

DEEP LEARNING

S2DS @ Westminster University

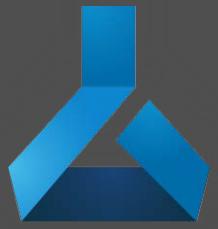
Tim Scarfe, Ph.D

Principal Software Engineer @ Microsoft
Commercial Software Engineering (CSE)
@ecsquendor
youtube.com/machinelearningatmicrosoft

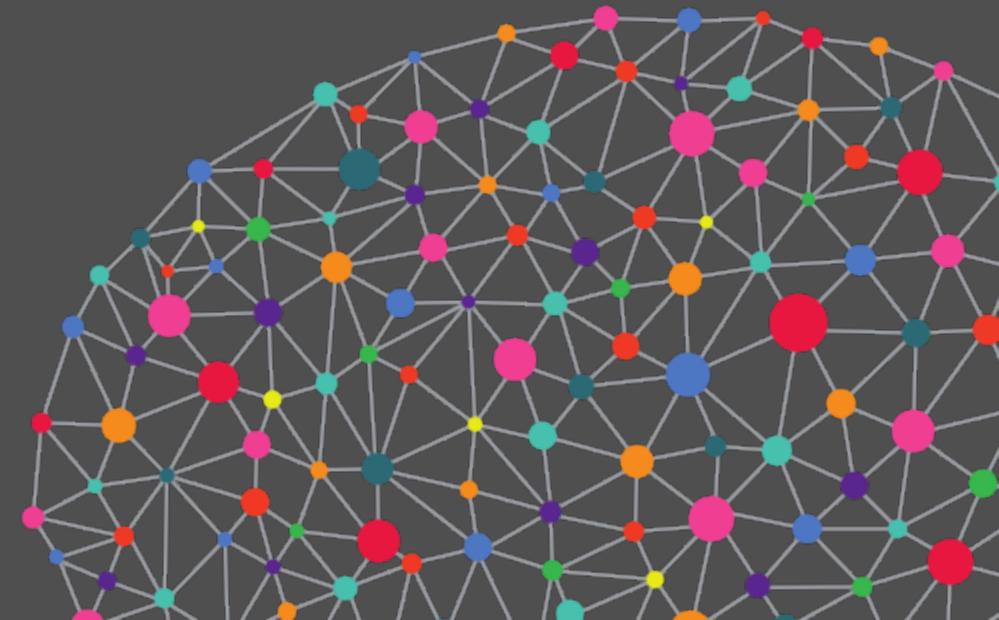
MACHINE LEARNING @ MICROSOFT



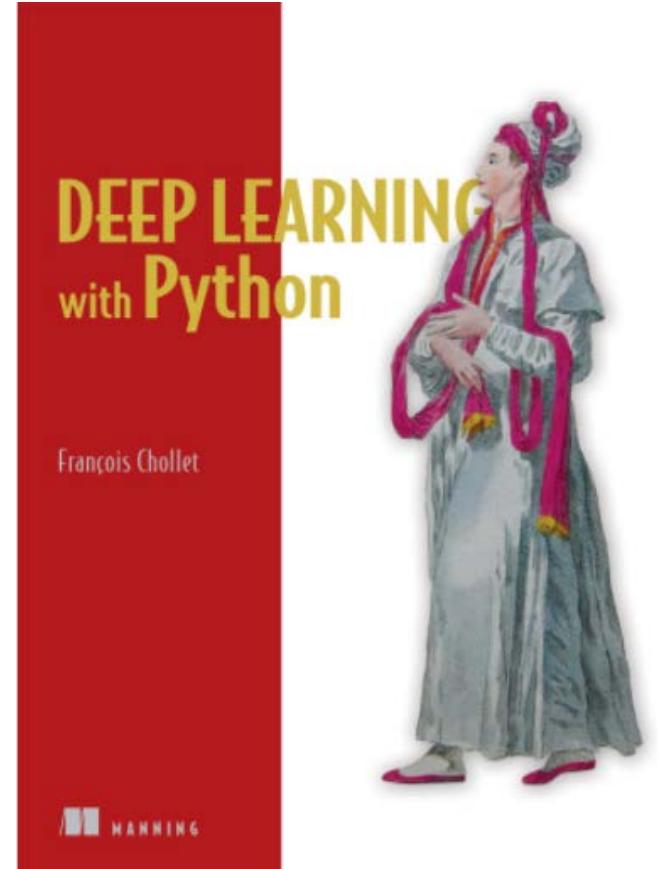
TALK OUTLINE



- Deep learning intro
- Distilled concepts of deep learning + Demo
- CNNs + Demo
- Sequence processing/NLP + Demo
- Confused.com machine learning story
- Questions



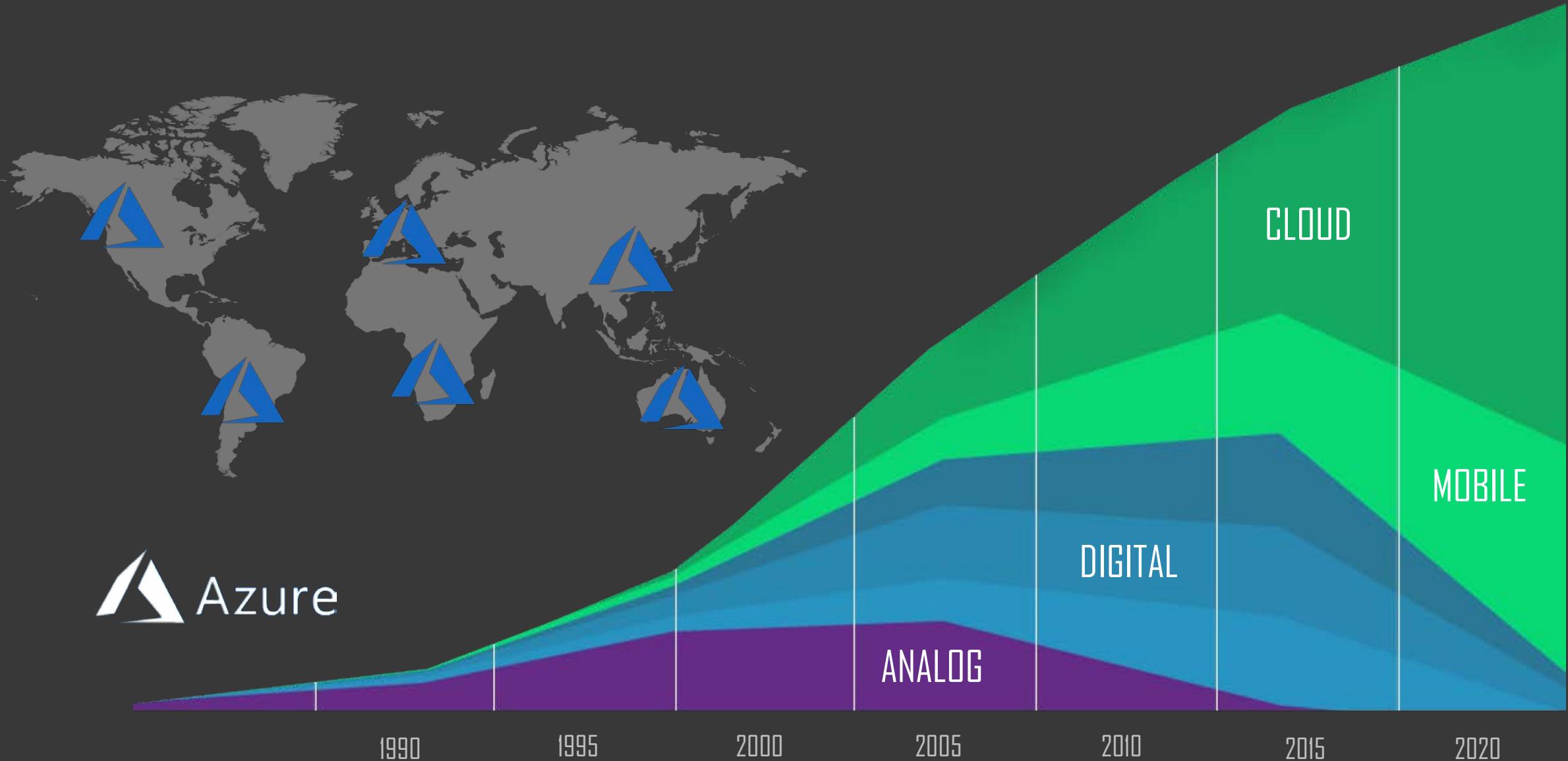
Francois Chollet is a living legend, please buy his book

