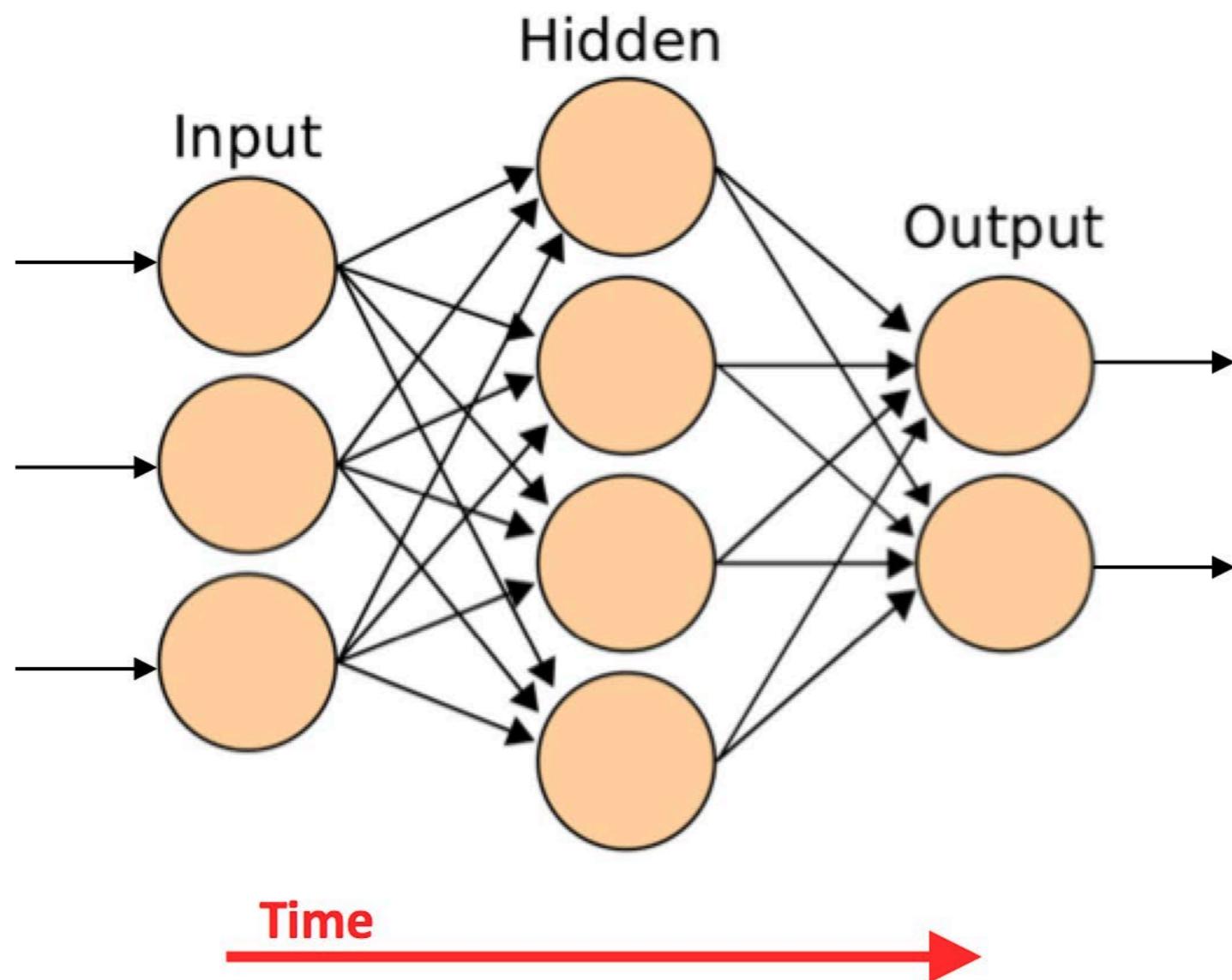
Hidden

Output



ARTIFICIAL NEURAL NETWORK

Feed-forward neural net



Adapted from: Cburnett, https://commons.wikimedia.org/wiki/File:Artificial_neural_network.svg



HOW DO WORDS GET INTO AND OUT OF THE NETWORK?

.Challenges for NMT:

- Input and output length not fixed, different sentence ordering in source and target languages
- Sentence alignment: how to align words between two different languages
- Context
- Training metrics



HOW DO WORDS GET INTO AND OUT OF THE NETWORK?

