

# Deep Learning

Philosophy, Working Principles and Algorithms

Presented by

**Fayyaz Minhas**

Department of Computer Science  
University of Warwick



@fayyazhere



# Philosophy



## Train PhD students to be thinkers not just specialists

*Many doctoral curricula aim to produce narrowly focused researchers rather than critical thinkers. That can and must change, says Gundula Bosch.*

Under pressure to turn out productive lab members quickly, many PhD programmes in the biomedical sciences have shortened their courses, squeezing out opportunities for putting research into its wider context. Consequently, most PhD curricula are unlikely to nurture the big thinkers and creative problem-solvers that society needs.

That means students are taught every detail of a microbe's life cycle but little about the life scientific. They need to be taught to recognize how errors can occur. Trainees should evaluate case studies derived from flawed real research, or use interdisciplinary detective games to find logical fallacies in the literature. Above all, students must be shown the scientific process as it is — with its limitations and potential pitfalls as well as its fun side, such as serendipitous discoveries and hilarious blunders.

This is exactly the gap that I am trying to fill at Johns Hopkins University in Baltimore, Maryland, where a new graduate science programme is entering its second year. Microbiologist Arturo Casadevall and I began pushing for reform in early 2015, citing the need to put the philosophy back into the doctorate of philosophy: that is, the 'Ph' back into the PhD. We call our programme R3, which means that our students learn to apply rigour to their design and conduct of experiments; view their work through the lens of social responsibility; and to think critically, communicate better, and thus improve reproducibility. Although we are aware of many innovative individual courses developed along these lines, we are striving for more-comprehensive reform.

Our offerings are different from others at the graduate level. We have critical-thinking assignments in which students analyse errors in reasoning in a *New York Times* opinion piece about 'big sugar', and the ethical implications of the arguments made in a *New Yorker* piece by surgeon Atul Gawande entitled 'The Mistrust of Science'. Our courses on rigorous research, scientific integrity, logic, and mathematical and programming skills are integrated into students' laboratory and fieldwork. Those studying the influenza virus, for example, work with real-life patient data sets and wrestle with the challenges of applied statistics.

A new curriculum starts by winning allies. Both students and faculty members must see value in moving off the standard track. We used informal interviews and focus groups to identify areas in which students and faculty members saw gaps in their training. Recurring themes included the inability to apply theoretical knowledge in statistical tests in the laboratory, frequent mistakes in choosing an appropriate set of experimental controls, and significant difficulty in explaining work to non-experts.

Introducing our programme to colleagues in the Johns Hopkins life-sciences departments was even more sensitive. I was startled by the oft-expressed opinion that scientific productivity depended more

on rote knowledge than on competence in critical thinking. Several principal investigators were uneasy about students committing more time to less conventional forms of education. The best way to gain their support was coffee: we repeatedly met lab heads to understand their concerns.

With the pilot so new, we could not provide data on students' performance, but we could address faculty members' scepticism. Some colleagues were apprehensive that students would take fewer courses in specialized content to make room for interdisciplinary courses on ethics, epistemology and quantitative skills. In particular, they worried that the R3 programme could lengthen the time required for students to complete their degree, leave them insufficiently knowledgeable in their subject areas and make them less productive in the lab.

We made the case that better critical thinking and fewer mandatory discipline-specific classes might actually position students to be more productive. We convinced several professors to try the new system and participate in structured evaluations on whether R3 courses contributed to students' performance.

So far, we have built 5 new courses from scratch and have enrolled 85 students from nearly a dozen departments and divisions. The courses cover the anatomy of errors and misconduct in scientific practice and teach students how to dissect the scientific literature. An interdisciplinary discussion series encourages broad and critical thinking about science. Our students learn to consider societal consequences of research advances, such

as the ability to genetically alter sperm and eggs.

Discussions about the bigger-picture problems of the scientific enterprise get students to reflect on the limits of science, and where science's ability to do something competes with what scientists should do from a moral point of view. In addition, we have seminars and workshops on professional skills, particularly leadership skills through effective communication, teaching and mentoring.

It is still early days for assessment. So far, however, trainees were repeatedly emphasized that gaining a broader perspective has been helpful. In future, we will collect information about the impact that the R3 approach has on graduates' career choices and achievements.

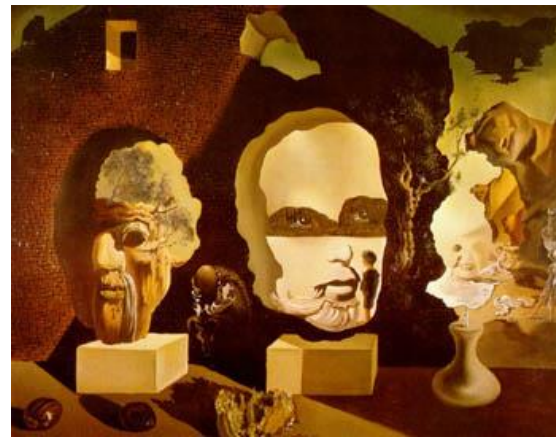
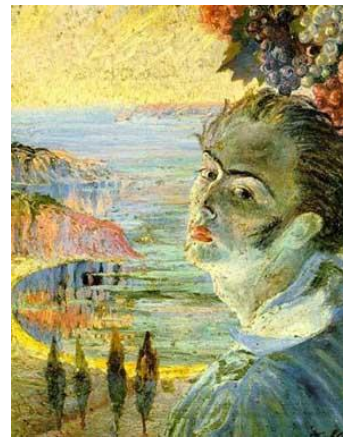
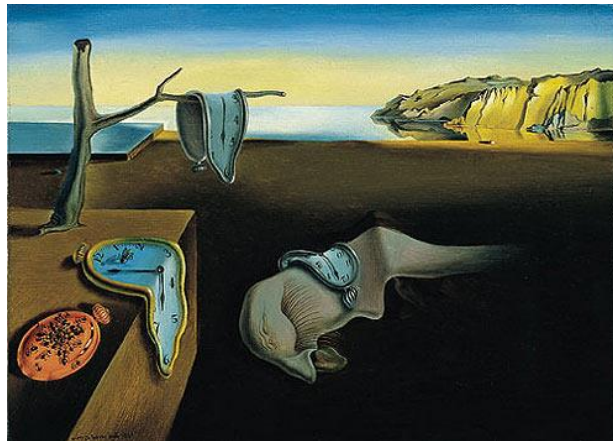
We believe that researchers who are educated more broadly will do science more thoughtfully, with the result that other scientists, and society at large, will be able to rely on this work for a better, more rational world. Science should strive to be self-improving, not just self-correcting. ■

**Gundula Bosch** directs the R3 Graduate Science Initiative at Johns Hopkins Bloomberg School of Public Health in Baltimore, Maryland. e-mail: gbosch@jhu.edu

PUT THE  
PHILOSOPHY  
BACK  
INTO THE  
DOCTORATE  
OF  
PHILOSOPHY.



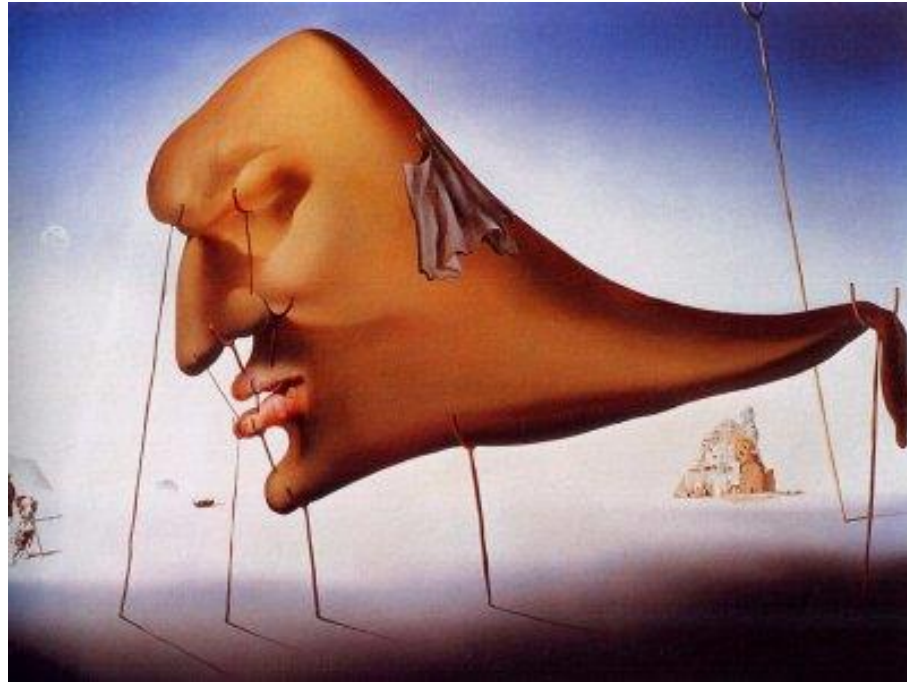
# Paintings by two different painters



Who's painting is this?

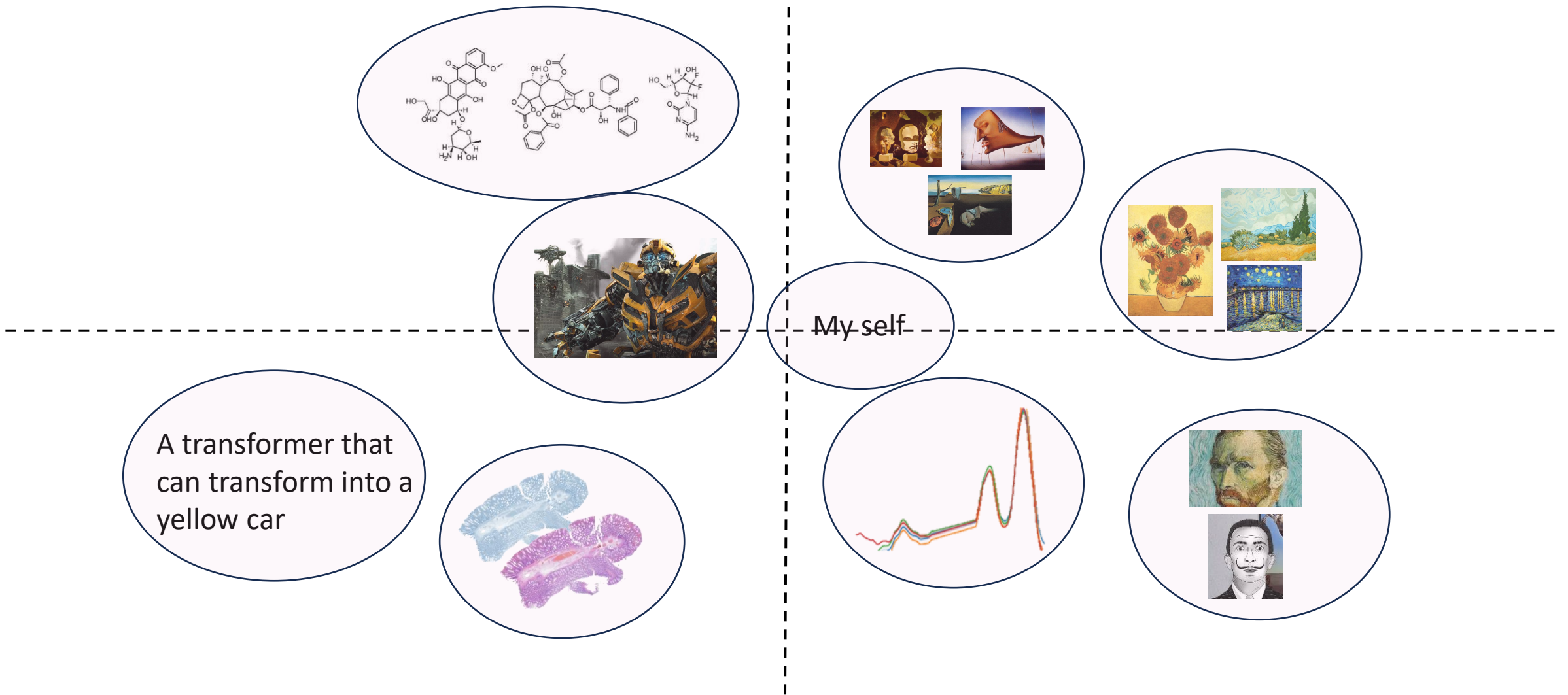


# And this?



*learning from data for generalization to unseen cases*





# I. Entities have (explicit or implicit) representations



“Bank” in which statement is more semantically related to the picture?

- **A:** As he walked by the **bank**, he saw some boats
- **B:** As he walked by the **bank**, he saw some tellers



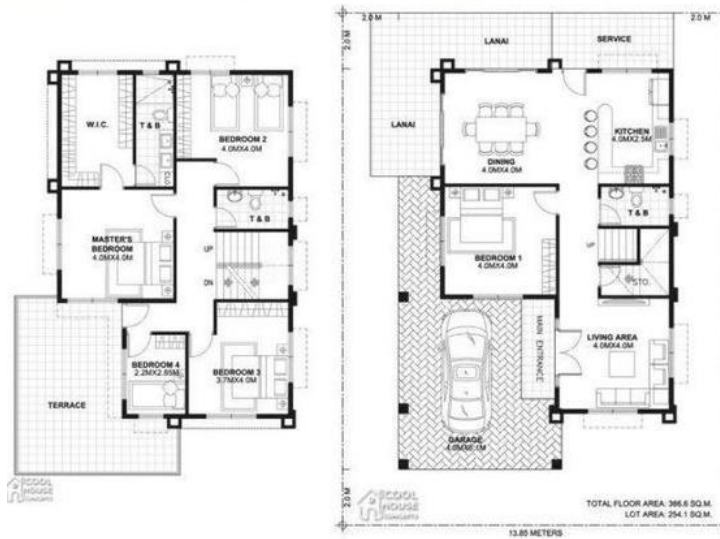
As he walked by the **bank**, he saw some boats



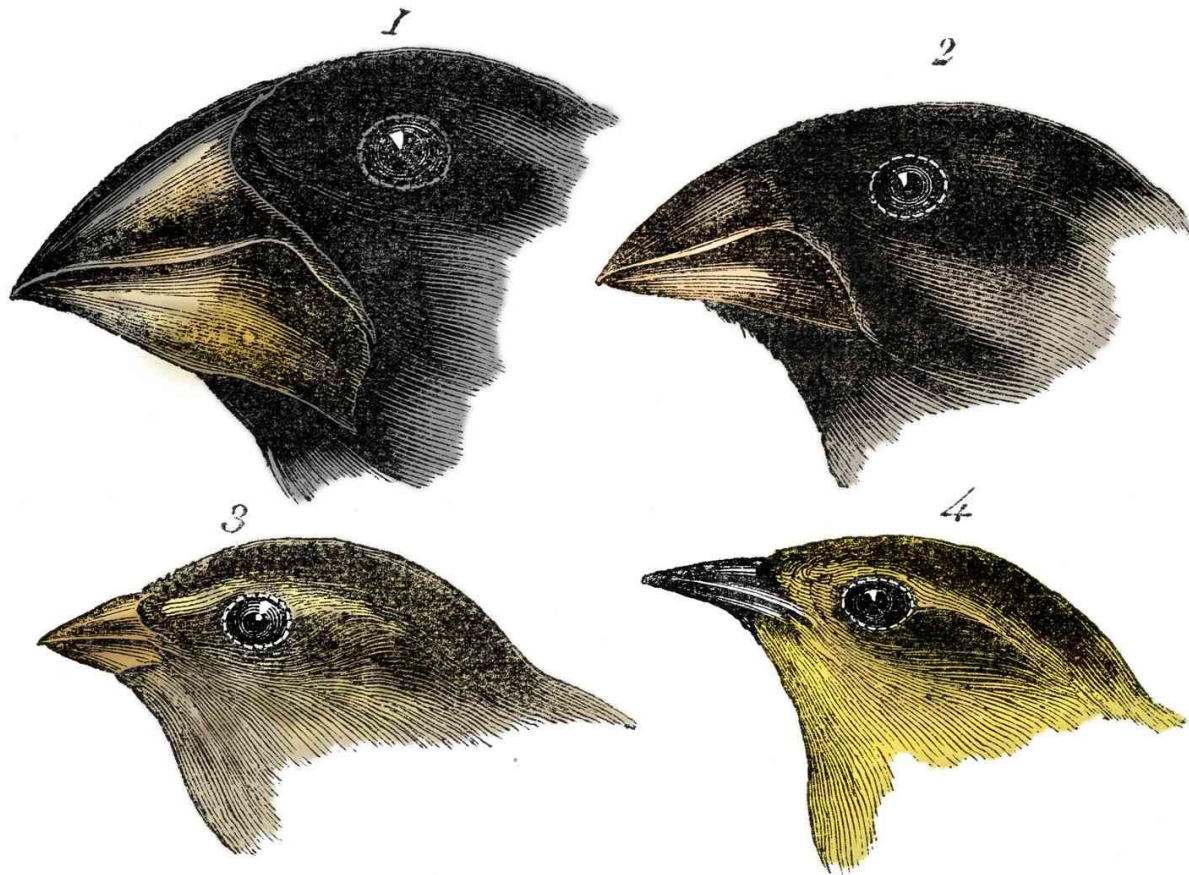
As he walked by the **bank**, he saw some tellers

II. Semantic relatedness of entities is context dependent and thus their representations are contextual





III. Representation of any entity can allow us to reconstruct or “generate” it



1. *Geospiza magnirostris*.  
3. *Geospiza parvula*.