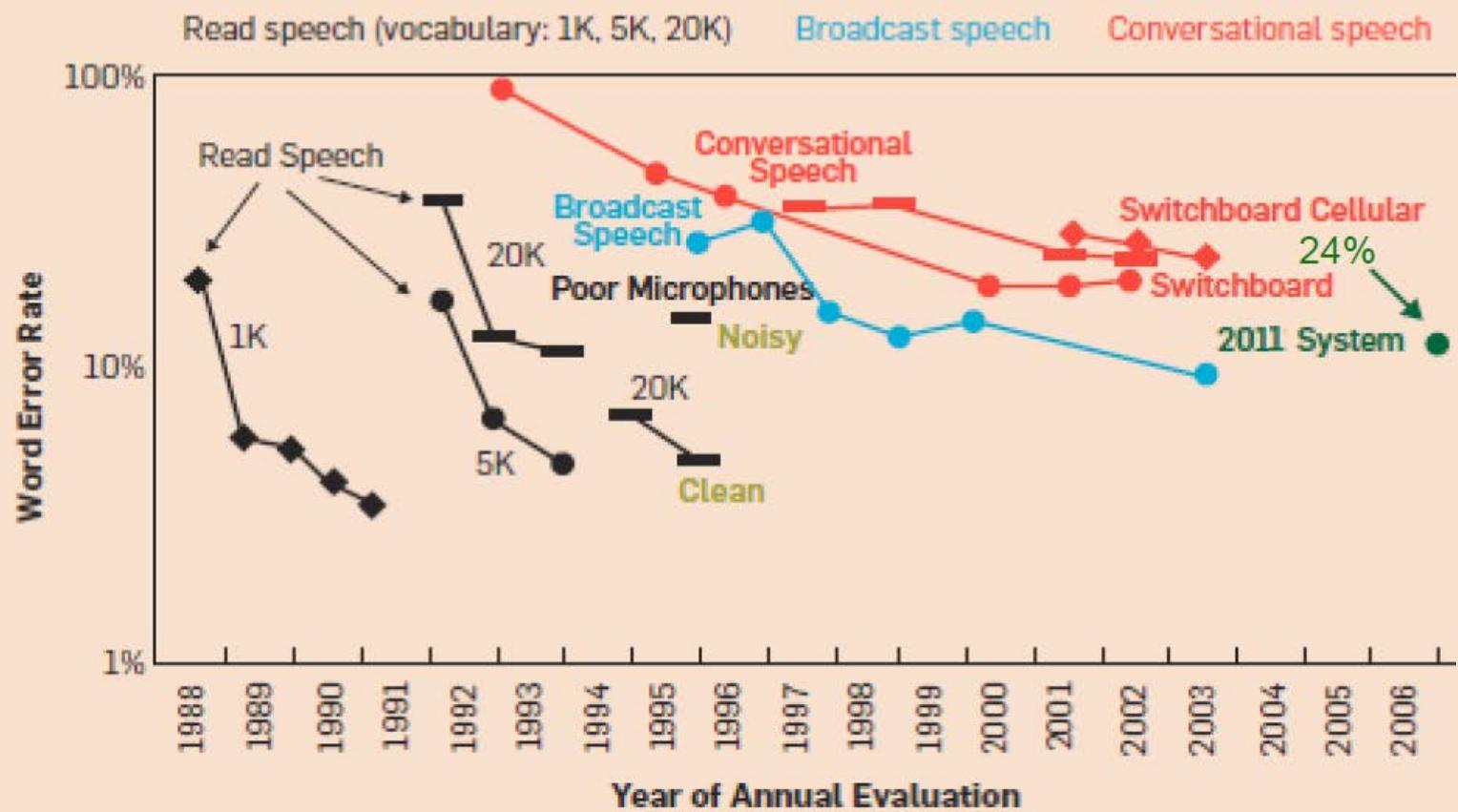


<https://www.manning.com/books/deep-learning-with-python>



COMPUTE + DATA EXPLOSION = SMART AI

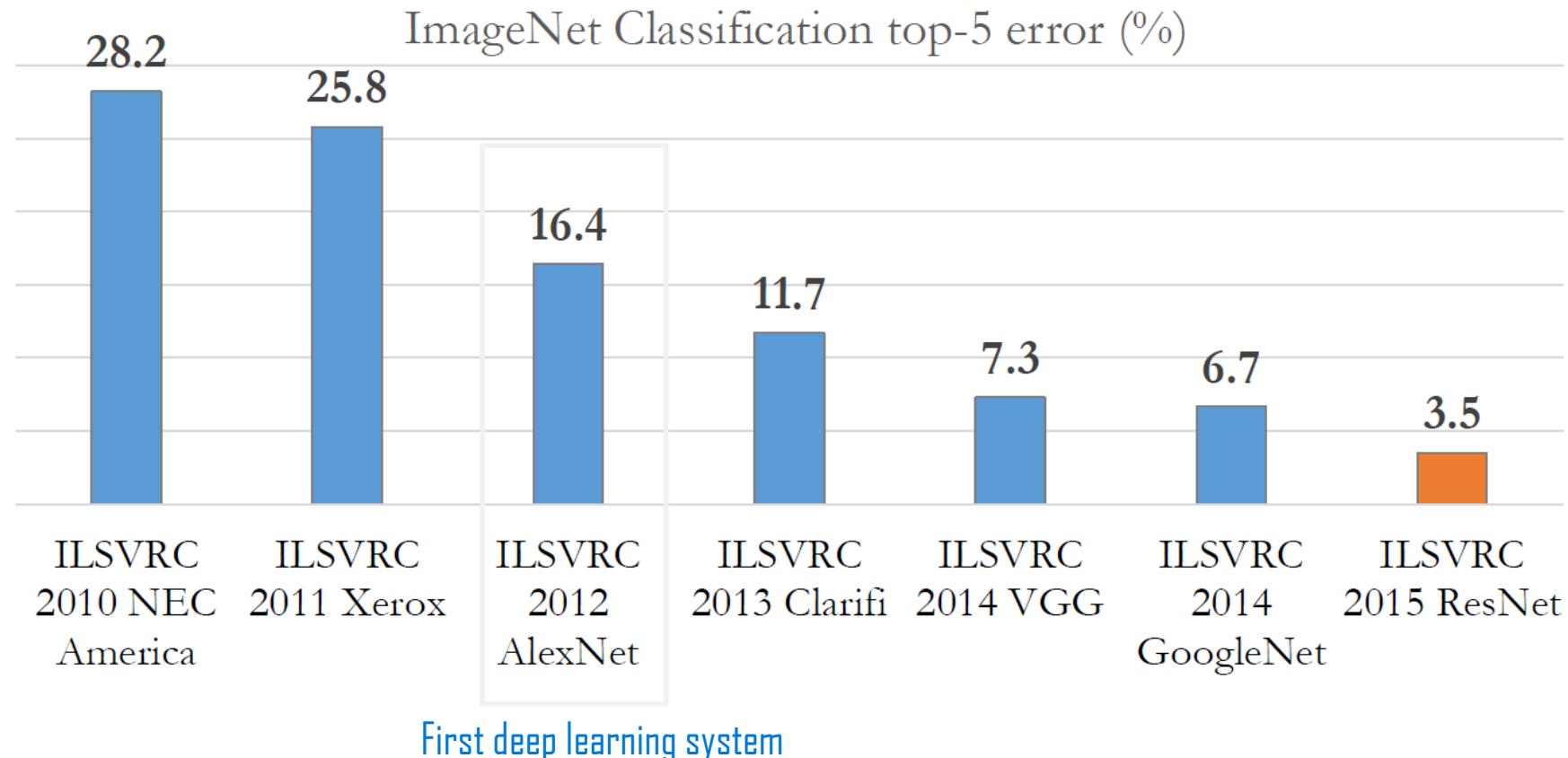




2017: ~5%!
Human error: ~5%

IMPROVEMENTS IN SPEECH RECOGNITION

IMPROVEMENTS IN COMPUTER VISION



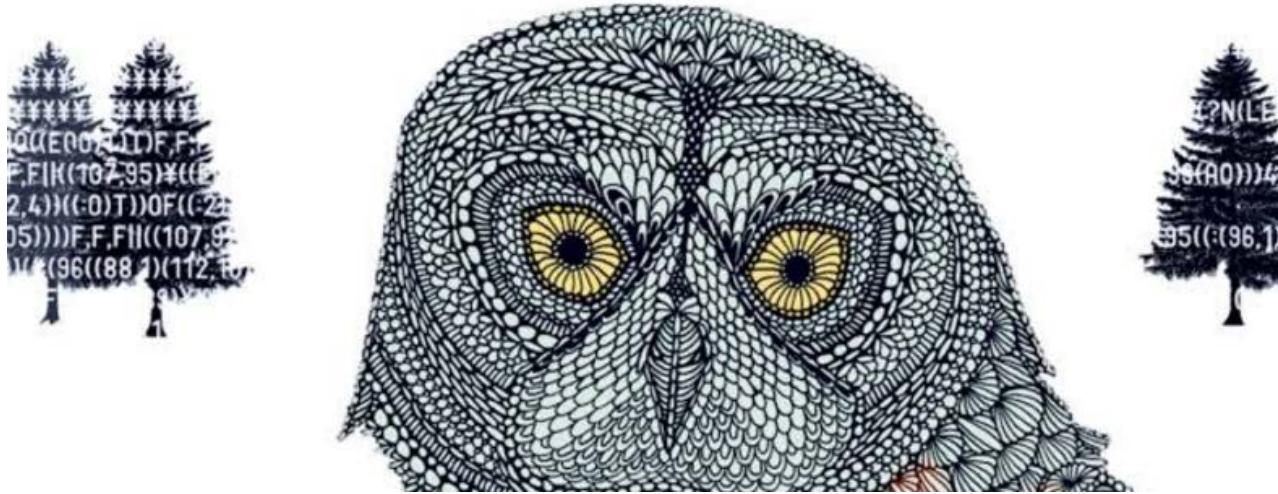
ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