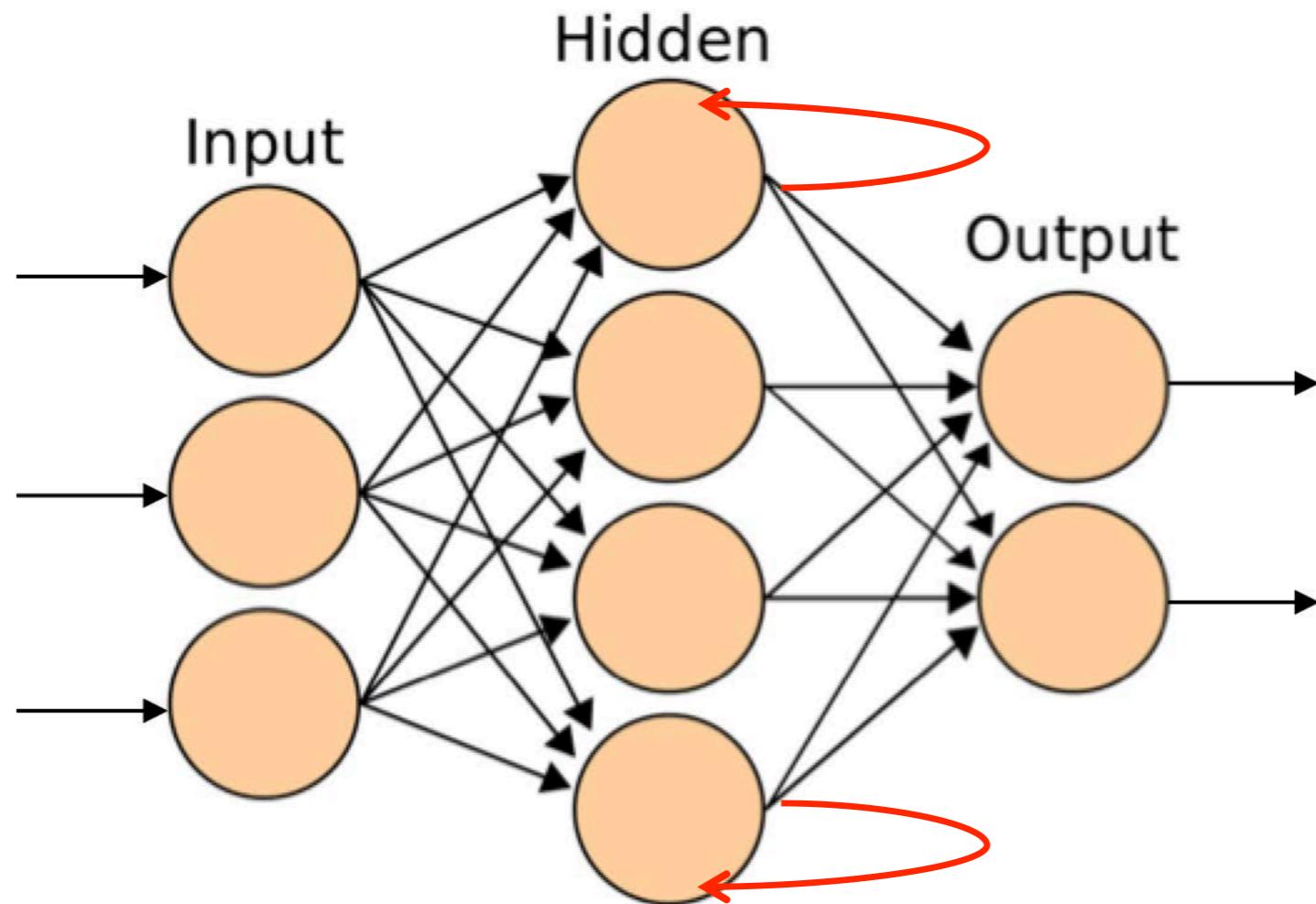
.Challenges for NMT:

- Input and output length not fixed, different sentence ordering in source and target languages =>
Use recurrent neural networks with attention or convolutional networks
- Context => document (or at least paragraph) level, not sentence level
- Training metrics



How Do Words Get Into AND Out OF THE NETWORK?