2. *Geospiza fortis*.  
4. *Certhidea olivacea*.



IV. It is possible to develop representations in an inductive manner (through empirical observations)





US Airways Flight 1549



**V. Intelligence is the capacity to develop and utilize causal representations of entities, enabling an organism or system to act effectively and adaptively.**



Only if we could have a mechanism that would enable developing such representations from empirical observations



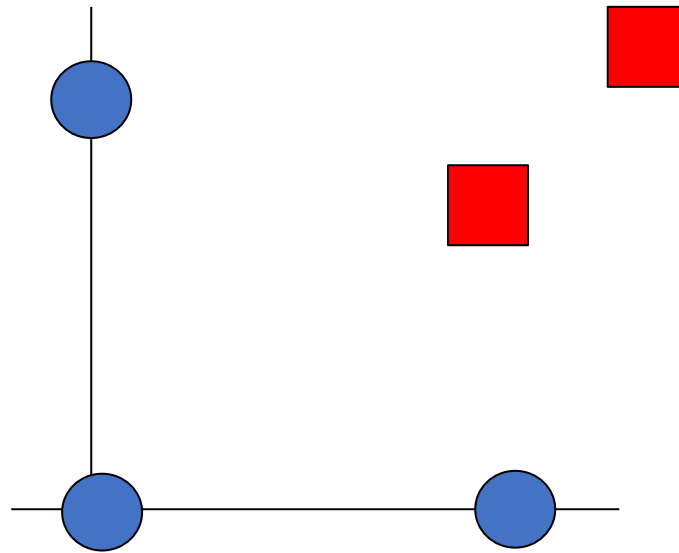
# Deep Learning

Learning Representations from training examples with “layers” of biologically inspired neurons

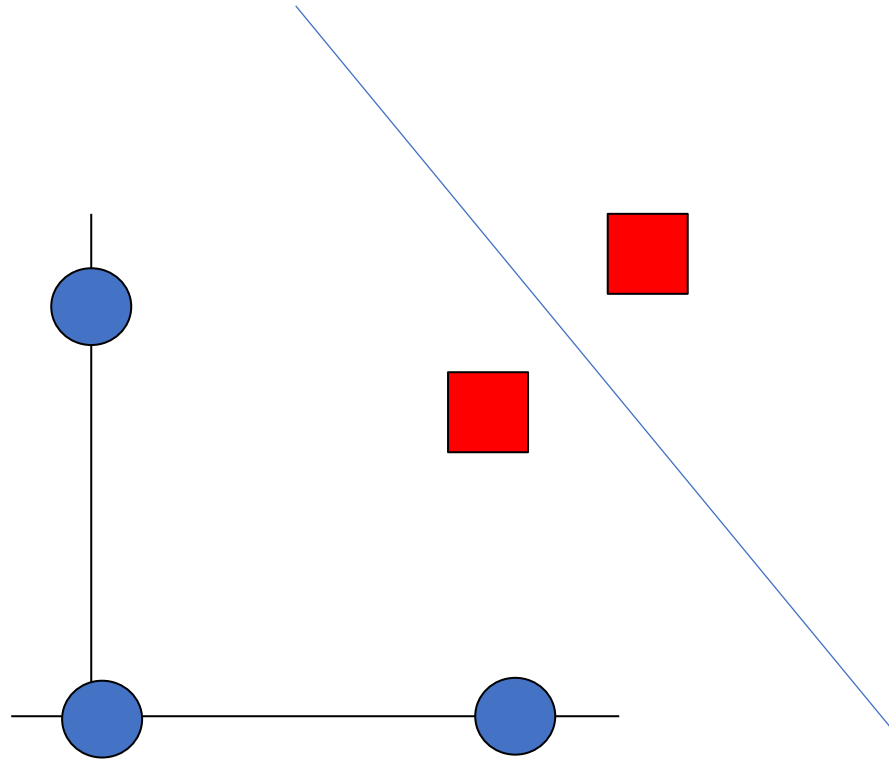
# Working Principles



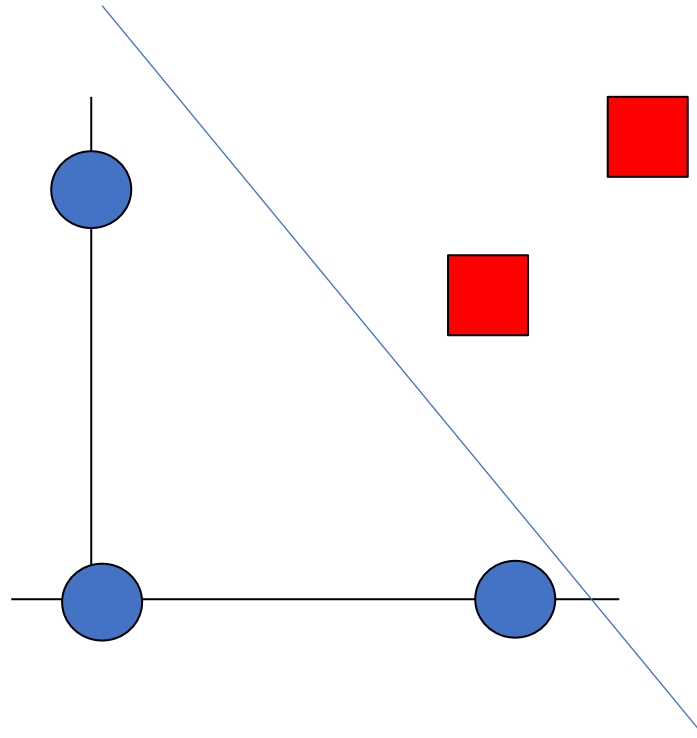
# Exercise



# Exercise

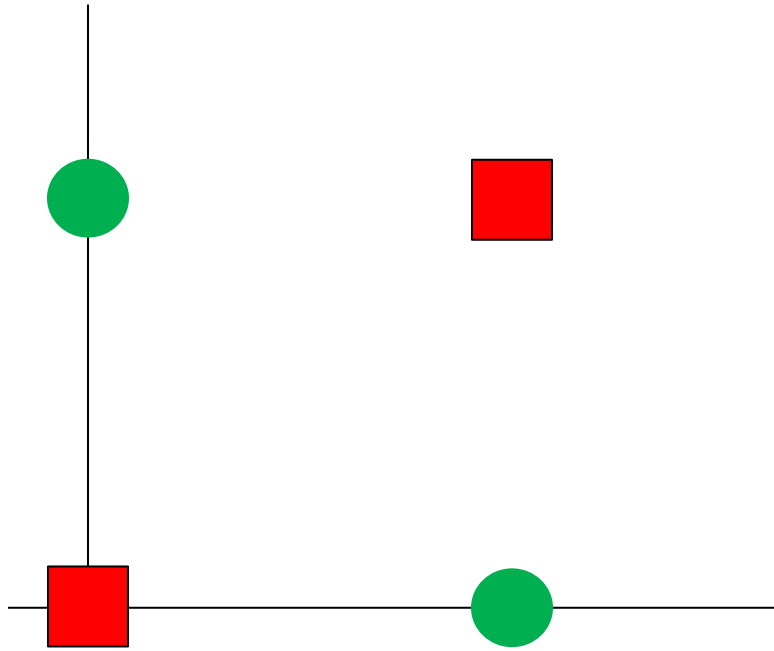


# Exercise

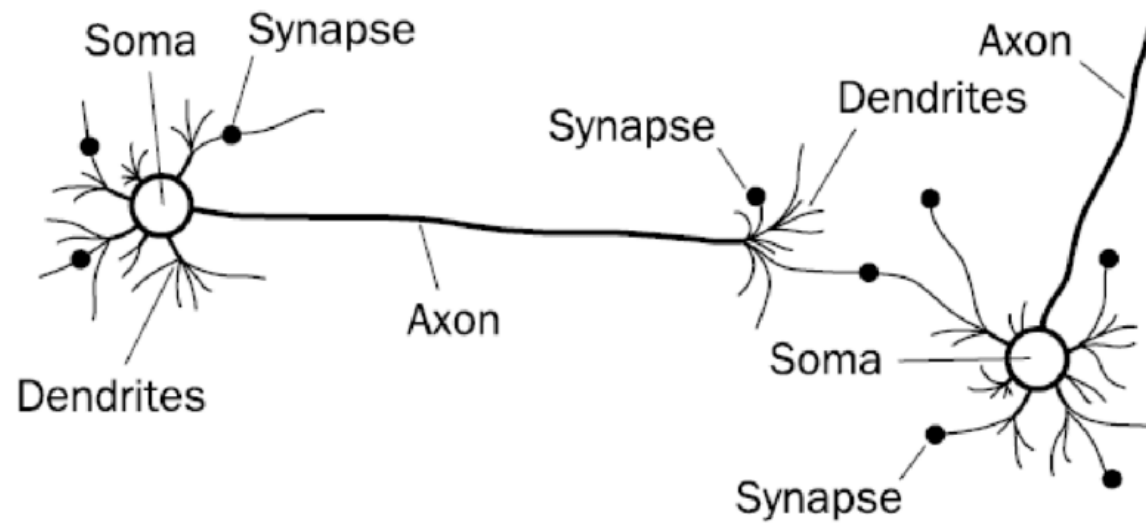




But what if we have a linearly inseparable problems?

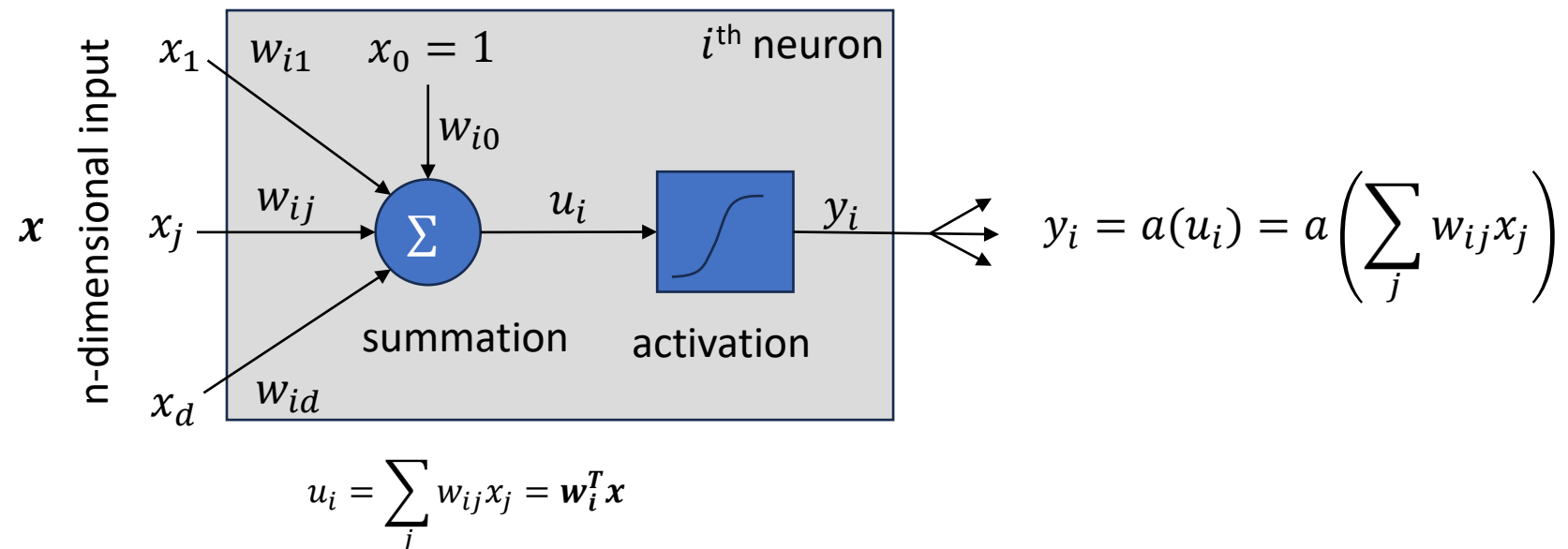
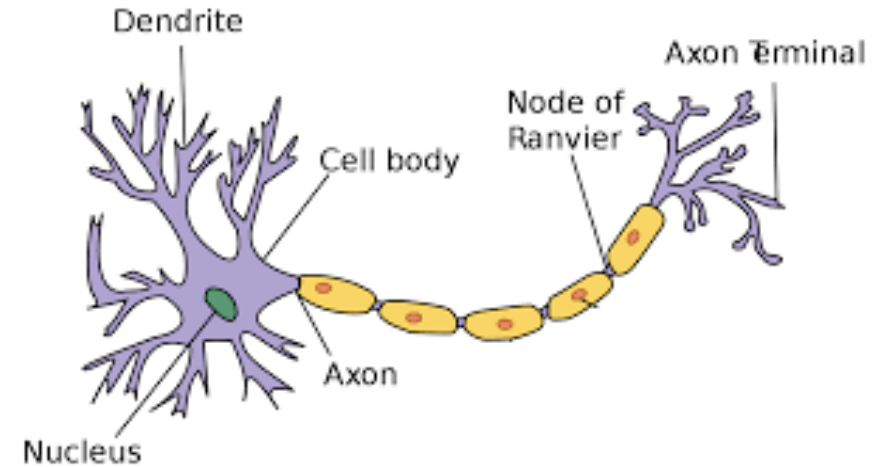