2017: ~2.2%

NICK BOSTROM

SUPERINTELLIGENCE

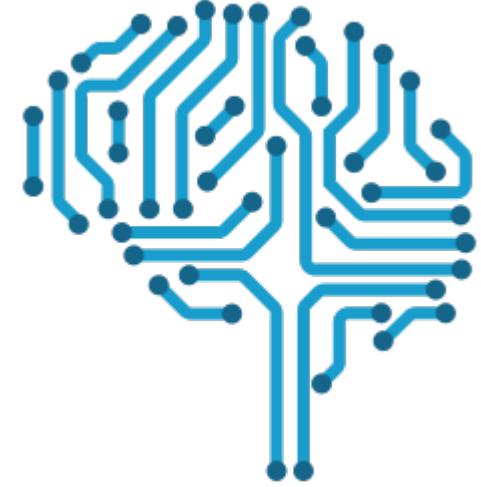
Paths, Dangers, Strategies



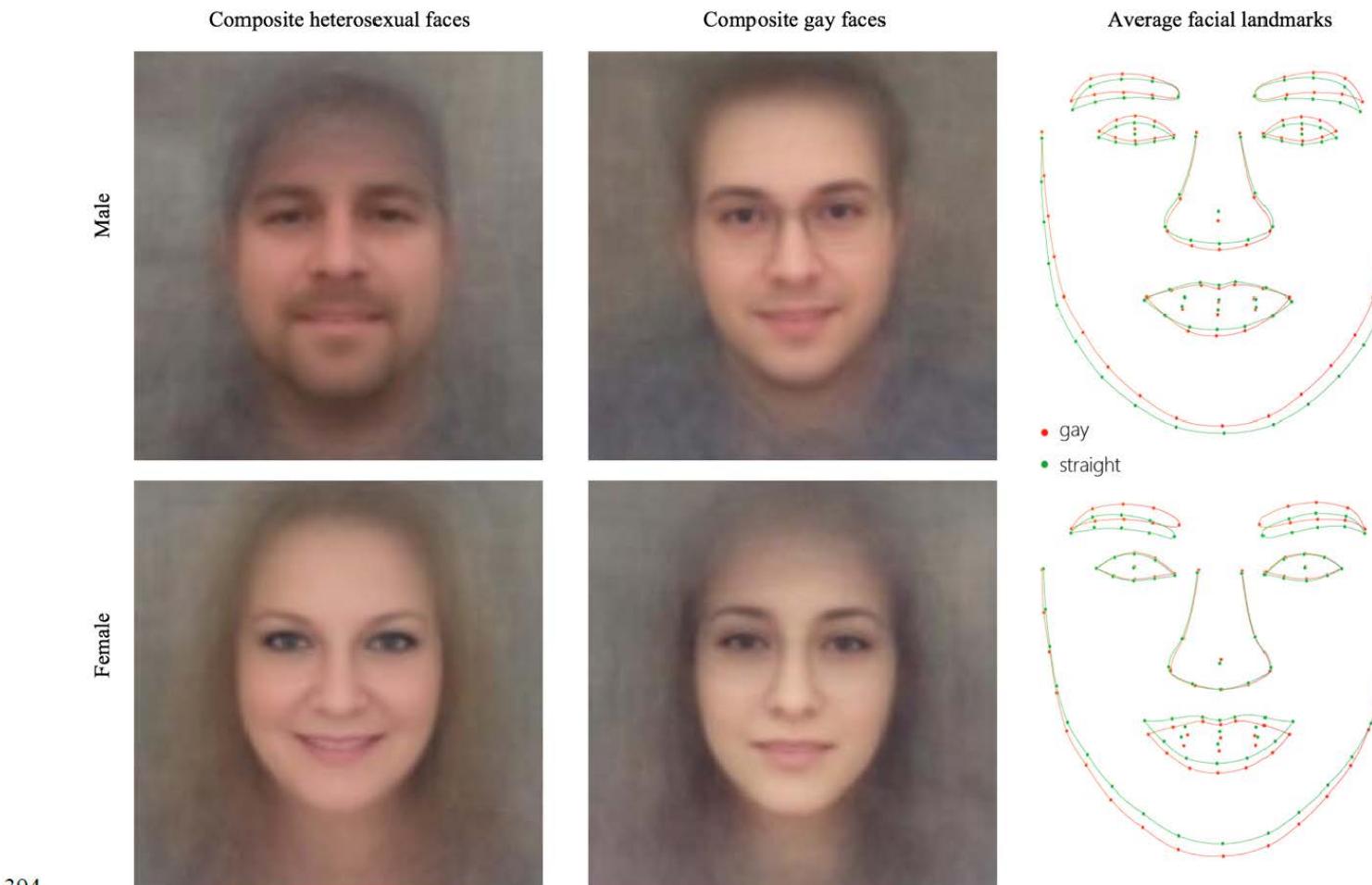
DON'T WORRY ABOUT GENERAL INTELLIGENCE



1. Discrimination
2. Opaque AI
3. Data is not neutral
4. Manipulating markets and consumers/voters
5. Lack of human connection
6. Engineers are not philosophers (moral reasoning)
7. Privacy
8. Innocent till proven guilty?



DO WORRY ABOUT ETHICS IN AI



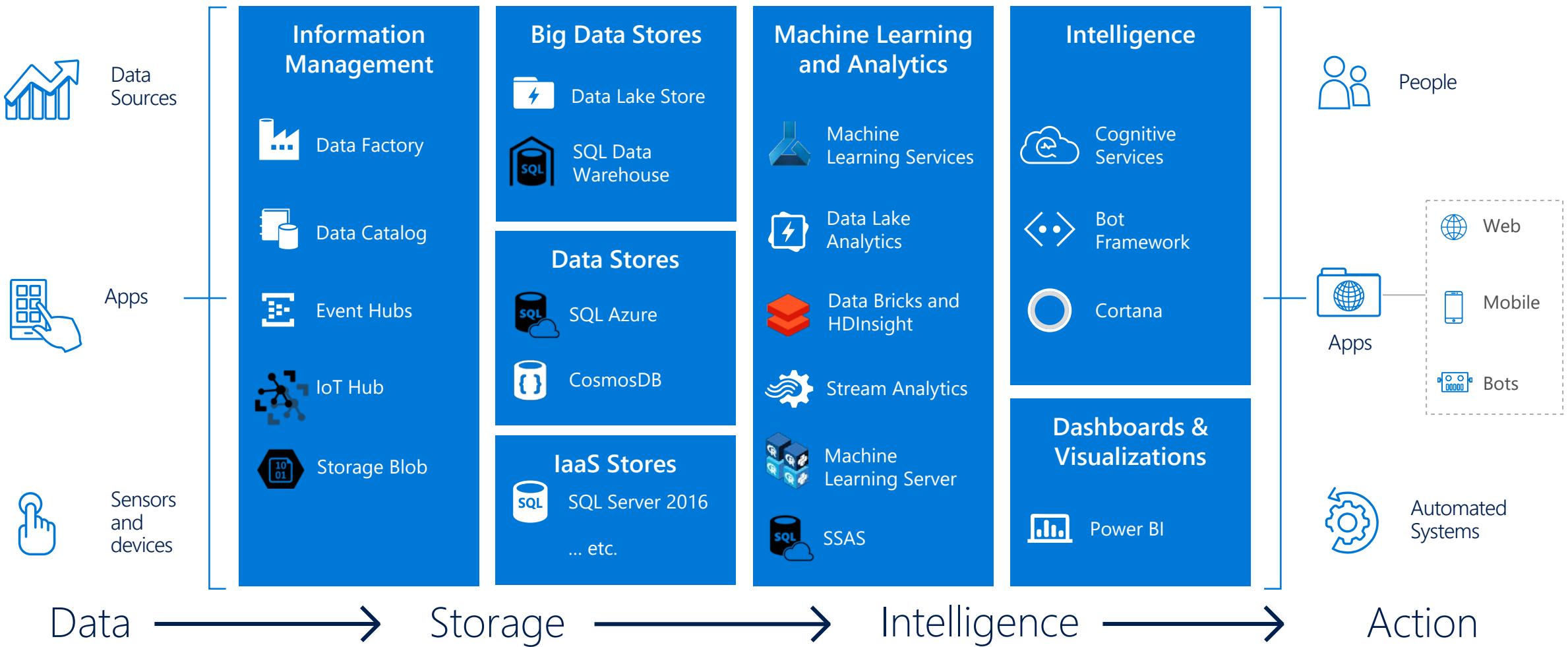
394

395 *Figure 4.* Composite faces and the average facial landmarks built by averaging faces classified as most and least likely to be gay.

2017 Wang, Y., & Kosinski, M.

PHYSIOGNOMY IS BACK?