Recurrent neural net

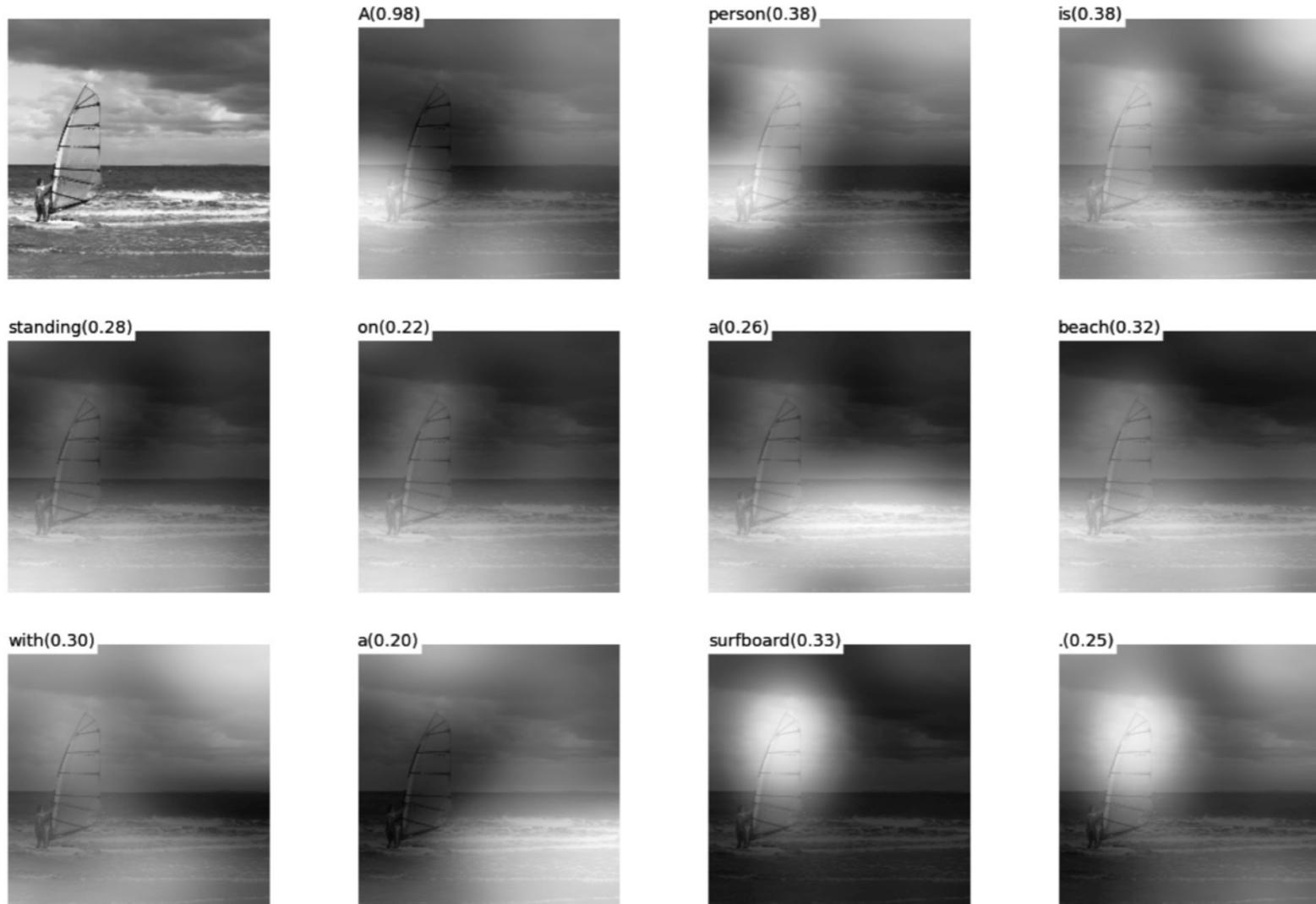


Adapted from: Cburnett, https://commons.wikimedia.org/wiki/File:Artificial_neural_network.svg



HOW DO WORDS GET INTO AND OUT OF THE NETWORK?