# Biological Neurons and Networks



# Single Neuron: Representation

- An abstraction of the biological neuron



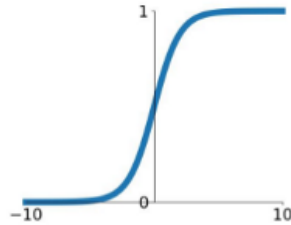
# Activation Functions

Can use any activation function

## Activation Functions

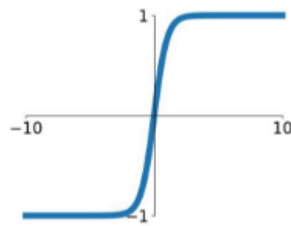
### Sigmoid

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



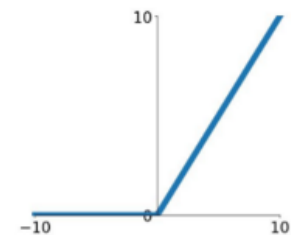
### tanh

$$\tanh(x)$$



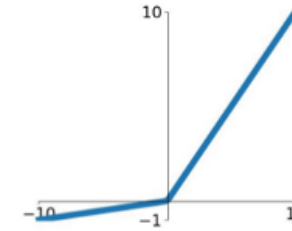
### ReLU

$$\max(0, x)$$



### Leaky ReLU

$$\max(0.1x, x)$$

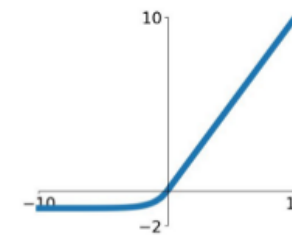


### Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

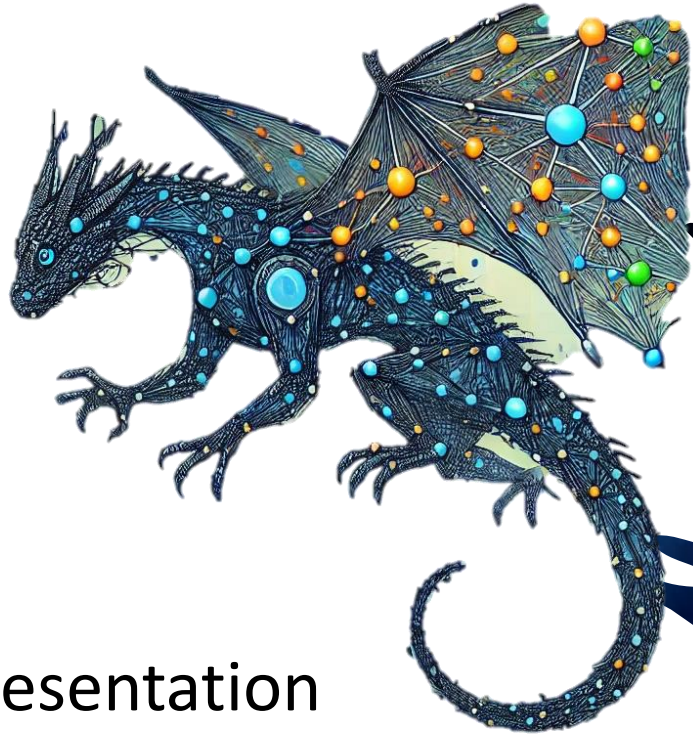
### ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$





# How to train your neural network

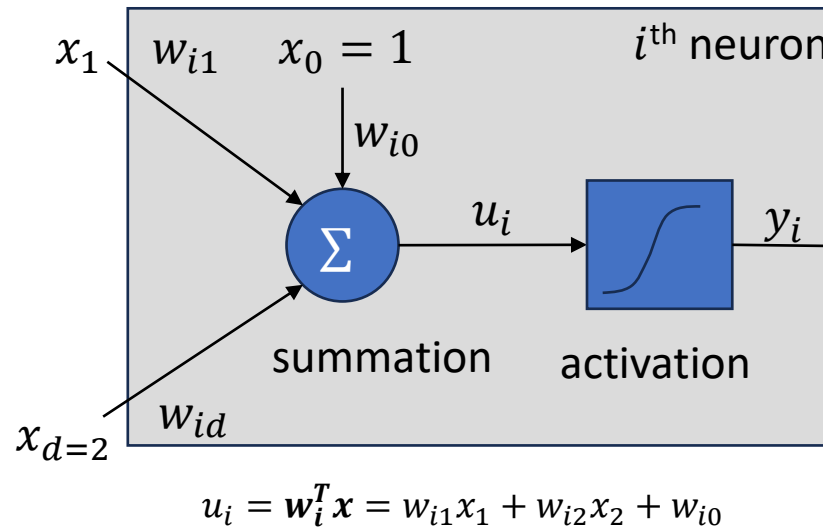
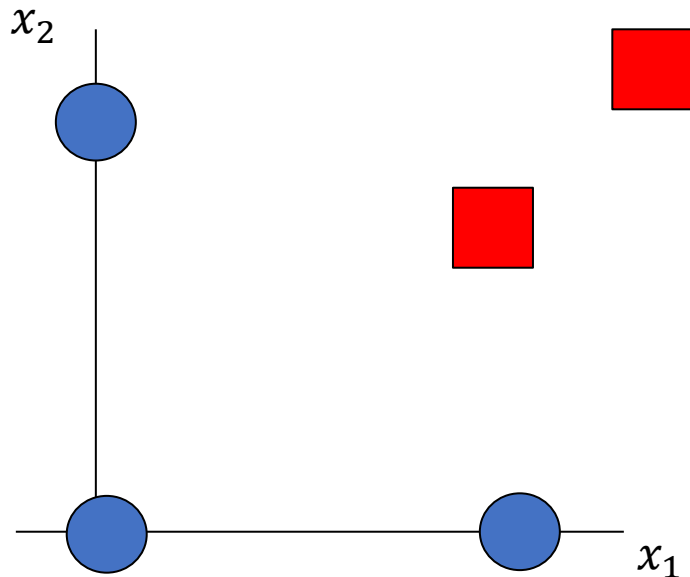


- Representation
- Evaluation
- Optimization





# Representation: How does the model produce its output?



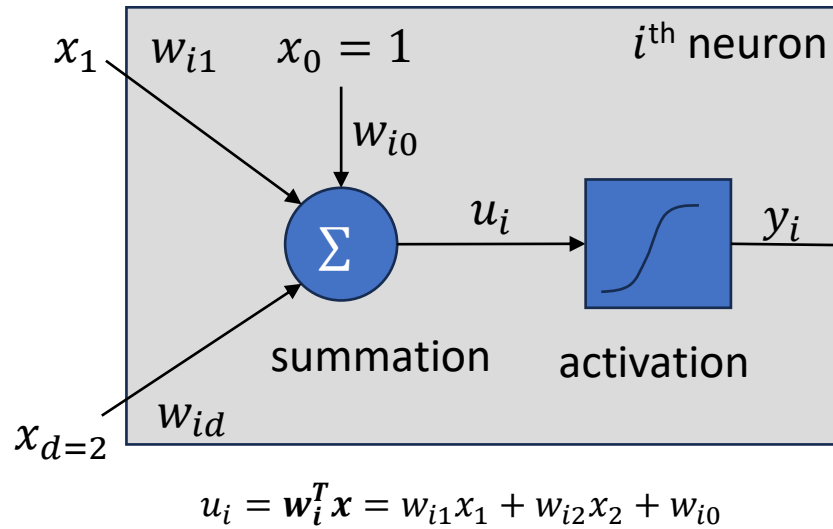
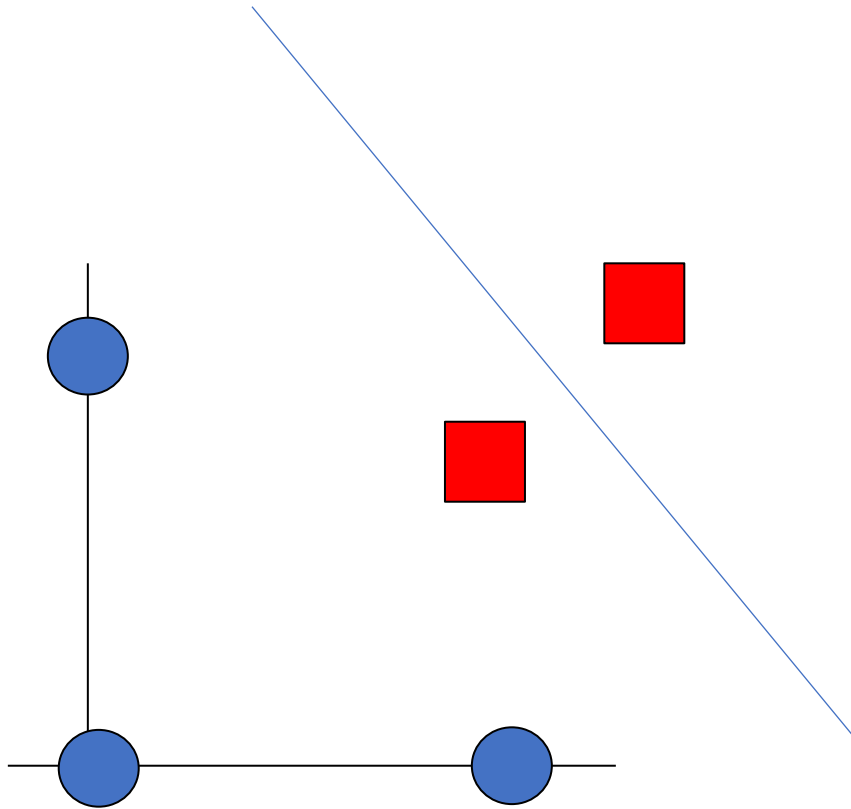
$$\begin{aligned} y_i &= a(u_i) = a\left(\sum_j w_{ij}x_j\right) \\ &= a(w_{i1}x_1 + w_{i2}x_2 + w_{i0}) \\ &= w_{i1}x_1 + w_{i2}x_2 + w_{i0} \end{aligned}$$

What do we want?

$$w_{i1}x_1 + w_{i2}x_2 + w_{i0} > 0$$

$$w_{i1}x_1 + w_{i2}x_2 + w_{i0} < 0$$

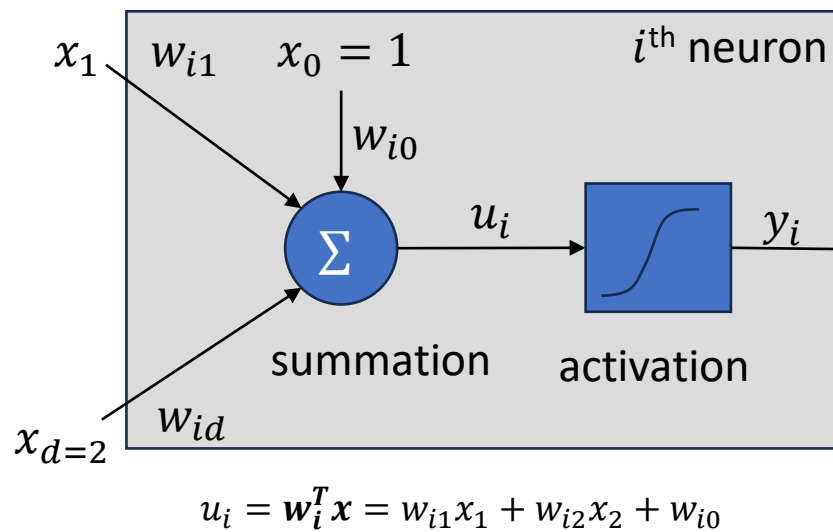
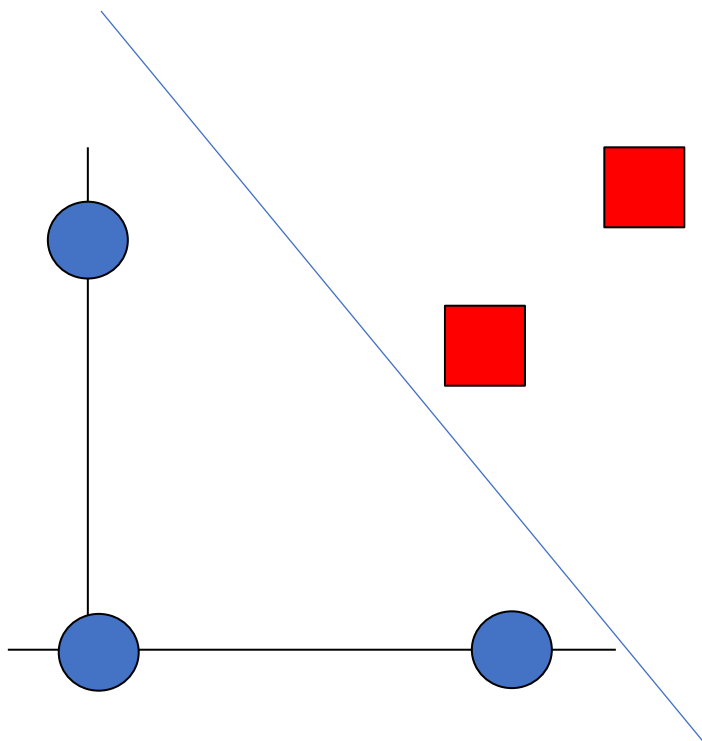
# Evaluation



$$\begin{aligned} y_i &= a(u_i) = a\left(\sum_j w_{ij}x_j\right) \\ &= a(w_{i1}x_1 + w_{i2}x_2 + w_{i0}) \\ &= w_{i1}x_1 + w_{i2}x_2 + w_{i0} \end{aligned}$$

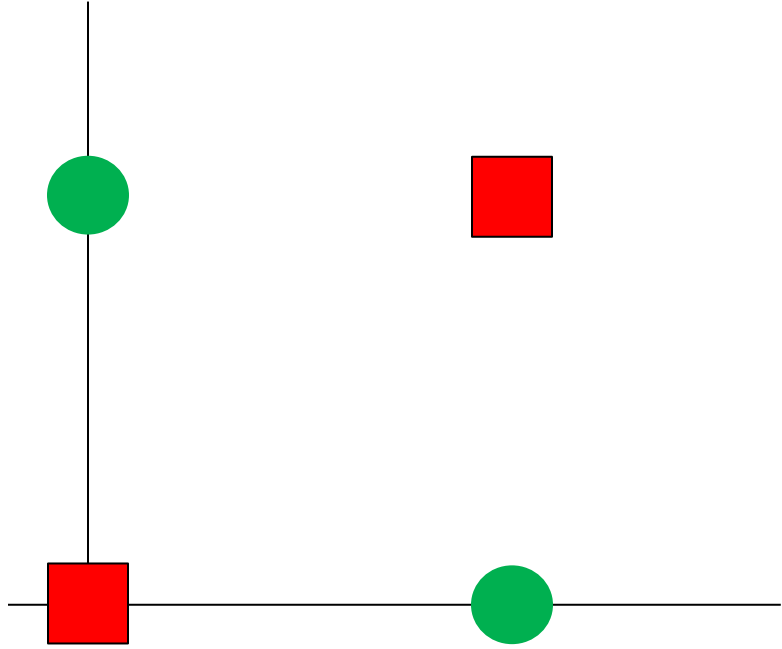
Calculate how much error the line produces

# Optimization



$$\begin{aligned} y_i &= a(u_i) = a\left(\sum_j w_{ij}x_j\right) \\ &= a(w_{i1}x_1 + w_{i2}x_2 + w_{i0}) \\ &= w_{i1}x_1 + w_{i2}x_2 + w_{i0} \end{aligned}$$

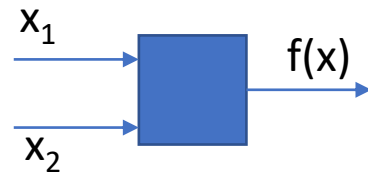
Update Weight parameters to reduce error (using gradient descent)





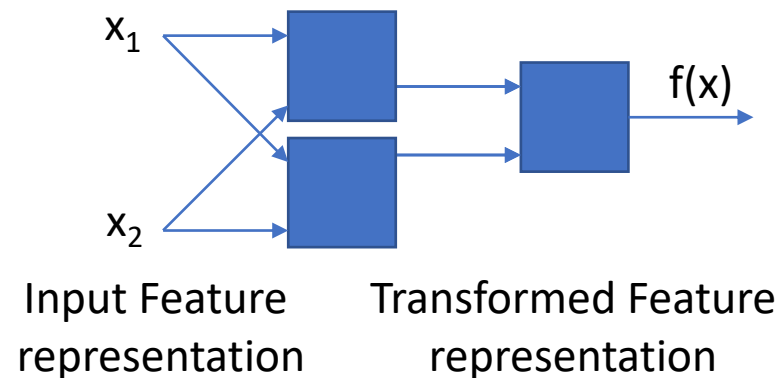
# How to fold the space?