AZURE INTELLIGENT CLOUD



AZURE MACHINE LEARNING STACK

Azure
Notebooks



CosmosDB / Data Lake Analytics / SQL Server

Pre-baked models



databricks

Cognitive Services
Software developer

Data Bricks
Data scientist

ML Studio
Citizen data scientist



Batch AI
Training



MMLSpark

Code-oriented



Data Science
virtual machine



Cognitive
Toolkit

SaaS

Notebooks

Discriminative

PaaS

Private
Containers

IaaS

Generative

Code-first

Onprem/Hybrid

Framework

ML Services
Data scientist

ML Server
Data scientist

CNTK
Applied Scientist

COGNITIVE SERVICES



Vision



Language



Speech



Search



Knowledge

Computer vision

Face

Emotion

Content Moderator

Video

Video Indexer

Cognitive Services Labs

Text analytics

Spell check

Web language model

Linguistic analysis

Translator

Speaker recognition

Speech

Web search

Image search

Video search

News search

Autosuggest

Academic knowledge

Entity linking service

Knowledge exploration

Recommendations

QnA maker

Custom

Vision Service

Custom

Language
Understanding

Custom

Speech Service

Custom

Search

Custom

Decision Service

MACHINE LEARNING PORTFOLIO

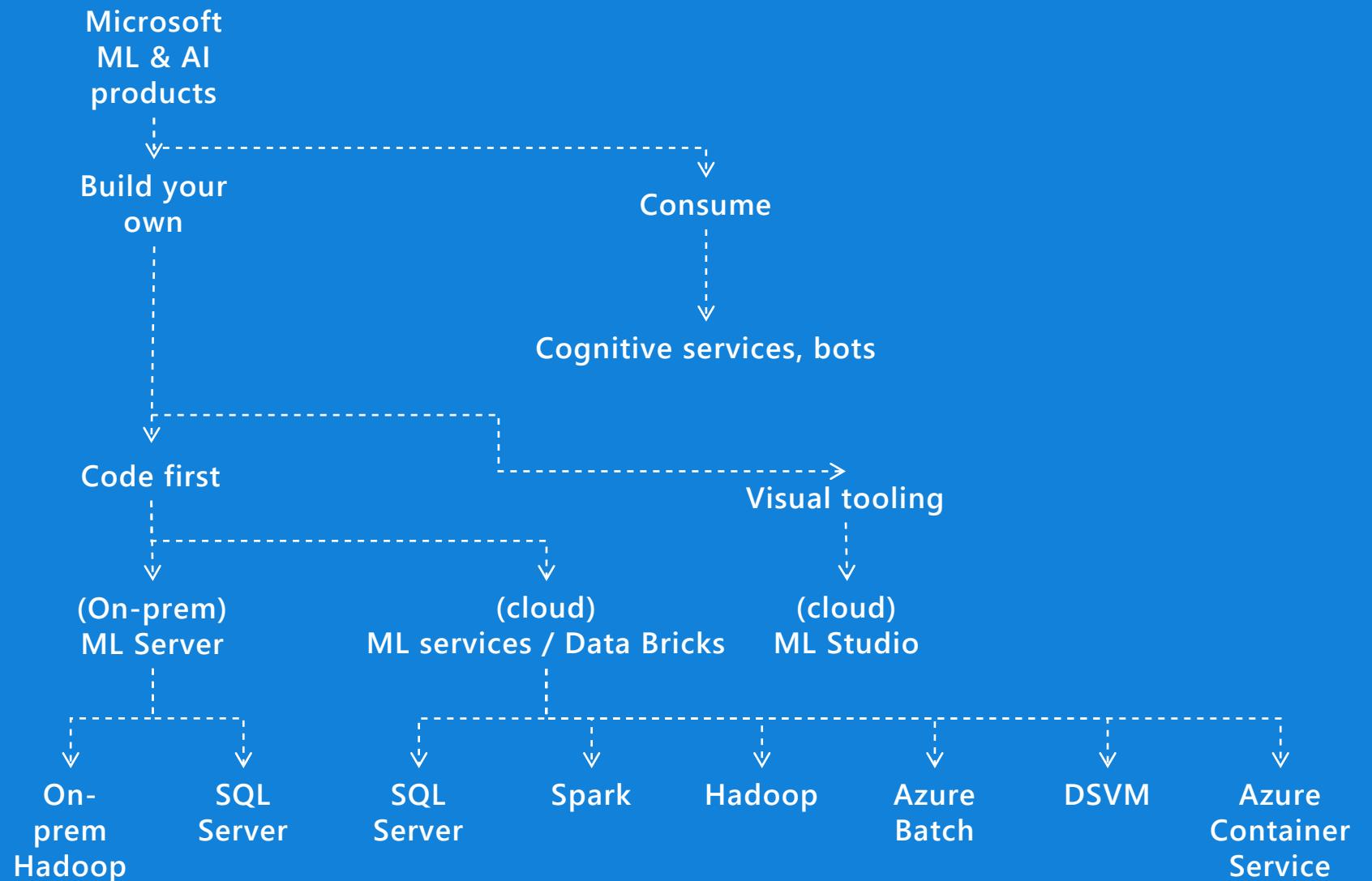


Build your own or consume pre-trained models?

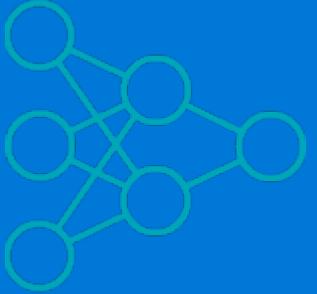
Which experience do you want?

Deployment target

What engine(s) do you want to use?



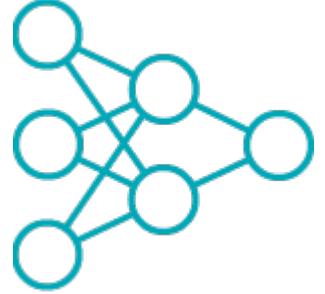
WHAT IS MACHINE LEARNING?



A machine learning from previous experience (data) so that it can perform some task better *based on that experience.*



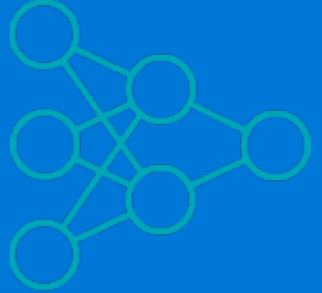
MACHINE LEARNING != DATA SCIENCE



The data scientist is trying to understand the data, the ML person is building a working model which may not be interpretable.