Attention mechanism



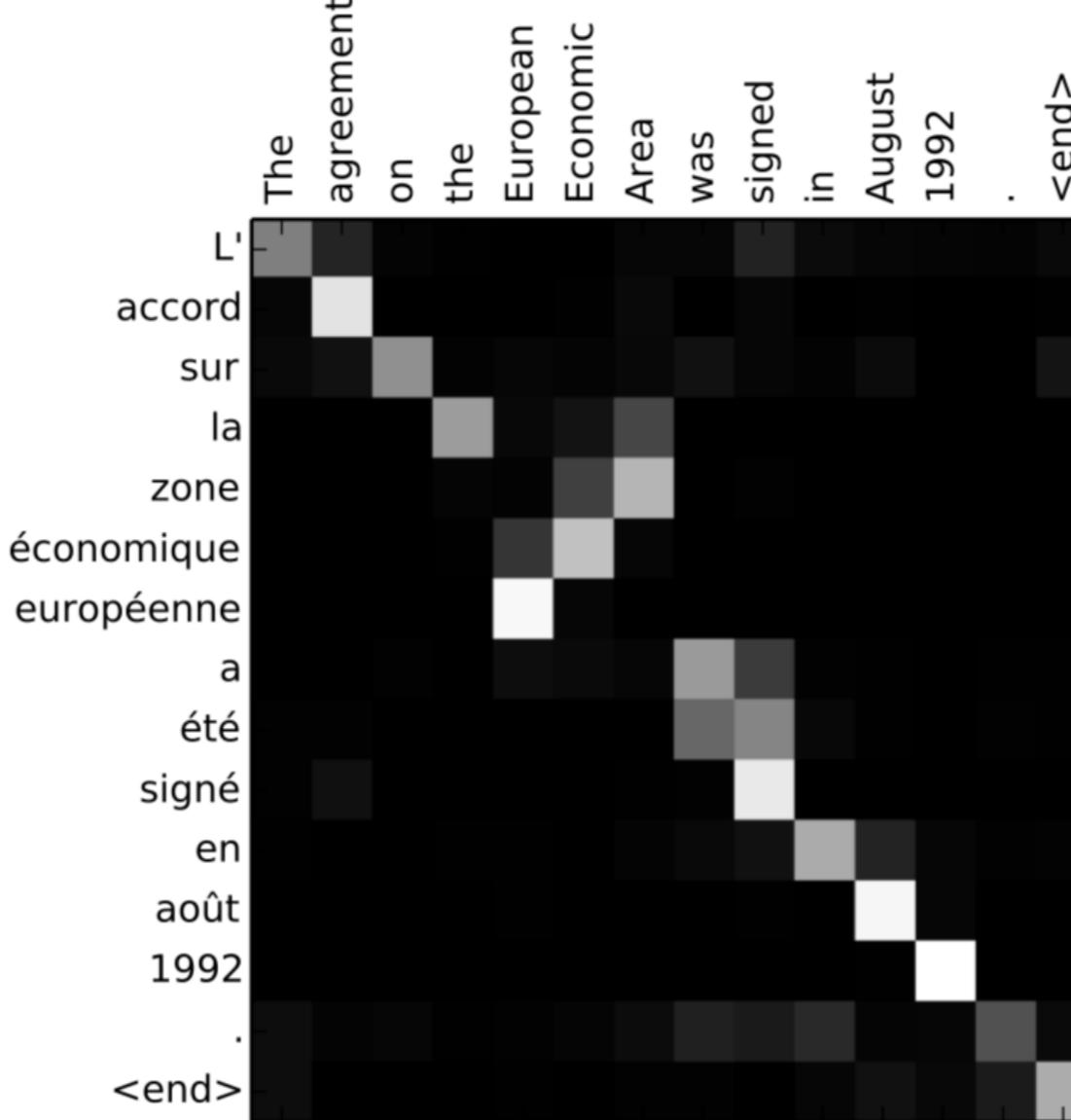
(b) A person is standing on a beach with a surfboard.

Source: K. Xu et al, <https://arxiv.org/abs/1502.03044>



HOW DO WORDS GET INTO AND OUT OF THE NETWORK?