- Change the definition of distance between points
- How to achieve this?
  - By transforming the features of the examples to another space



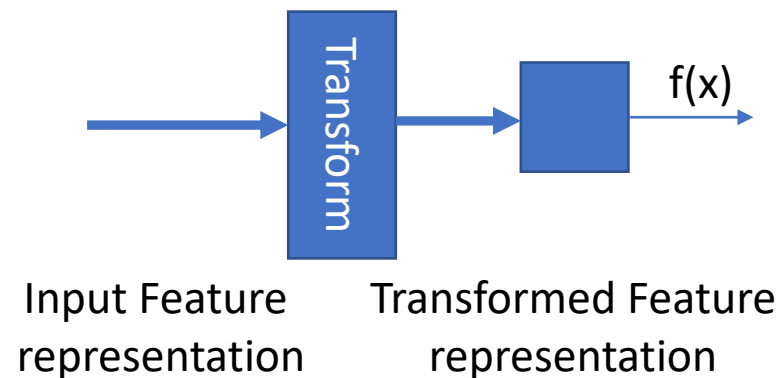
# How to fold the space?

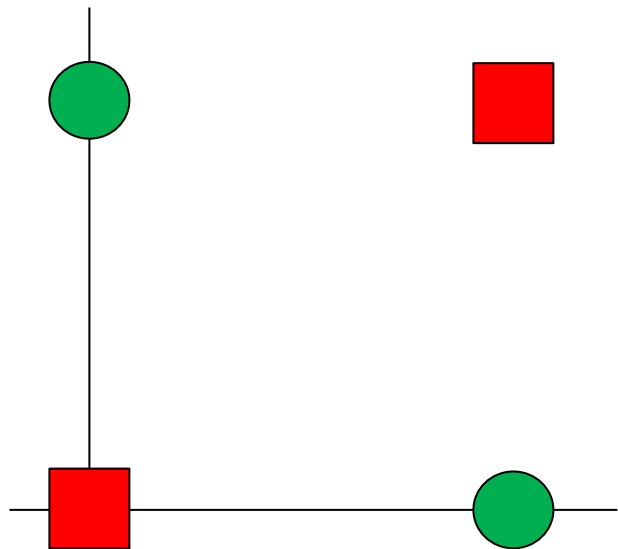
- Change the definition of distance between points
- How to achieve this?
  - By transforming the features of the examples to another space



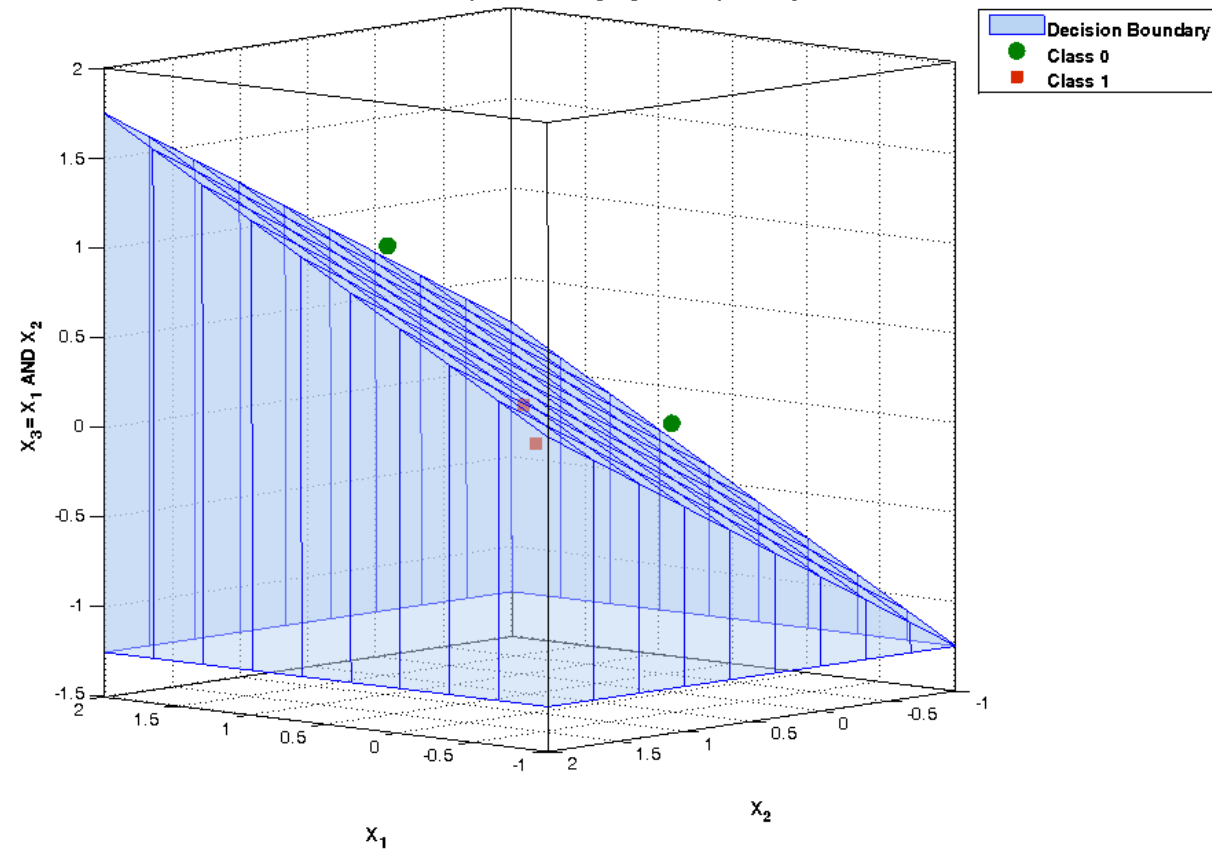
# How to fold the space?

- Change the definition of distance between points
- How to achieve this?
  - By transforming the features of the examples to another space
  - At an abstract level





Linear Discrimination Boundary for XOR using Higher Subspace Projection



$x^{(1)}$	$x^{(2)}$	$y$
0	0	-1
0	1	+1
1	0	+1
1	1	-1

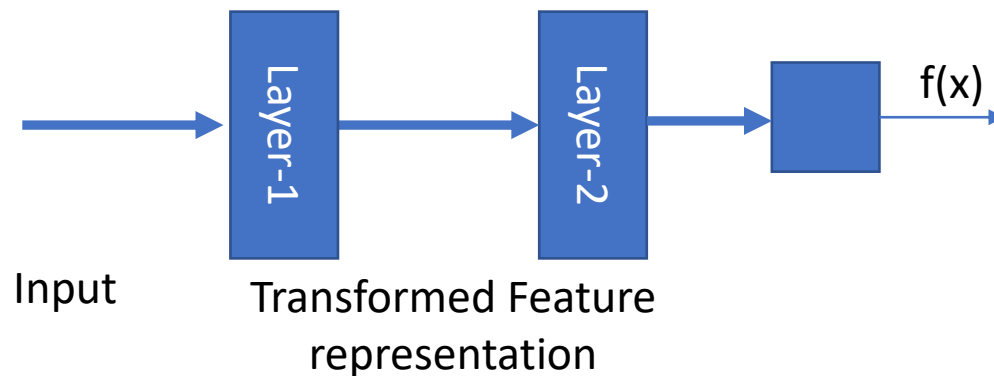
$$\phi \left( \begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix} \right) = \begin{bmatrix} x^{(1)2} \\ x^{(2)2} \\ \sqrt{2}x^{(1)}x^{(2)} \end{bmatrix}$$



$x'^{(1)}$	$x'^{(2)}$	$x'^{(3)}$	$y$
0	0	0	-1
0	1	0	+1
1	0	0	+1
1	1	$\sqrt{2}$	-1

# How to fold the space?

- Change the definition of distance between points
- How to achieve this?
  - By transforming the features of the examples to another space
  - We can do it multiple times





# Why deep learning?

- By multiple layers of neurons, we can achieve
  - Drawing lines
  - Linear separability
  - Implicit representation learning
- Deeper architectures are more “efficient at learning representations”
  - We need fewer cuts to cut a shape if we fold many times

# Tinker With a **Neural Network** Right Here in Your Browser. Don't Worry, You Can't Break It. We Promise.



Epoch  
000,246

Learning rate

0.03



Activation

Tanh



Regularization

None



Regularization rate

0



Problem type

Classification



## DATA

Which dataset do you want to use?



Ratio of training to test data: 50%



Noise: 0



Batch size: 10



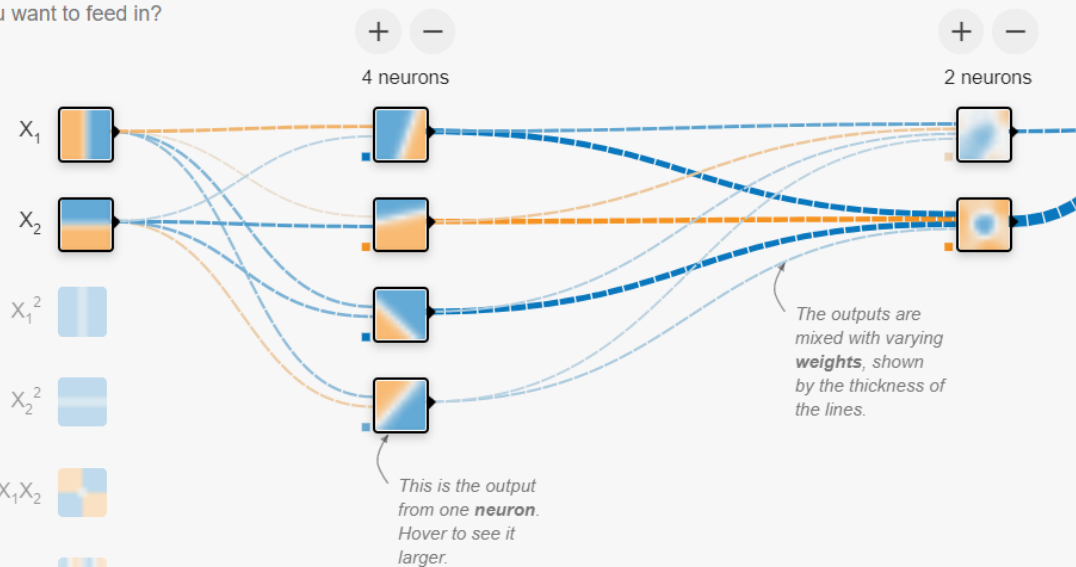
REGENERATE

## FEATURES

Which properties do you want to feed in?

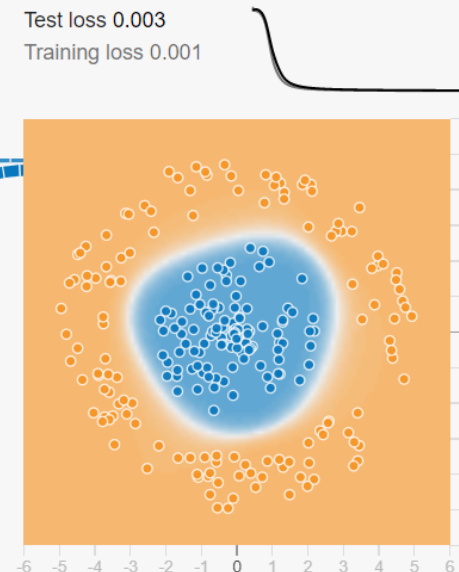
- $X_1$
- $X_2$
- $X_1^2$
- $X_2^2$
- $X_1 X_2$
- $\sin(X_1)$
- $\sin(X_2)$

## 2 HIDDEN LAYERS

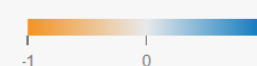


## OUTPUT

Test loss 0.003  
Training loss 0.001



Colors shows data, neuron and weight values.



Show test data

Discretize output

<https://playground.tensorflow.org>

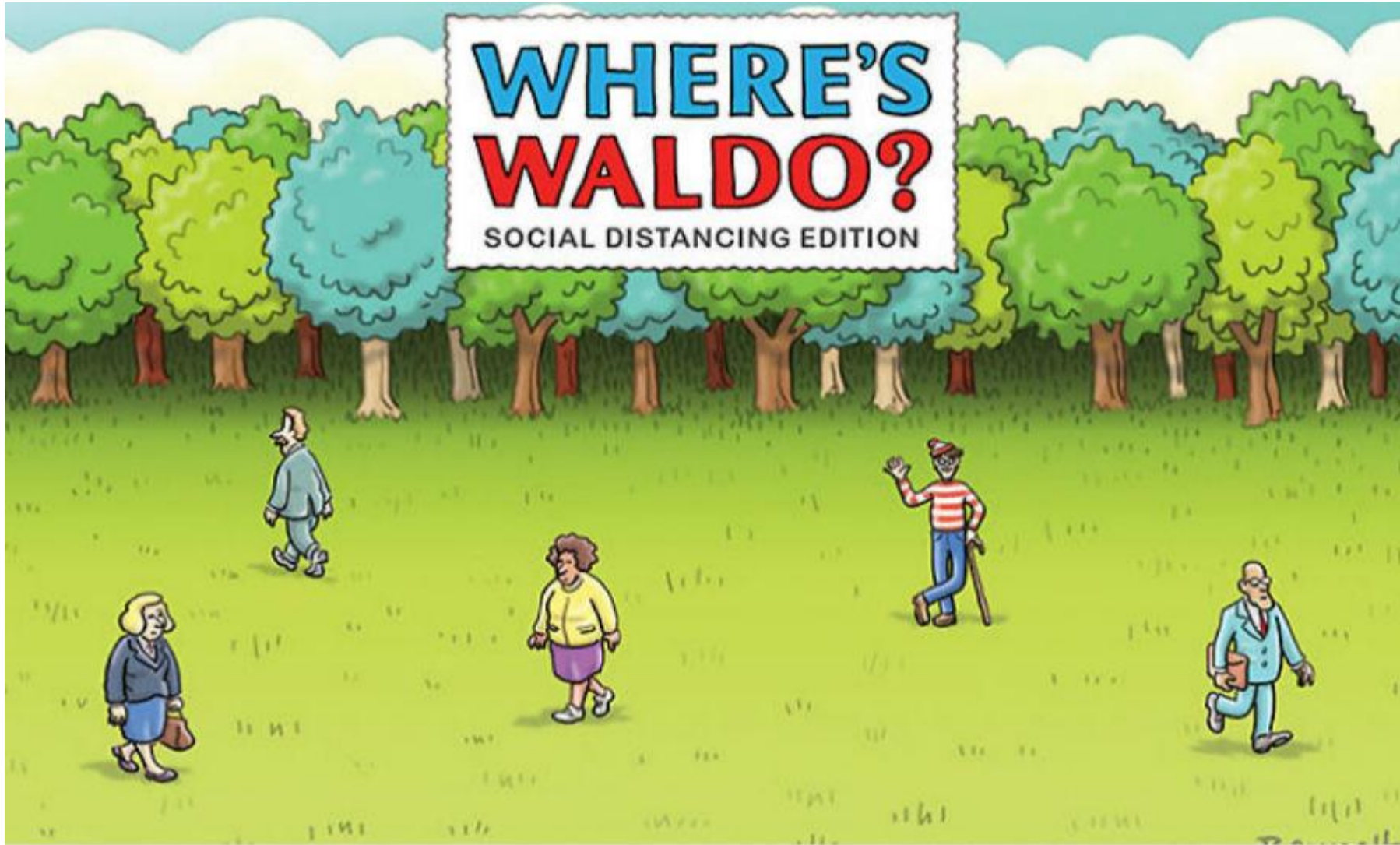
# Convolutional Neural Networks



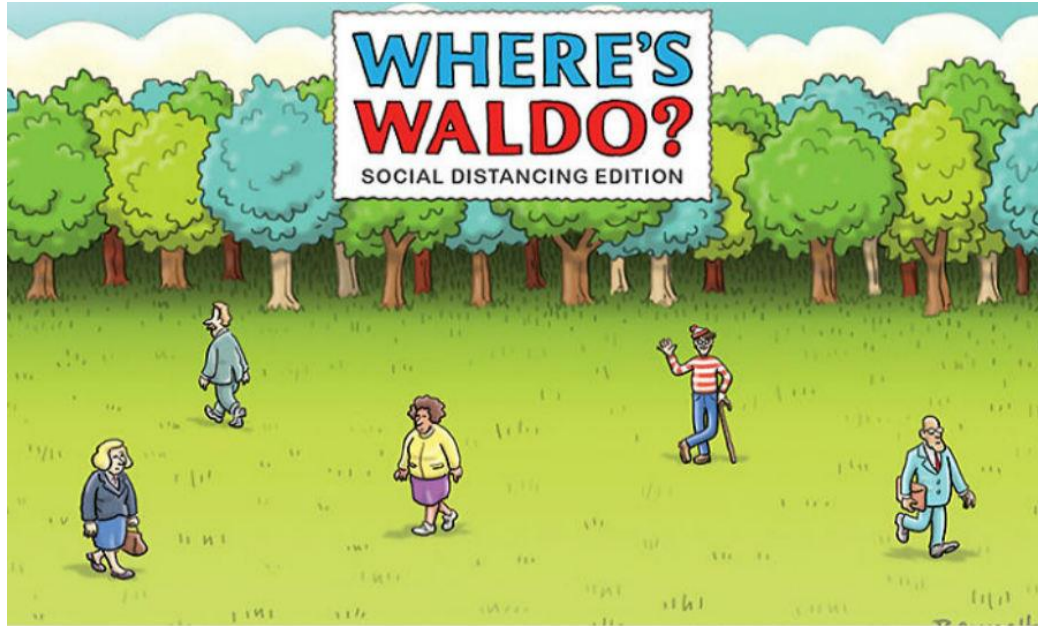
# Where's Waldo?