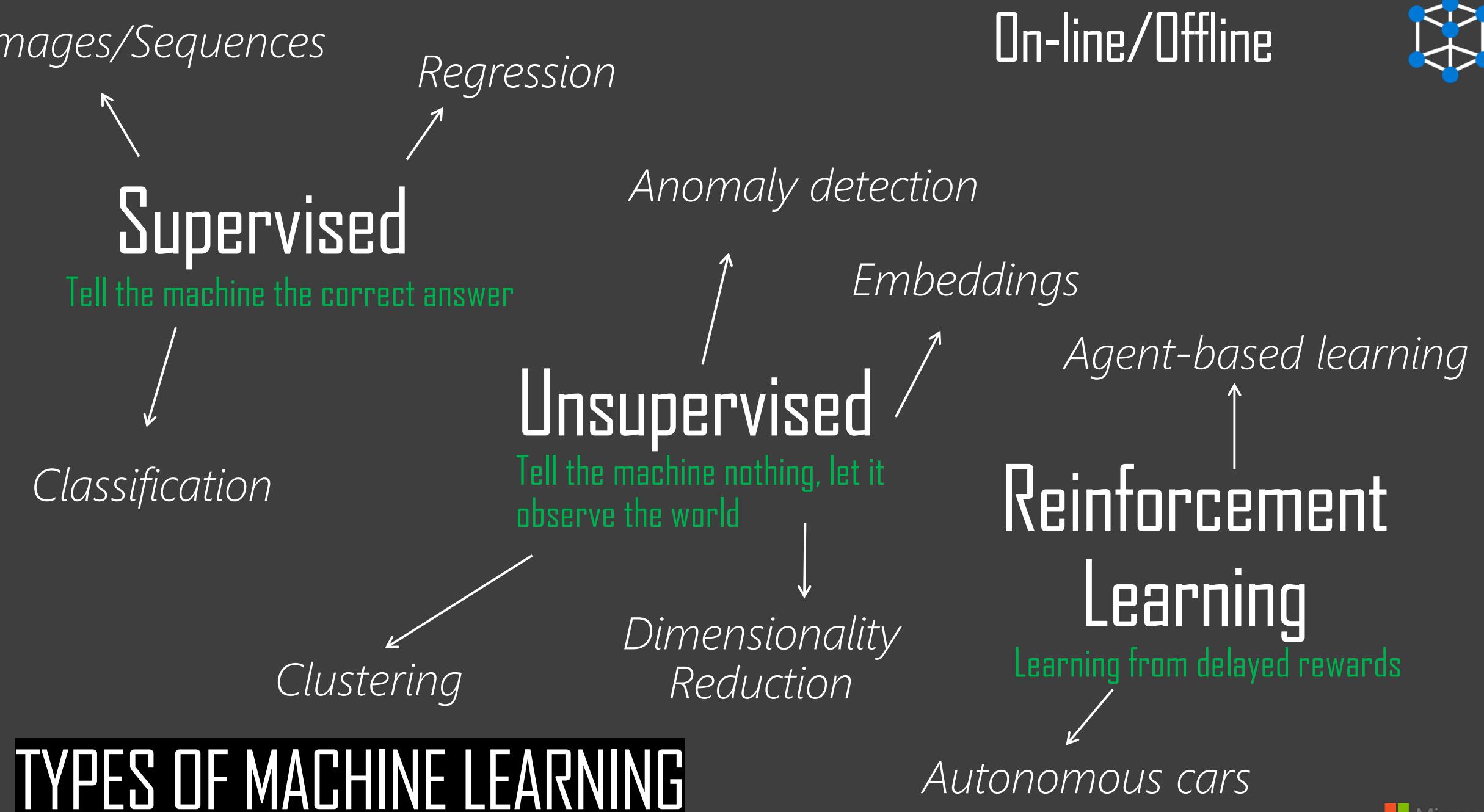


INDUCTIVE INFERENCE



ML algorithms often "induce" a generic and separate artefact or "decision function", learned from data.





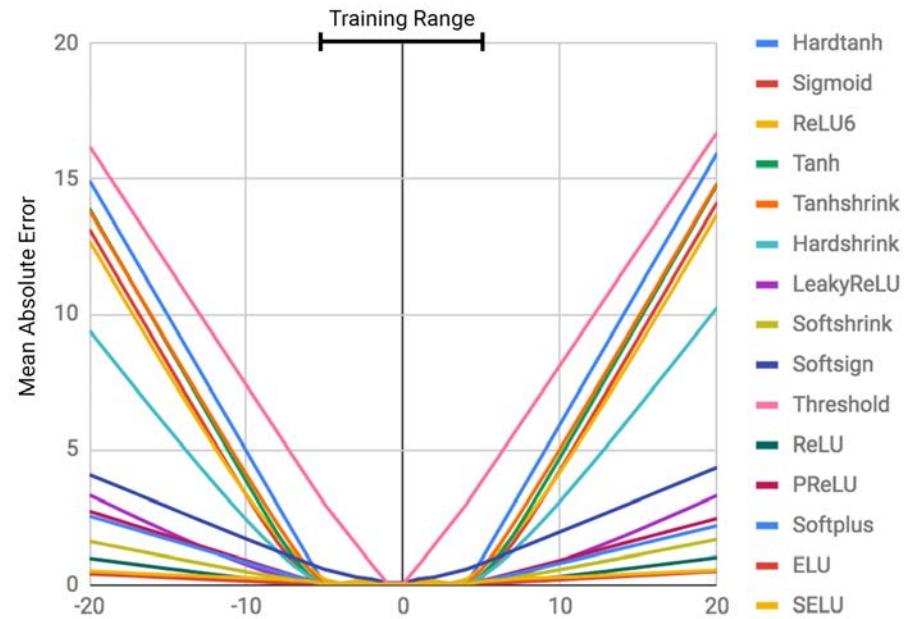