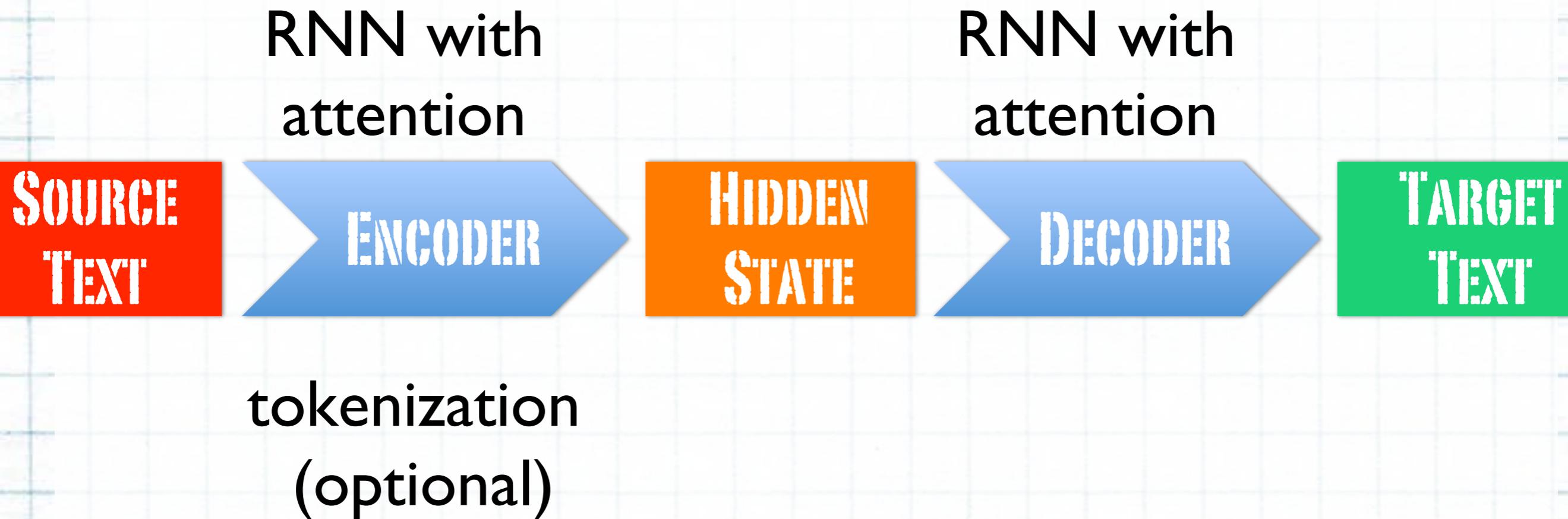
Attention mechanism



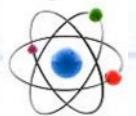
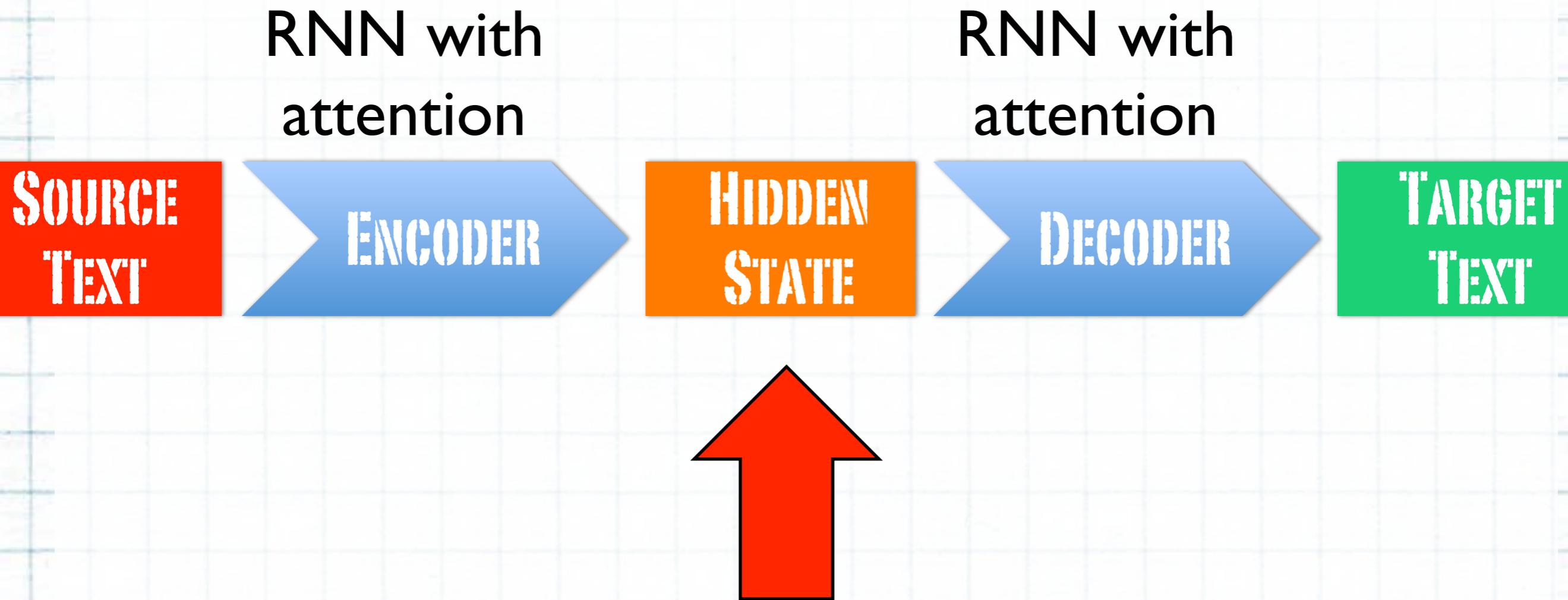
Source: D. Bahdanau et al, ICLR conference proceedings, <https://arxiv.org/pdf/1409.0473.pdf>



HOW DO WORDS GET INTO AND OUT OF THE NETWORK?