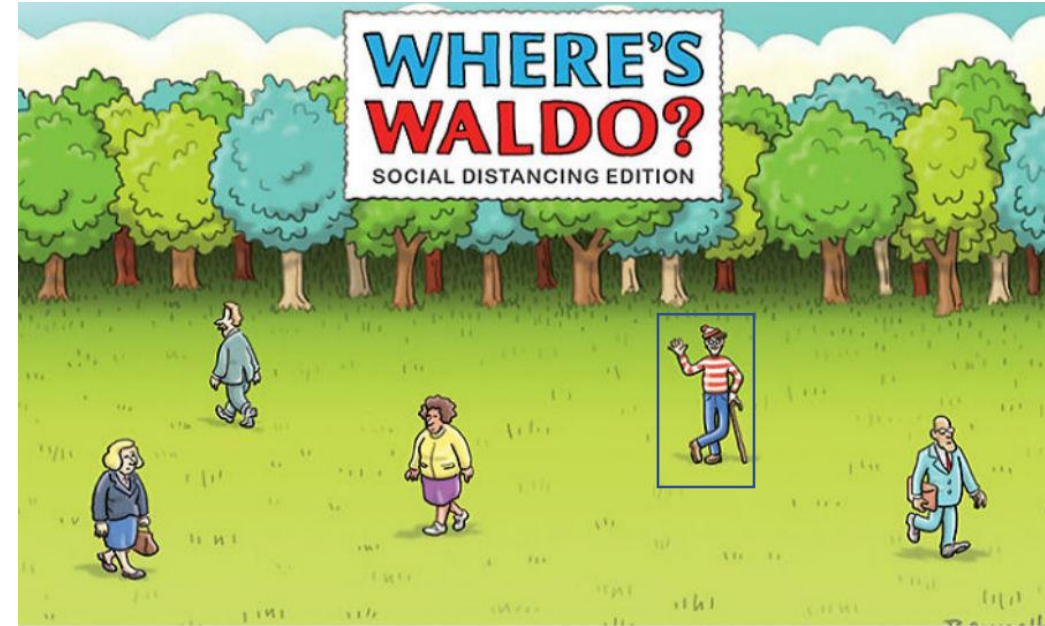


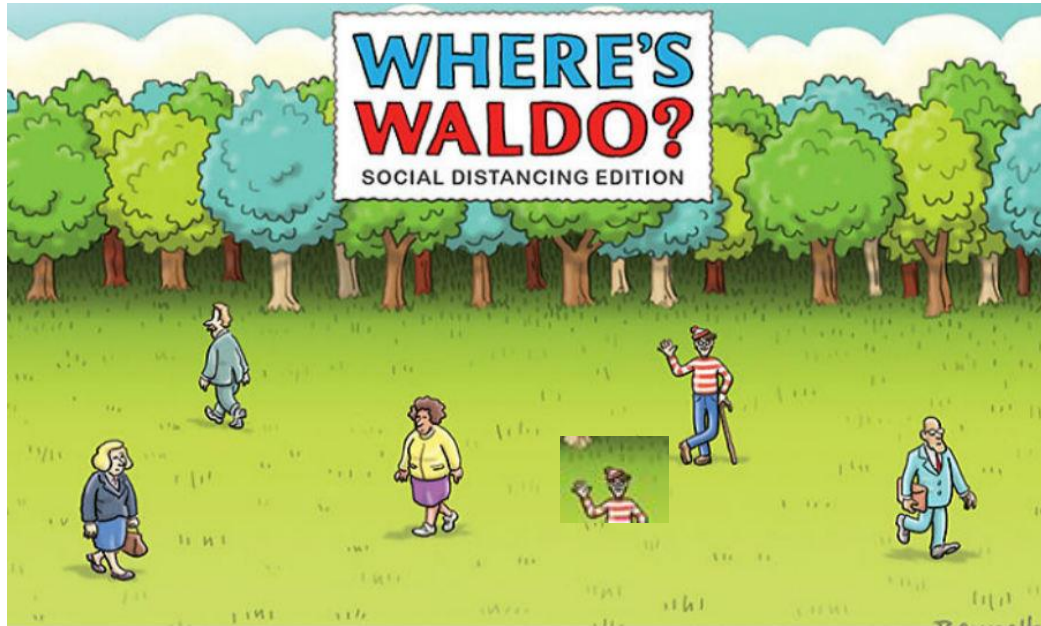






Model

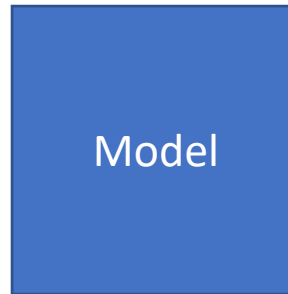
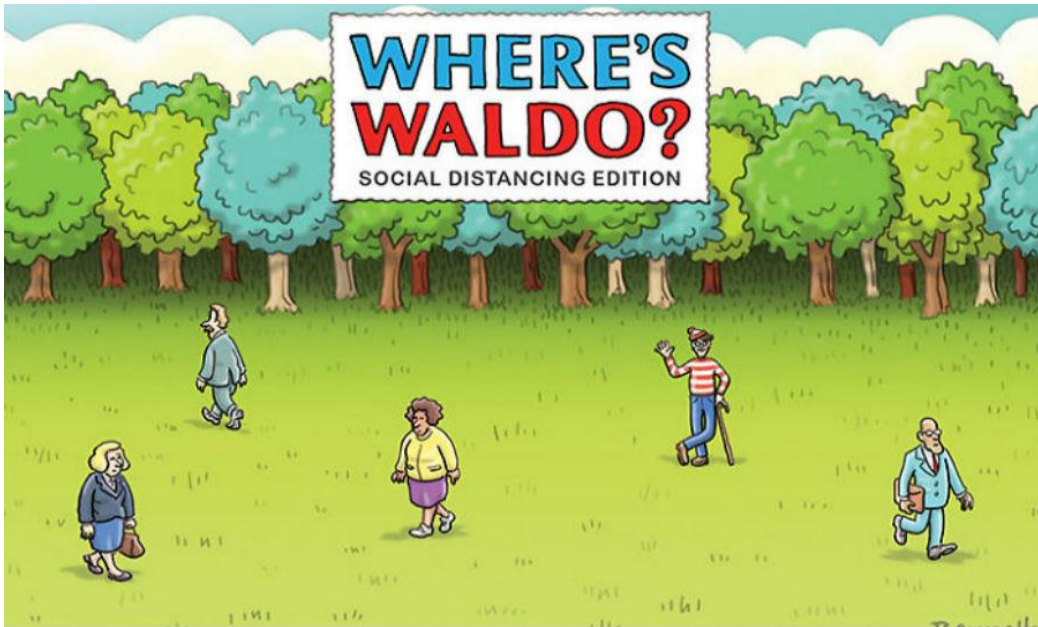




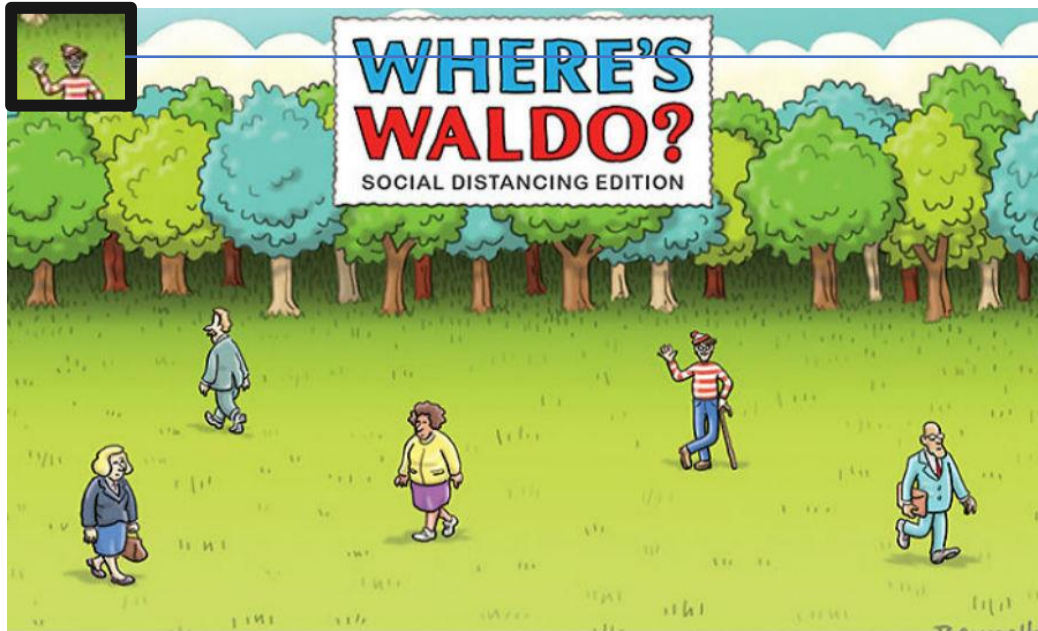
Model



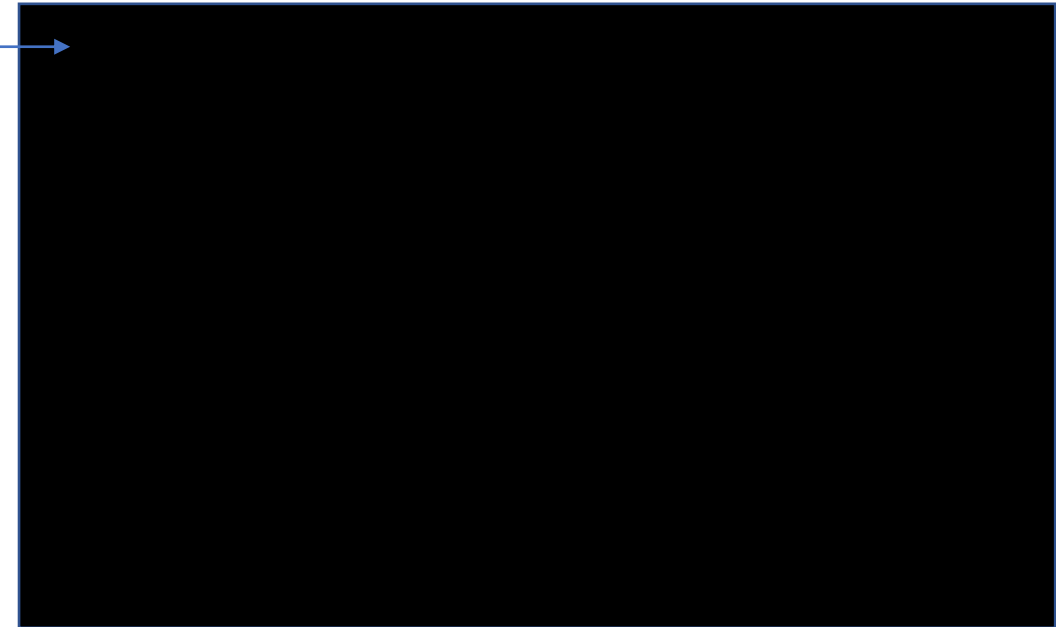
- Assume we have a cutout of what Waldo looks like
- And if we “scan” (formally called correlate or convolve) the cutout against the input image – we should see a peak at the location where Waldo occurs in the input image



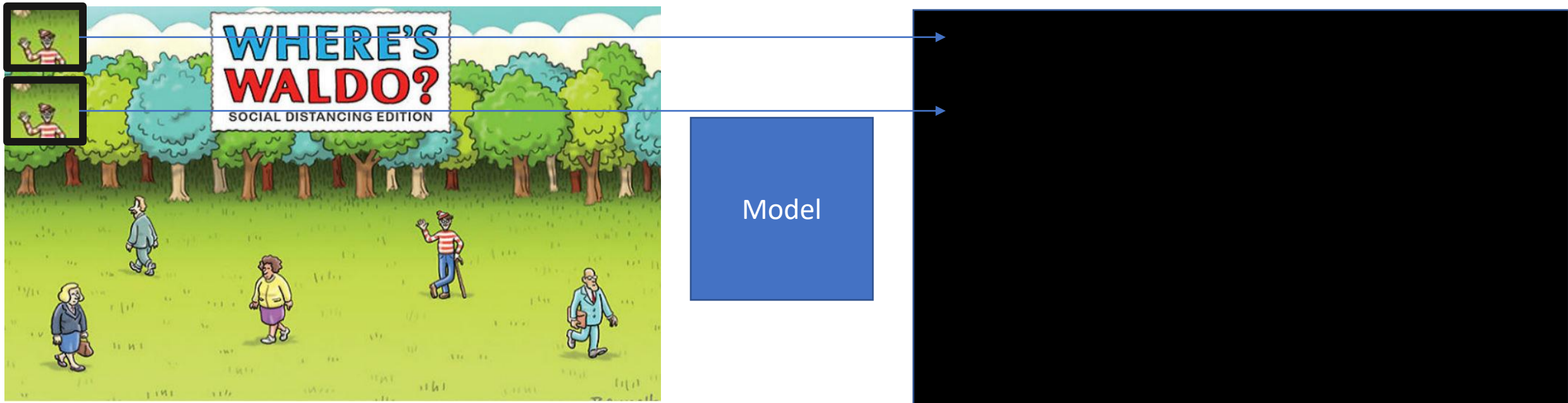
- Assume we have a cutout of what Waldo looks like
- And if we “scan” (formally called correlate or convolve) the cutout against the input image – we should see a peak at the location where Waldo occurs in the input image