Machine learning approximates a function, which we may only have a few examples of.

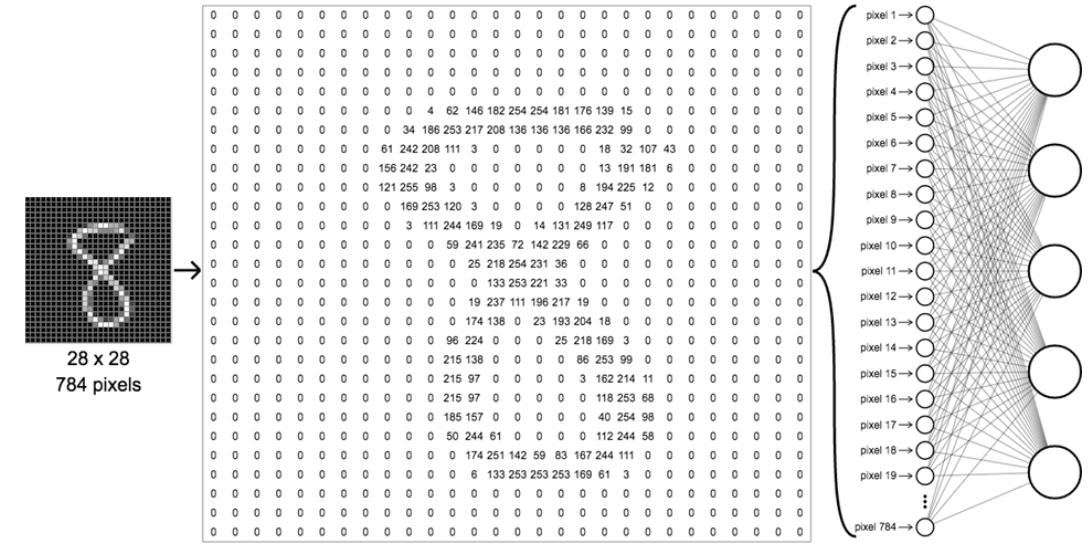
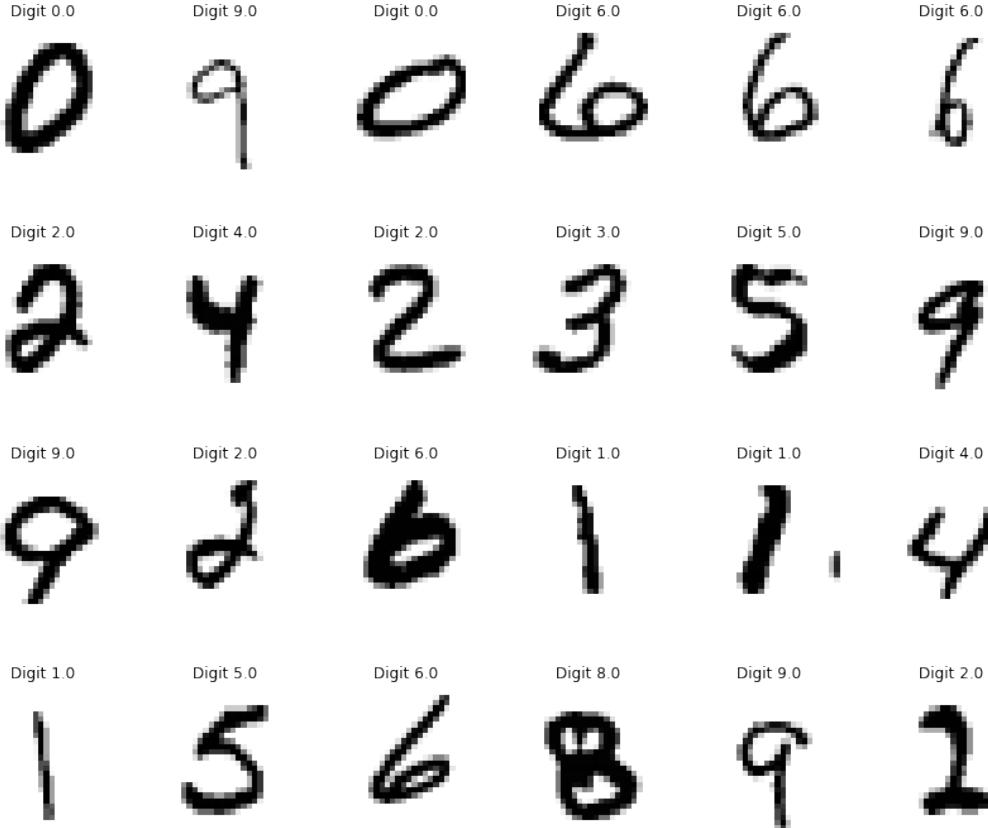
A SIMPLE SINEFUNCTION