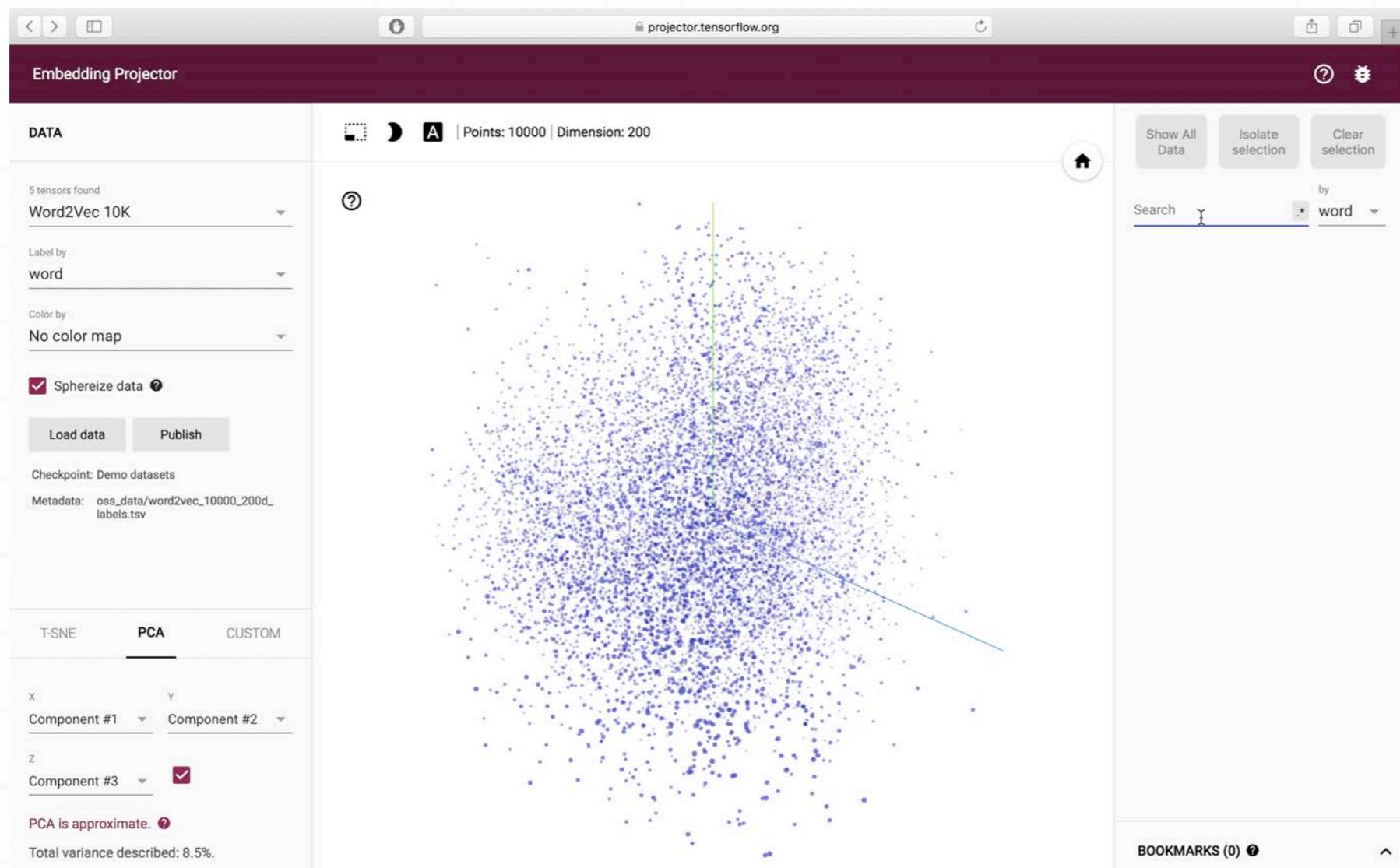


How Do Words Get Into and Out of the Network?