Model

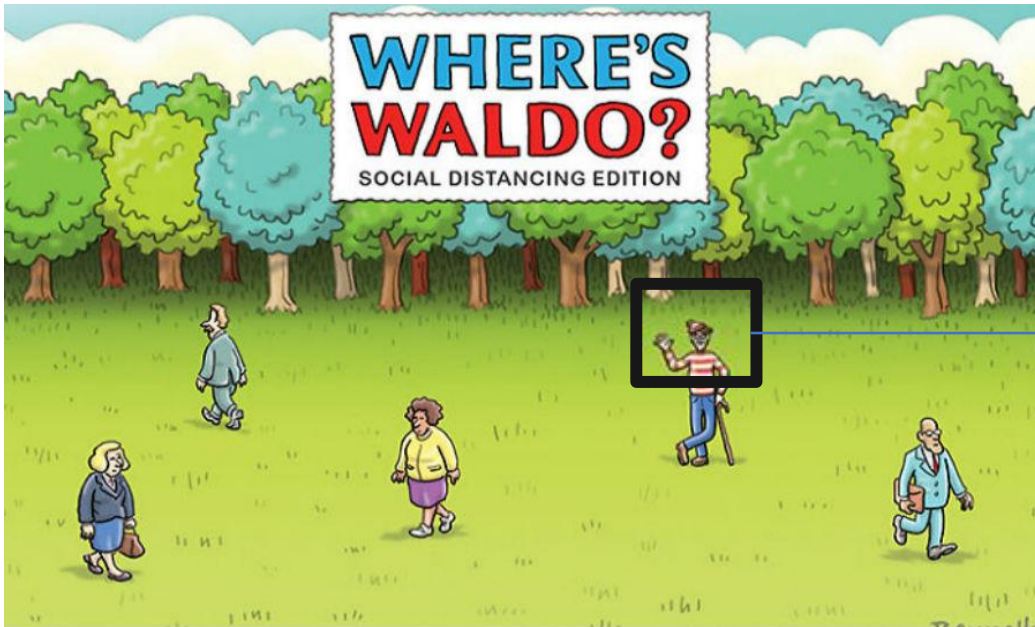


- Assume we have a cutout of what Waldo looks like
- And if we “scan” (formally called correlate or convolve) the cutout against the input image – we should see a peak at the location where Waldo occurs in the input image

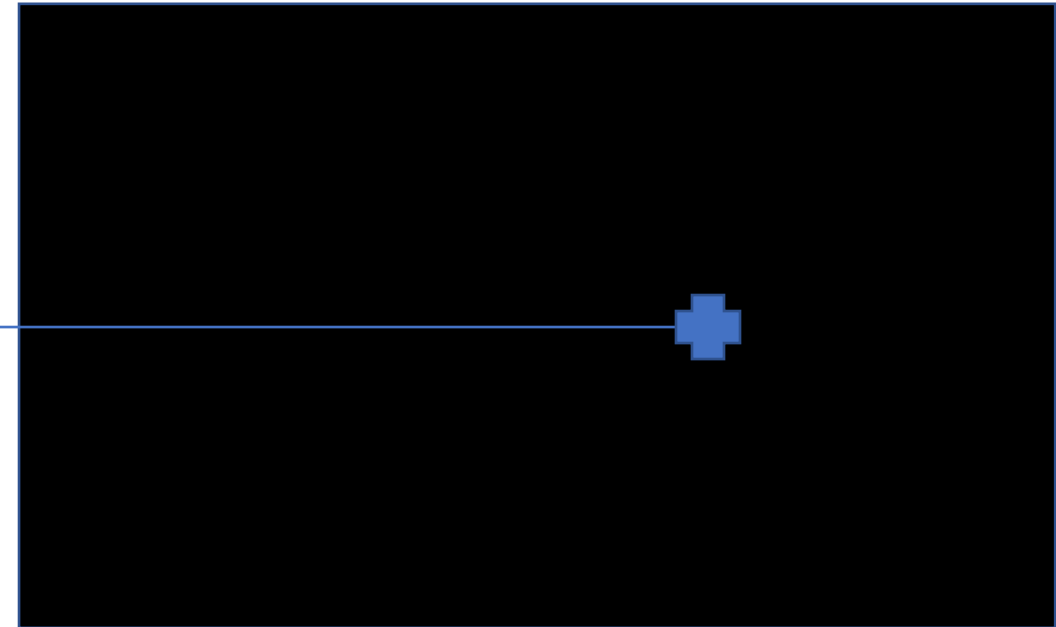




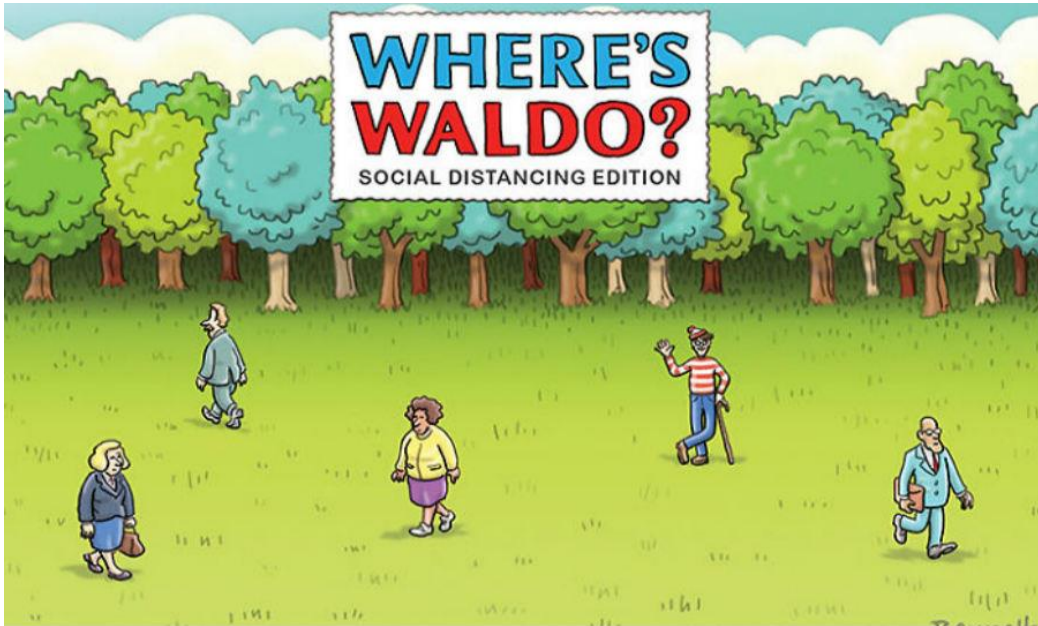
- Assume we have a cutout of what Waldo looks like
- And if we “scan” (formally called correlate or convolve) the cutout against the input image – we should see a peak at the location where Waldo occurs in the input image



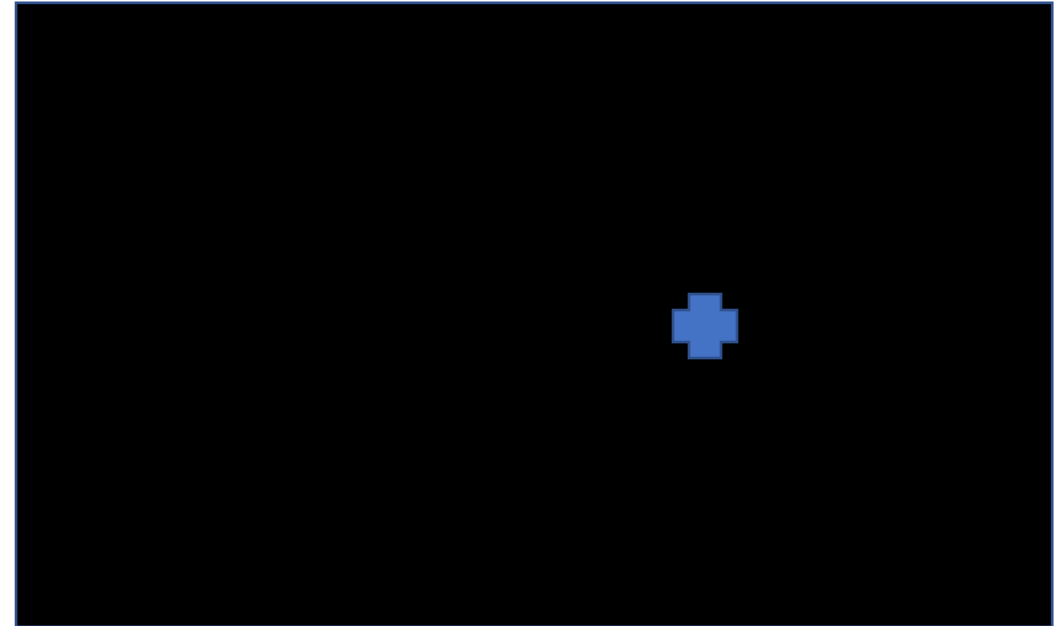
Model



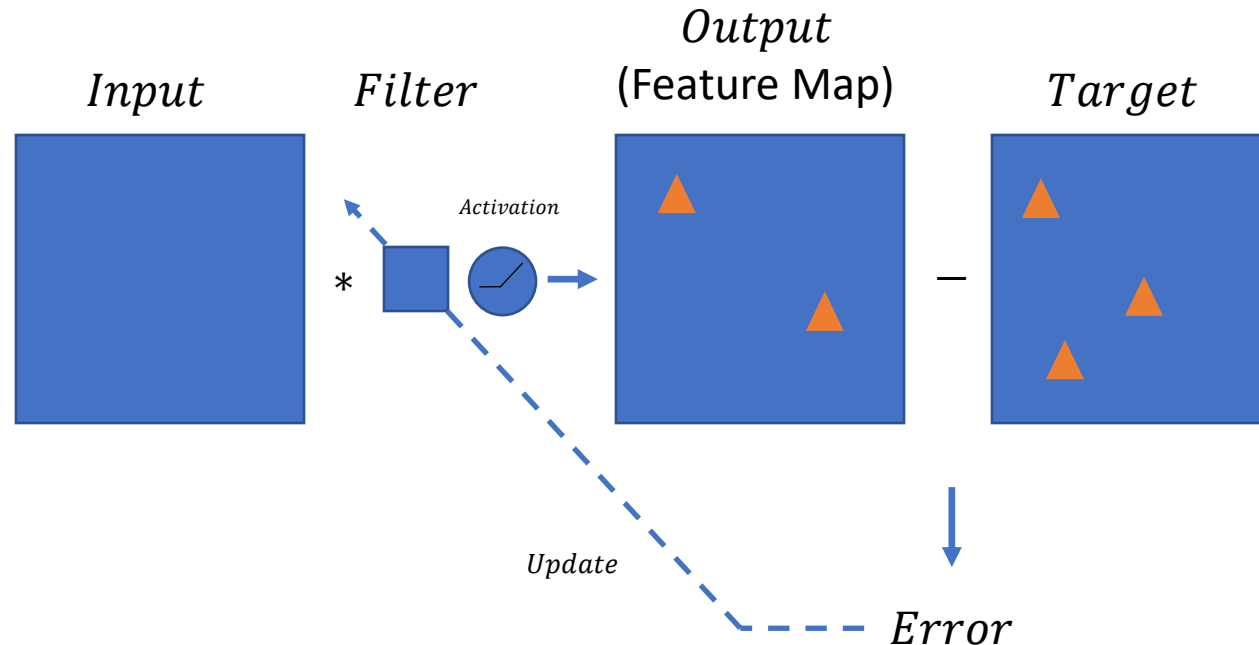
- What if I don't have a cut out of Waldo?
- Can we find him still?



Model



# (Very) basic “convolutional neural network”

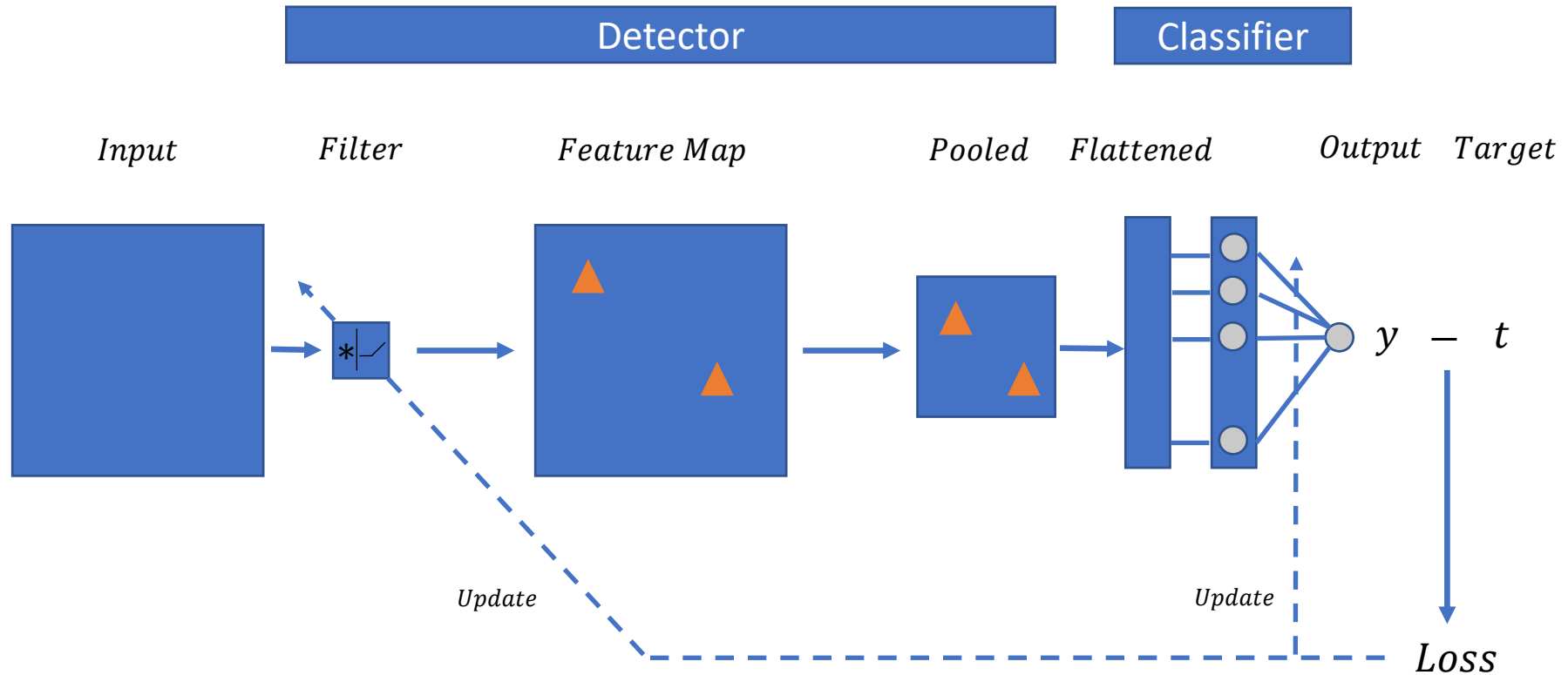


- Acts as a “detection” or “feature extraction” unit

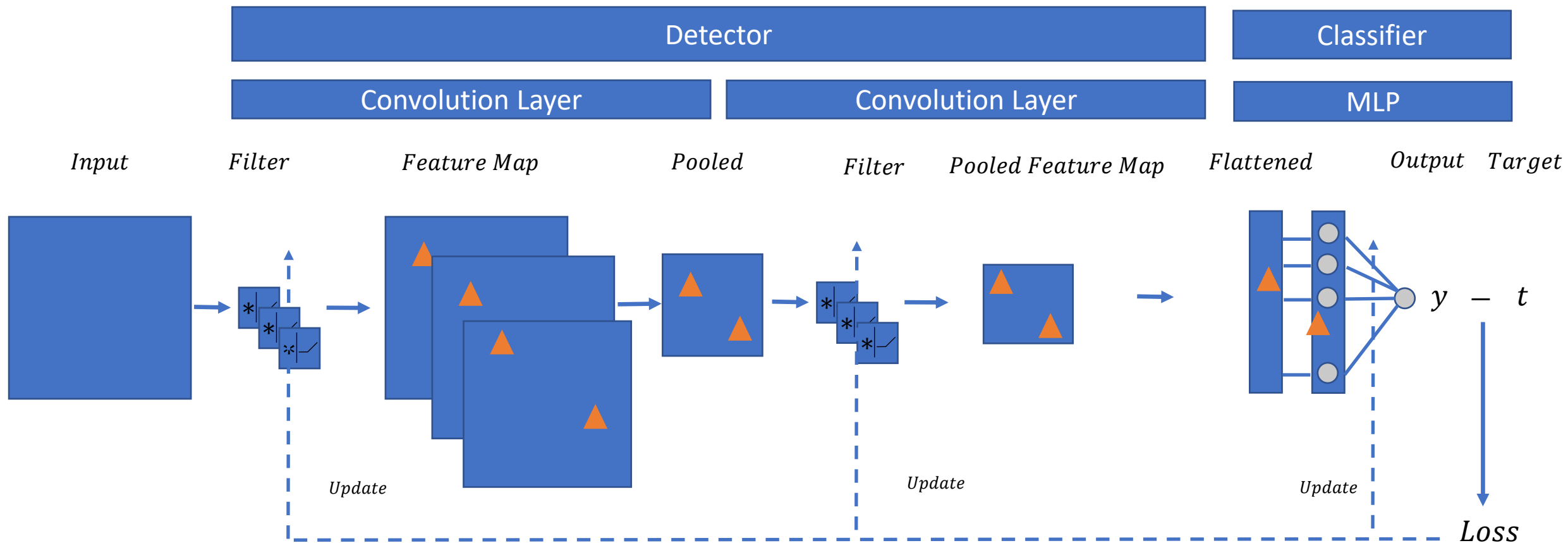
[https://github.com/foxtrotmike/CS909/blob/master/learn\\_filters.ipynb](https://github.com/foxtrotmike/CS909/blob/master/learn_filters.ipynb)