INTERPOLATION VS EXTRAPOLATION

ML algorithms perform very, very badly on data out of range. They don't generalise as well as you think. Trask et al 2018 (DeepMind)



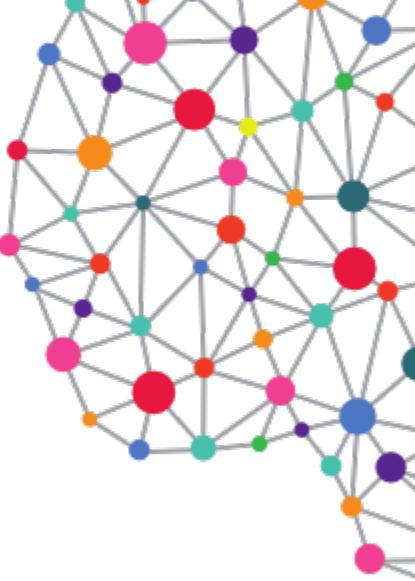
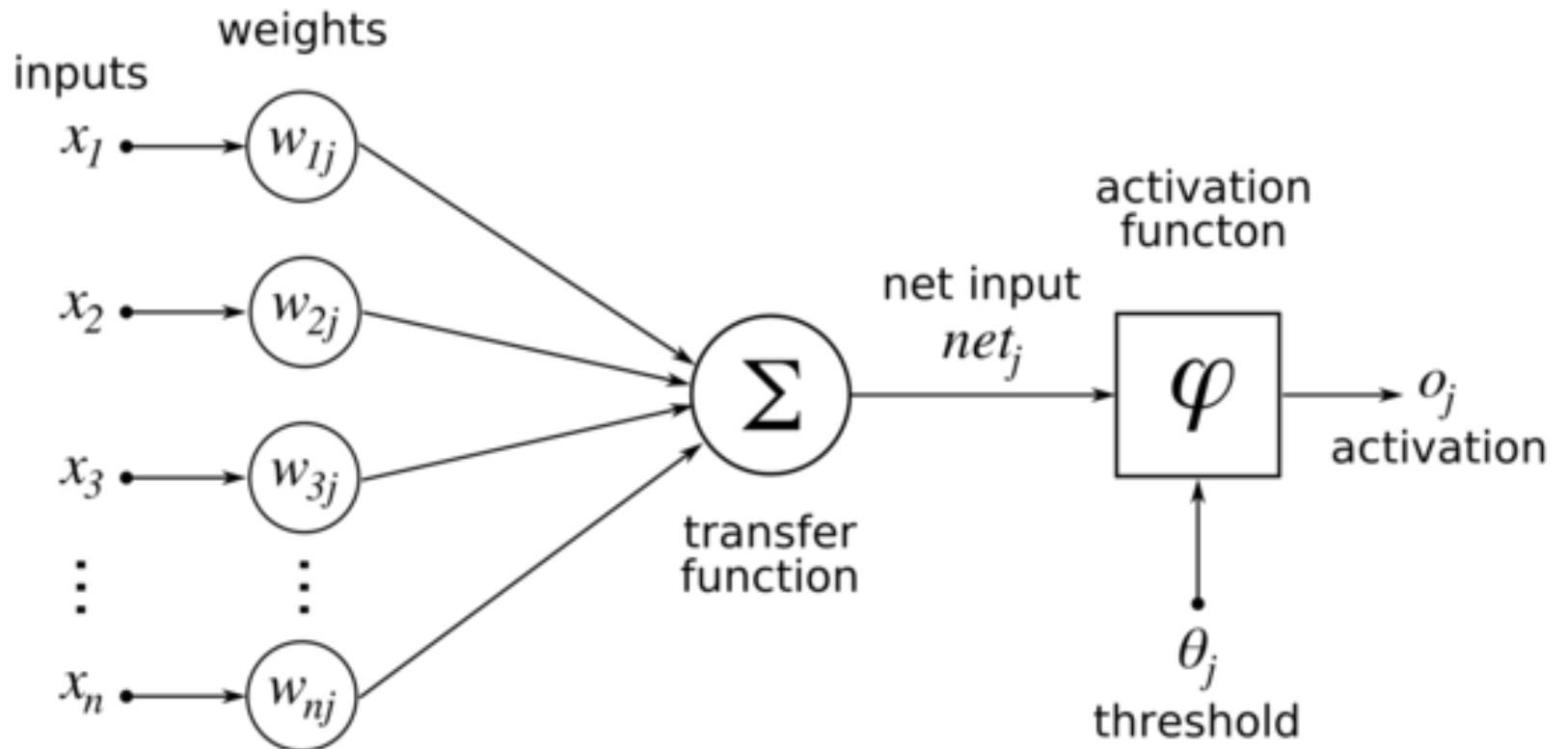
MNIST DIGIT CLASSIFICATION