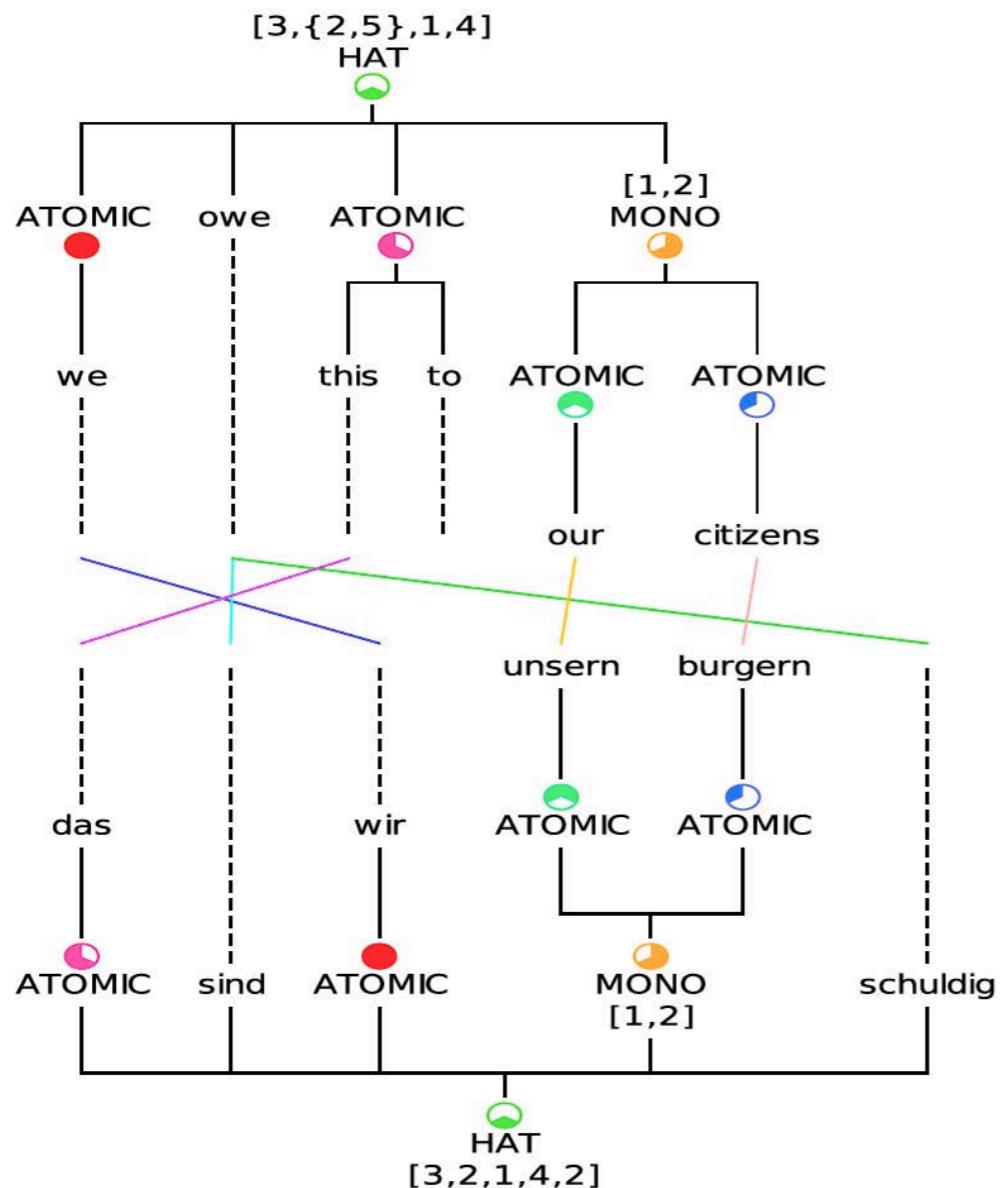
How Do Words Get Into and Out of the Network?

<https://projector.tensorflow.org>



HOW DO WORDS GET INTO AND OUT OF THE NETWORK?

Recall: SMT



From G. M. de Buy Wenninger, K. Sima'an,
PBML No. 101, April 2014, pp. 43



NEURAL NETS - RECAP

- ✓ Training = extraction of "features" (=patterns) from training data



NEURAL NETS - RECAP

- ✓ Training = extraction of "features" (=patterns) from training data
- ✓ The more hidden layers and hidden units, the more parameters (possible overfitting!)



NEURAL NETS - RECAP

- ✓ Training = extraction of "features" (=patterns) from training data
- ✓ The more hidden layers and hidden units, the more parameters (possible overfitting!)
- ✓ Beware: Garbage in -> worse garbage out!



NEURAL NETS - RECAP

- ✓ Training = extraction of "features" (=patterns) from training data
- ✓ The more hidden layers and hidden units, the more parameters (possible overfitting!)
- ✓ Beware: Garbage in -> worse garbage out!
- ✓ ANNs work well for pattern recognition after training, including "context"