[https://github.com/foxtrotmike/CS909/blob/master/cnn\\_mnist\\_pytorch.ipynb](https://github.com/foxtrotmike/CS909/blob/master/cnn_mnist_pytorch.ipynb)

# Basic convolutional neural network for ML



# CNNs





A mostly complete chart of

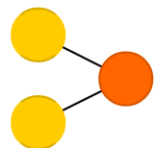
# Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

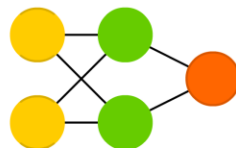
<http://www.asimovinstitute.org/neural-network-zoo/>

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

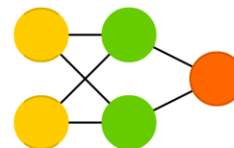
Perceptron (P)



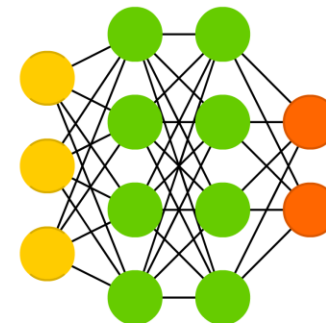
Feed Forward (FF)



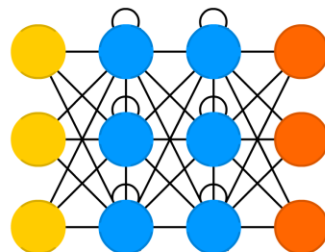
Radial Basis Network (RBF)



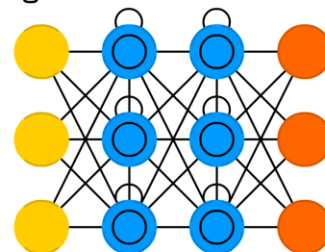
Deep Feed Forward (DFF)



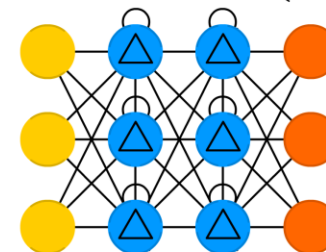
Recurrent Neural Network (RNN)



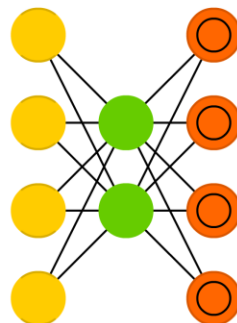
Long / Short Term Memory (LSTM)



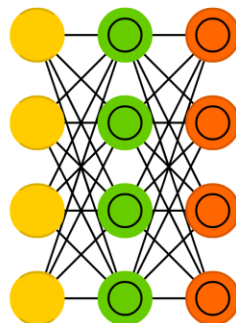
Gated Recurrent Unit (GRU)



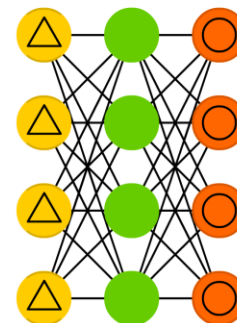
Auto Encoder (AE)



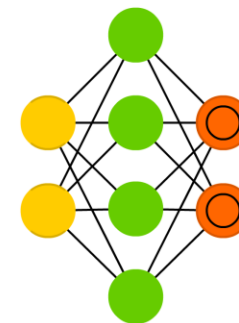
Variational AE (VAE)



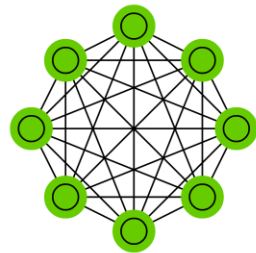
Denosing AE (DAE)



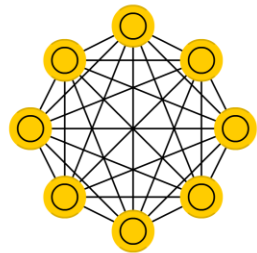
Sparse AE (SAE)



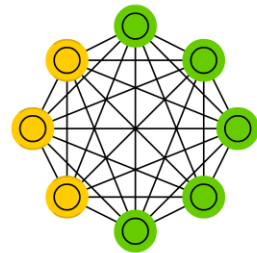
Markov Chain (MC)



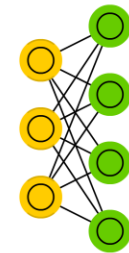
Hopfield Network (HN)



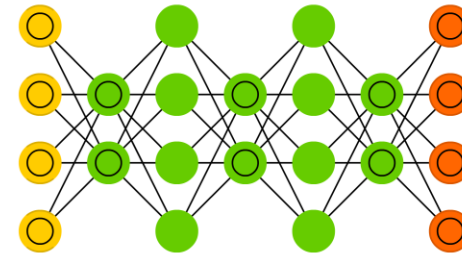
Boltzmann Machine (BM)



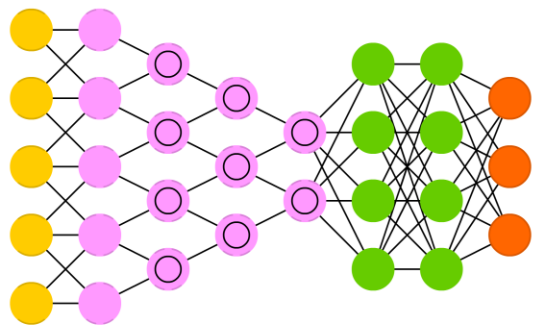
Restricted BM (RBM)



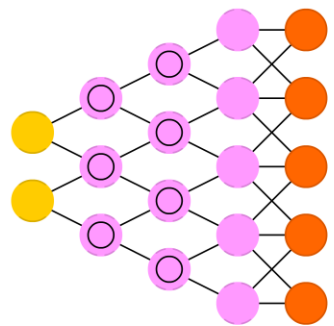
Deep Belief Network (DBN)



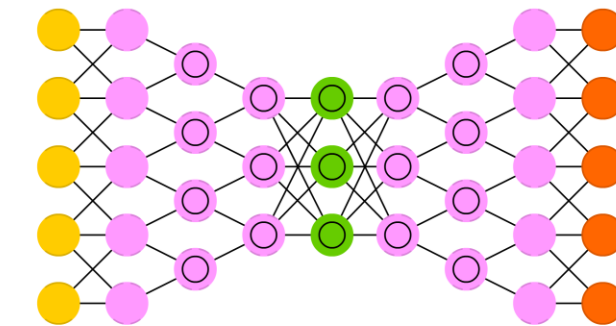
Deep Convolutional Network (DCN)



Deconvolutional Network (DN)

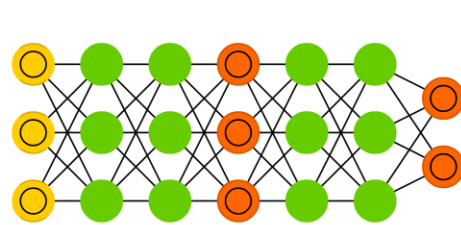


Deep Convolutional Inverse Graphics Network (DCIGN)

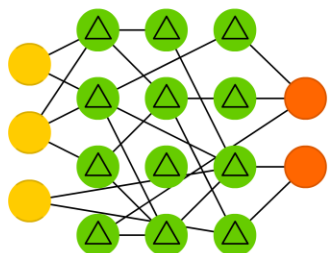


- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool

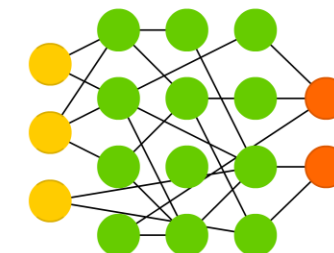
Generative Adversarial Network (GAN)



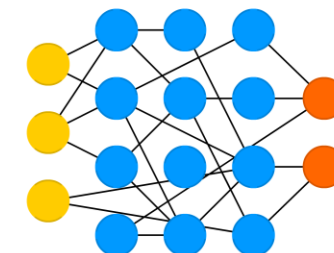
Liquid State Machine (LSM)



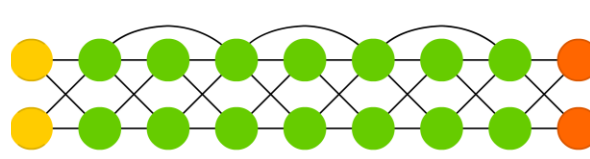
Extreme Learning Machine (ELM)



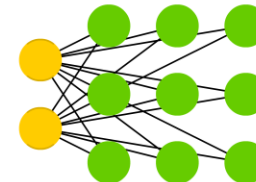
Echo State Network (ESN)



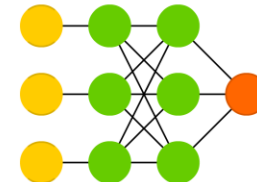
Deep Residual Network (DRN)



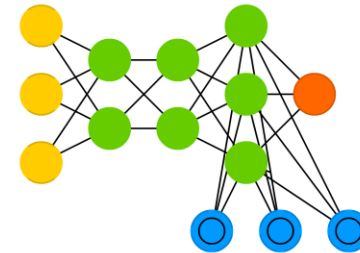
Kohonen Network (KN)

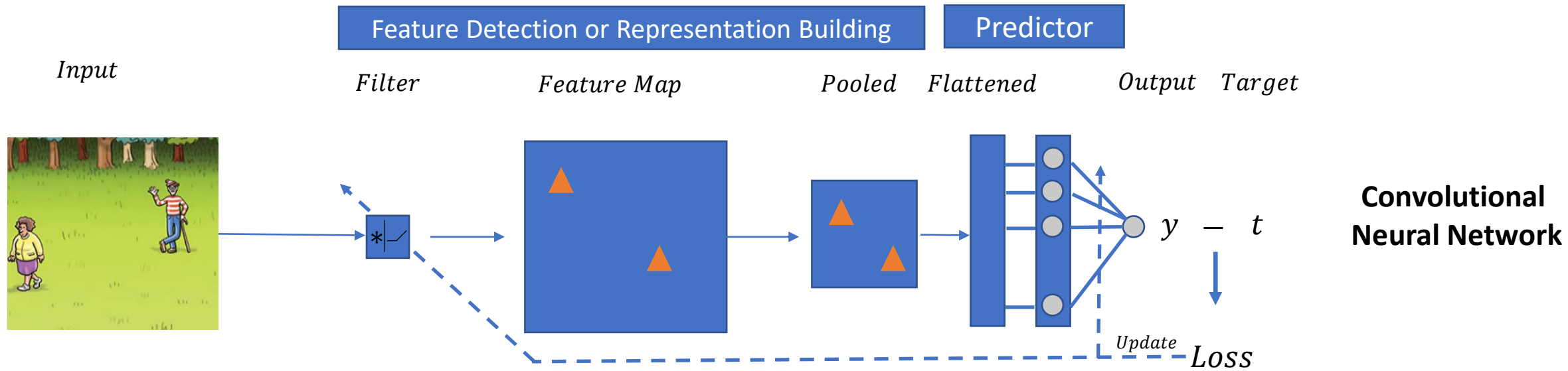


Support Vector Machine (SVM)

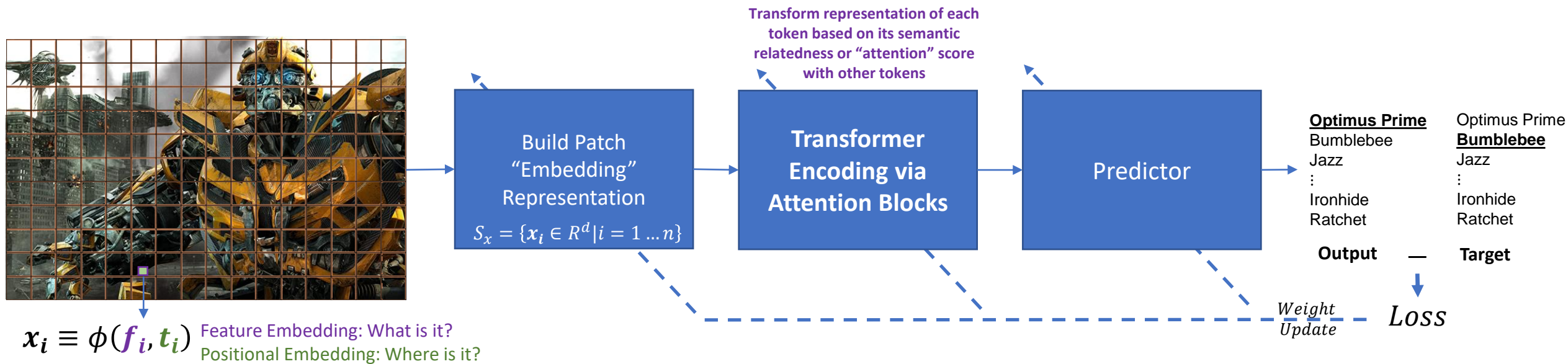


Neural Turing Machine (NTM)



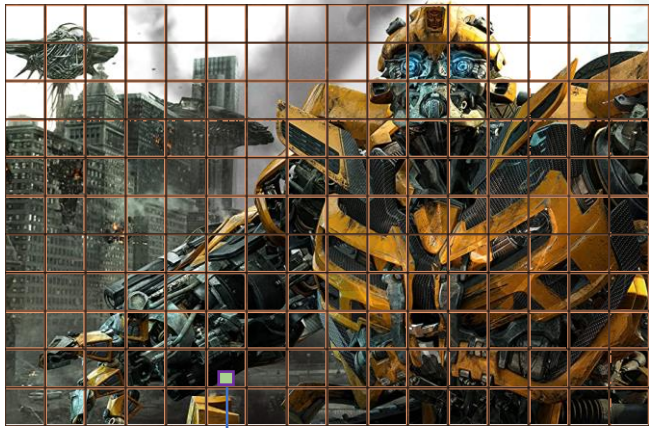


### (Vision) Transformers



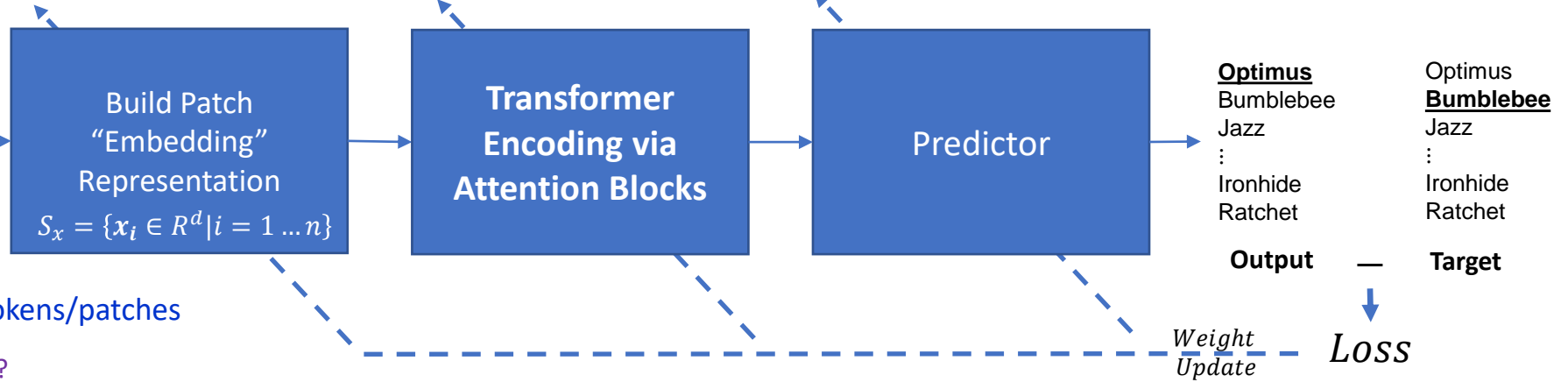
# (Vision) Transformers (for classification)

1 2 ...



$n$  tokens/patches

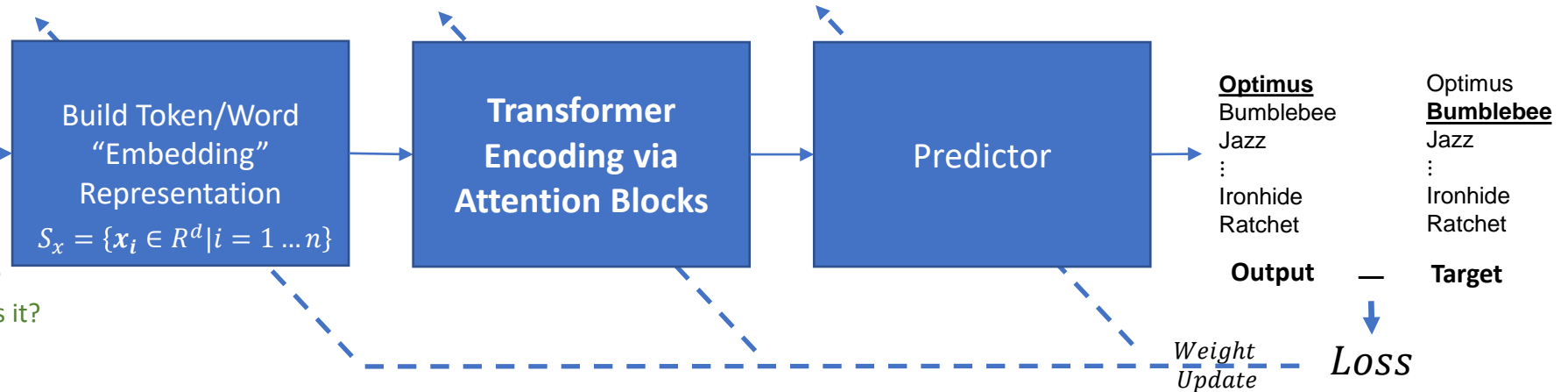
$x_i \equiv \phi(f_i, t_i)$  Feature Embedding: What is it?  
Positional Embedding: Where is it?



Transform representation of each token based on its semantic relatedness or "attention" score with other tokens

A transformer that can transform into a yellow car is called \_\_\_\_\_.

$x_i \equiv \phi(f_i, t_i)$  Feature Embedding: What is it?  
Positional Embedding: Where is it?



# (NLP) Transformers (for next word prediction)

Simplest:  $\phi(f_i, t_i) = f_i + t_i$



# What is attention and why do you need it?

[Submitted on 12 Jun 2017 (v1), last revised 6 Dec 2017 (this version, v5)]

## Attention Is All You Need

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder-decoder configuration. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

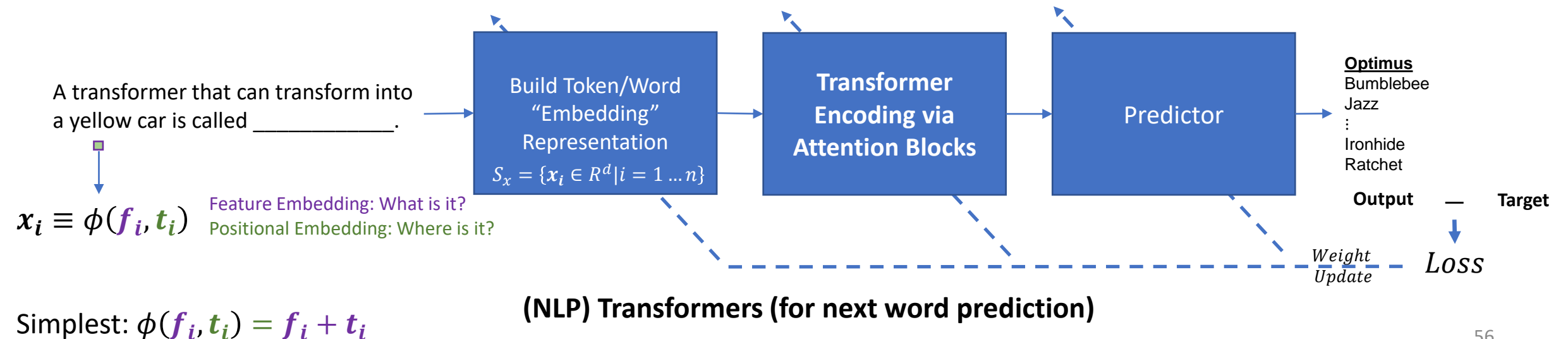
$$x'_q = A(x_q; M(x_q, S_{x_k}); \theta) = \sum_{x_k \in M(x_q, S_{x_k})} a(x_q, x_k; M(x_q, S_{x_k}), \theta_k) v(x_k; \theta_v) = \sum_{x_k \in M(x_q, S_{x_k})} \frac{k(x_q, x_k; \theta_k)}{\sum_{x_{k'} \in M(x_q, S_{x_k})} k(x_q, x_{k'}; \theta_k)} v(x_k; \theta_v)$$

Tsai, Yao-Hung Hubert, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. “**Transformer Dissection: A Unified Understanding of Transformer’s Attention via the Lens of Kernel.**” *ArXiv:1908.11775 [Cs, Stat]*, November 11, 2019. <http://arxiv.org/abs/1908.11775>.



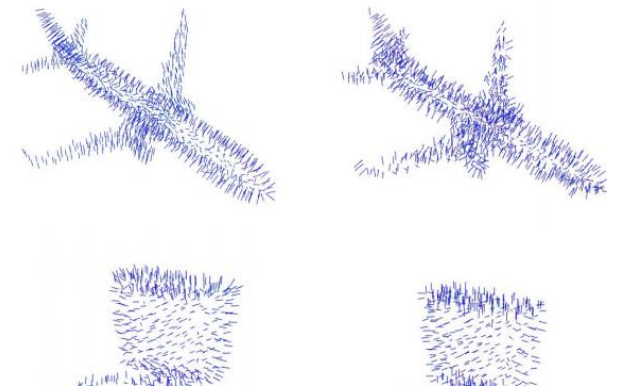
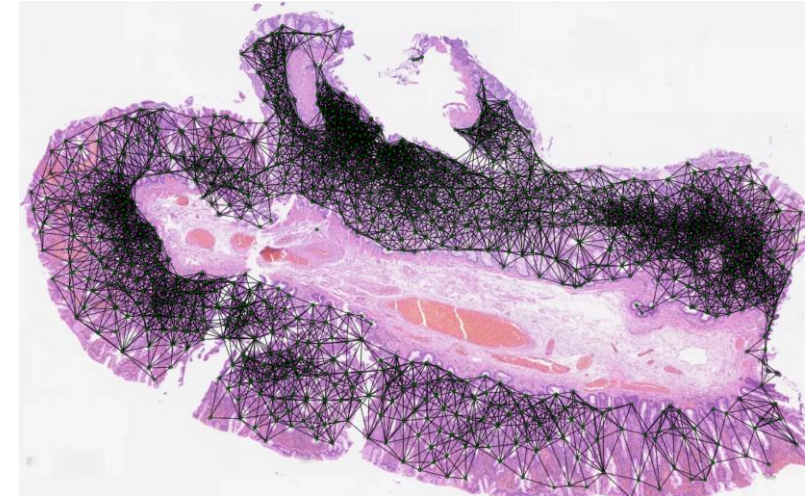
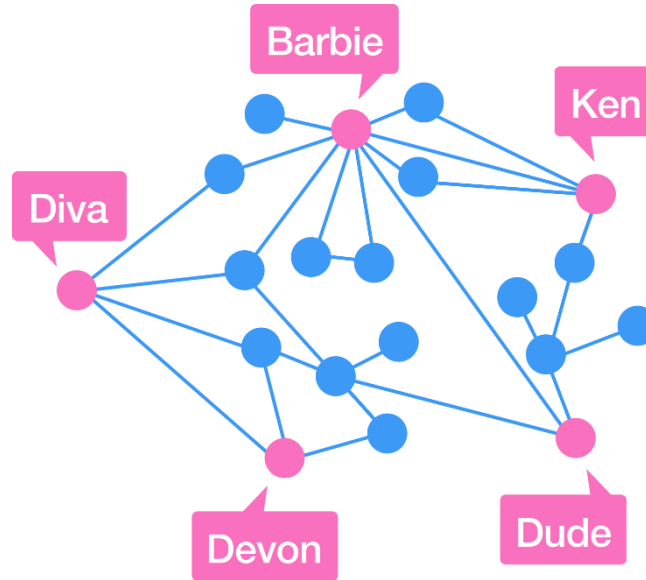
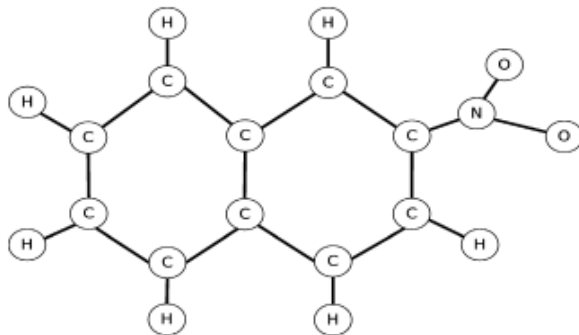
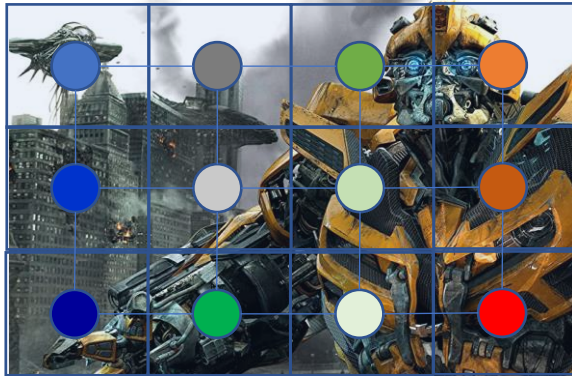
# Large Language Models: ChatGPT

- Simply: A transformer that generates the next word given some context
  - Multiple (>100) transformer layers with over a billion weights
  - Trained over the entire internet corpus
  - Fine-tuned (and controlled) with human feedback prompting



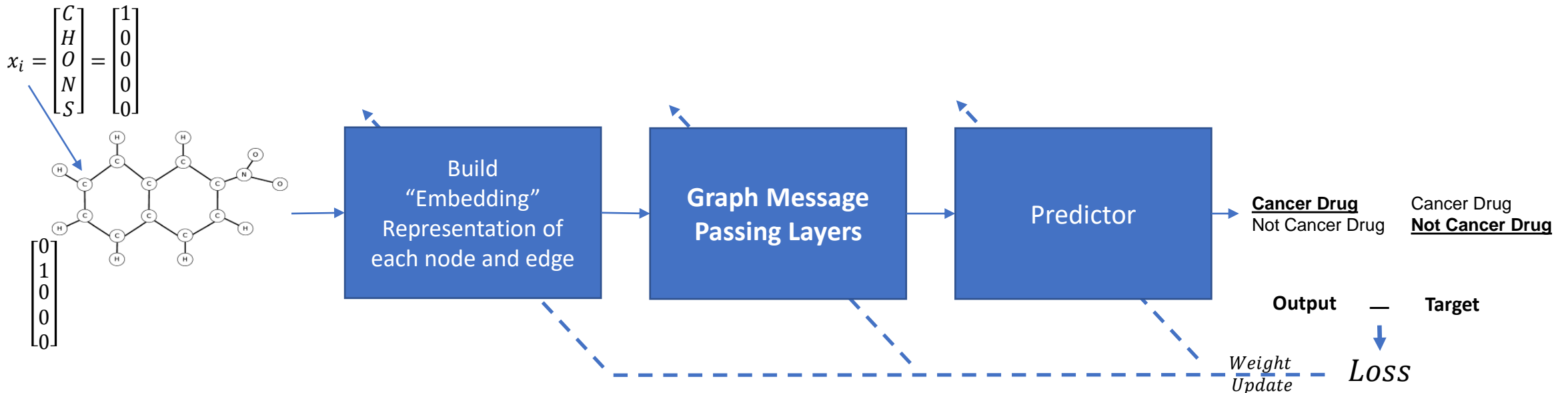
# Graphs and their neural networks

This is a graph



# Graph Neural Networks

- Simple Graph Classification Example
  - Node and edge level prediction problems also possible



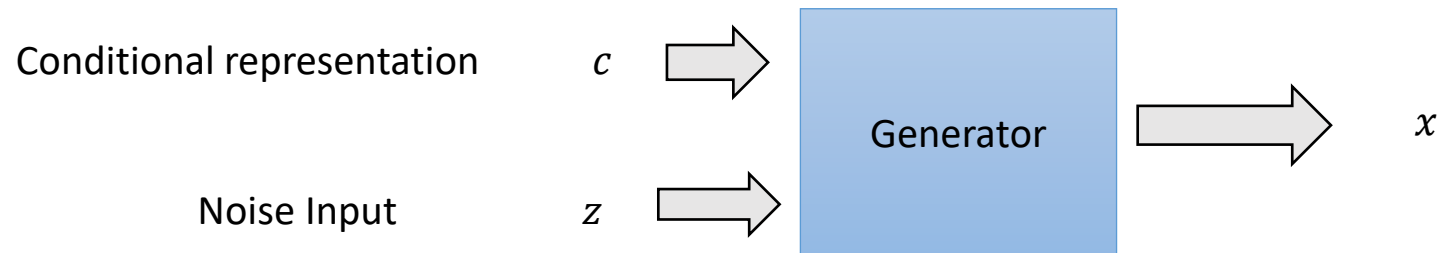
Input: Graph consisting of

Node set: what are things (each node has feature representation)

Edge set: how are they connected (each edge can have a feature representation but, in the very least, it tells us what nodes are connected by an edge)

# Generative Machine Learning

*“Image showing “It's coming home”  
in the context of the English team  
winning the football final”*



<https://github.com/foxtrotmike/CS909/blob/master/simpleGAN.ipynb>

<https://github.com/wgrgwrght/Simple-Diffusion/blob/main/SimpleDiffusion.ipynb>





# IT'S COMING HOME





## Philosophical basis

- I. Entities have (explicit or implicit) representations
- II. Semantic relatedness of entities is context dependent and thus their representations are contextual
- III. Representation of any entity can allow us to reconstruct or “generate” it
- IV. It is possible to develop representations in an inductive manner (through empirical observations)
- V. Intelligence is the capacity to develop and utilize causal representations of entities, enabling an organism or system to act effectively and adaptively.
- VI. Empirical observations may not be enough**

## Algorithms

- Learning representations
- Using Convolutions, Transformers or Graph Layers
- Generative Machine Learning: GANs, Latent Diffusion Models
- Learning Algorithm: Optimization of model parameters through gradient descent based on existing data  
Learning mechanisms: Self Supervised Learning, Next word prediction
- Deep Reinforcement Learning?
- Causal Machine Learning**

# Next Steps

- Data Mining and Machine Learning Lecture Series
  - Assumes knowledge of programming
  - Covers Python
  - Assumes no knowledge of ML
  - Starts from simple classification
  - Ends at Generative Adversarial Networks



@fayyazhere

<https://youtu.be/CjHp5HGckQk>

<https://sites.google.com/view/fayyaz/courses/data-mining>