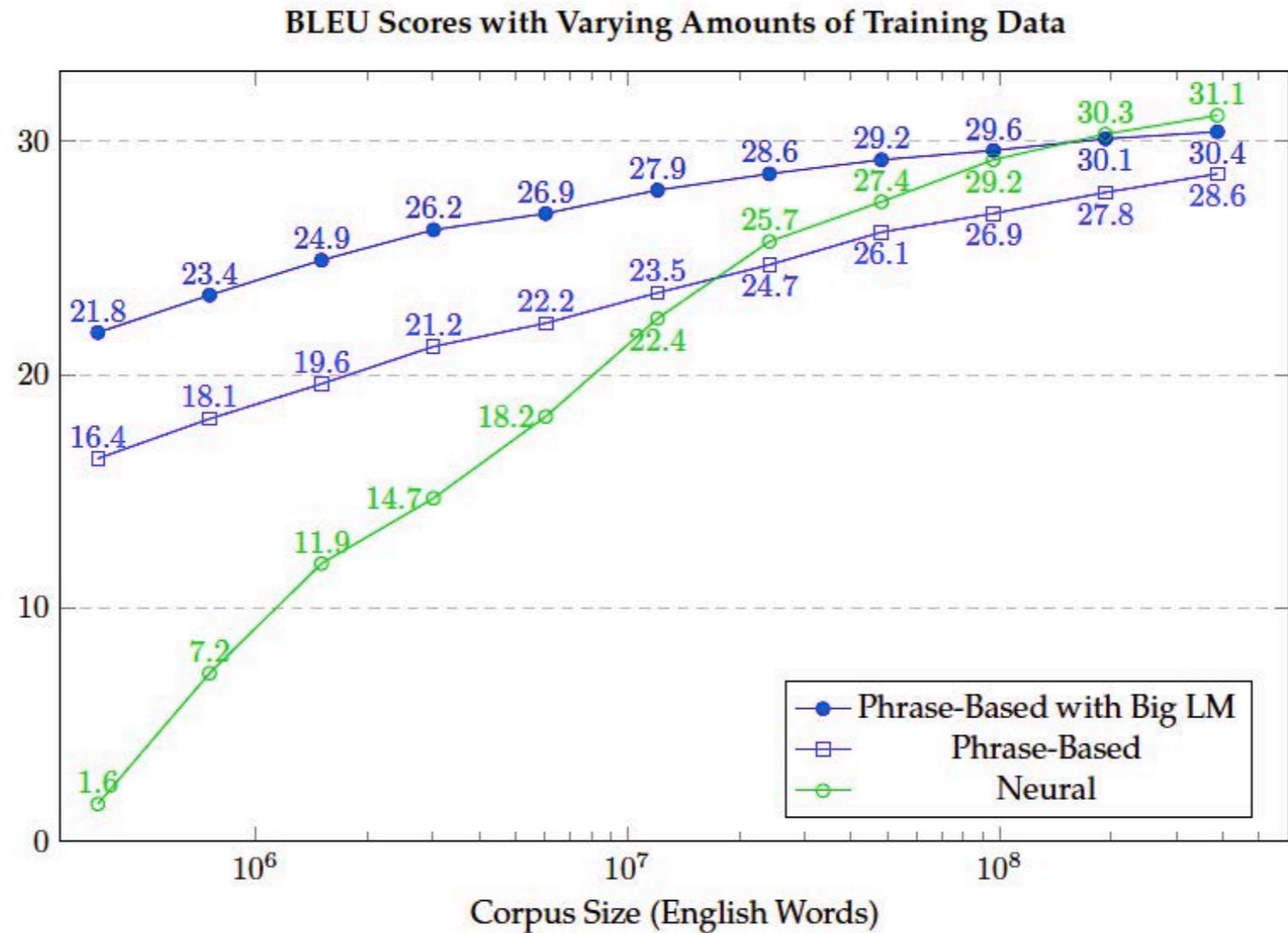


NEURAL NETS - RECAP

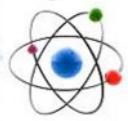
- ✓ Training = extraction of "features" (=patterns) from training data
- ✓ The more hidden layers and hidden units, the more parameters (possible overfitting!)
- ✓ Beware: Garbage in -> worse garbage out!
- ✓ ANNs work well for pattern recognition after training, including "context"
- ✓ Completely unpredictable when confronted with new, hitherto unknown data



TRAINING DATA



Source: P. Koehn, <https://arxiv.org/abs/1709.07809>



UNPREDICTABILITY

English Spanish Maori Detect language ▾



dog dog dog dog dog

X

German English Spanish ▾

Translate

dog dog dog dog



UNPREDICTABILITY

English Spanish Maori Detect language ▾



dog dog dog dog dog dog|



German English Spanish ▾

Translate

dog dog dog - reader email



Suggest an edit



Carola F. Berger, An Introduction to NMT, ATA59

UNPREDICTABILITY

English Spanish Maori Detect language ▾



dog |



German English Spanish ▾

Translate

Doomsday Clock is three minutes at twelve We are experiencing characters and a dramatic developments in the world



UNPREDICTABILITY

English Spanish

Maori

Detect language



dog dog



German English

Spanish

Translate

Doomsday Clock is three minutes at twelve We are experiencing characters and a dramatic developments in the world, which indicate that we are approaching the end times and Jesus' return



UNPREDICTABILITY

English Spanish Maori Detect language ▾

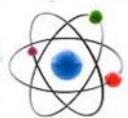


dog dog| X

German English Spanish ▾

Translate

krist dog dog dog dog dog dog dog dog dog dog



REFERENCES & FURTHER READING

- Slides at: <https://www.CFBtranslations.com>
- Handouts in app and also at
<https://www.CFBtranslations.com>
- A. Ng, *Machine Learning*, Coursera,
<https://www.coursera.org/learn/machine-learning>
- Google's Tensorflow: <https://www.tensorflow.org/>
- Madly Ambiguous - game to illustrate how NMT deals with context: <http://madlyambiguous.osu.edu>