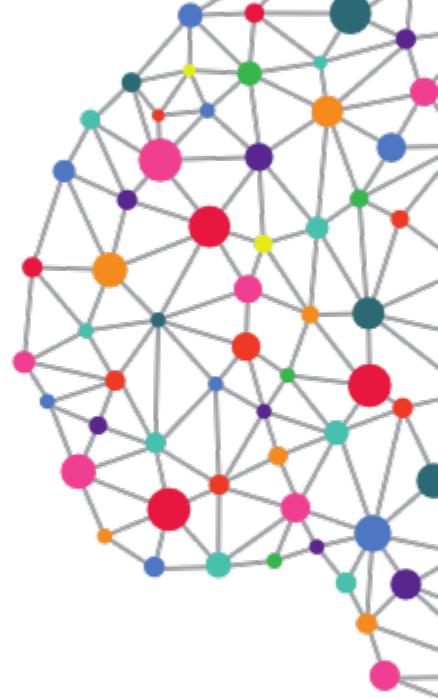
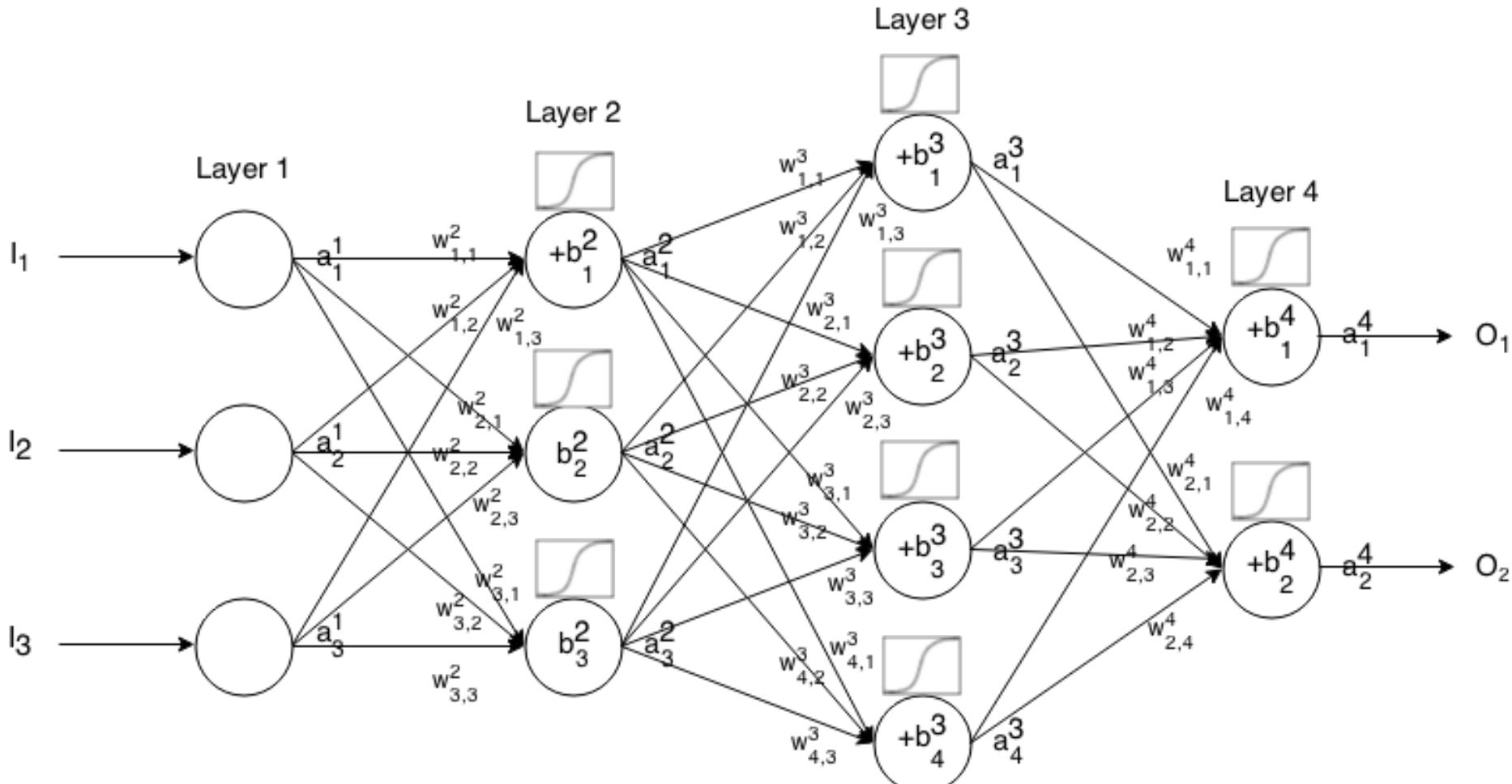
Represent input image as a vector $\mathbf{x} \in \mathbb{R}^{784}$
 Learn a classifier $f(\mathbf{x})$ such that,

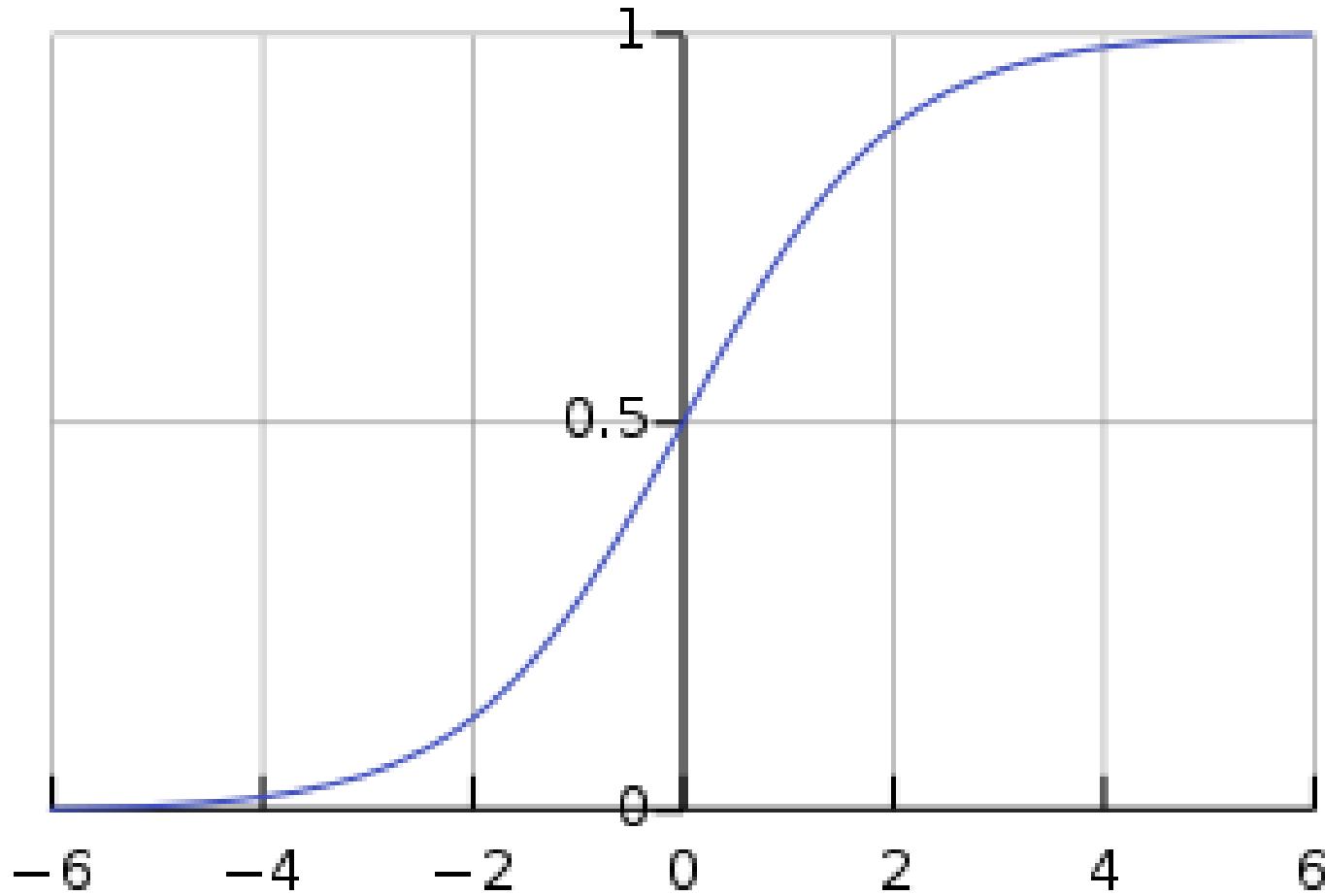
$$f : \mathbf{x} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

WHAT IS A NEURAL NETWORK?



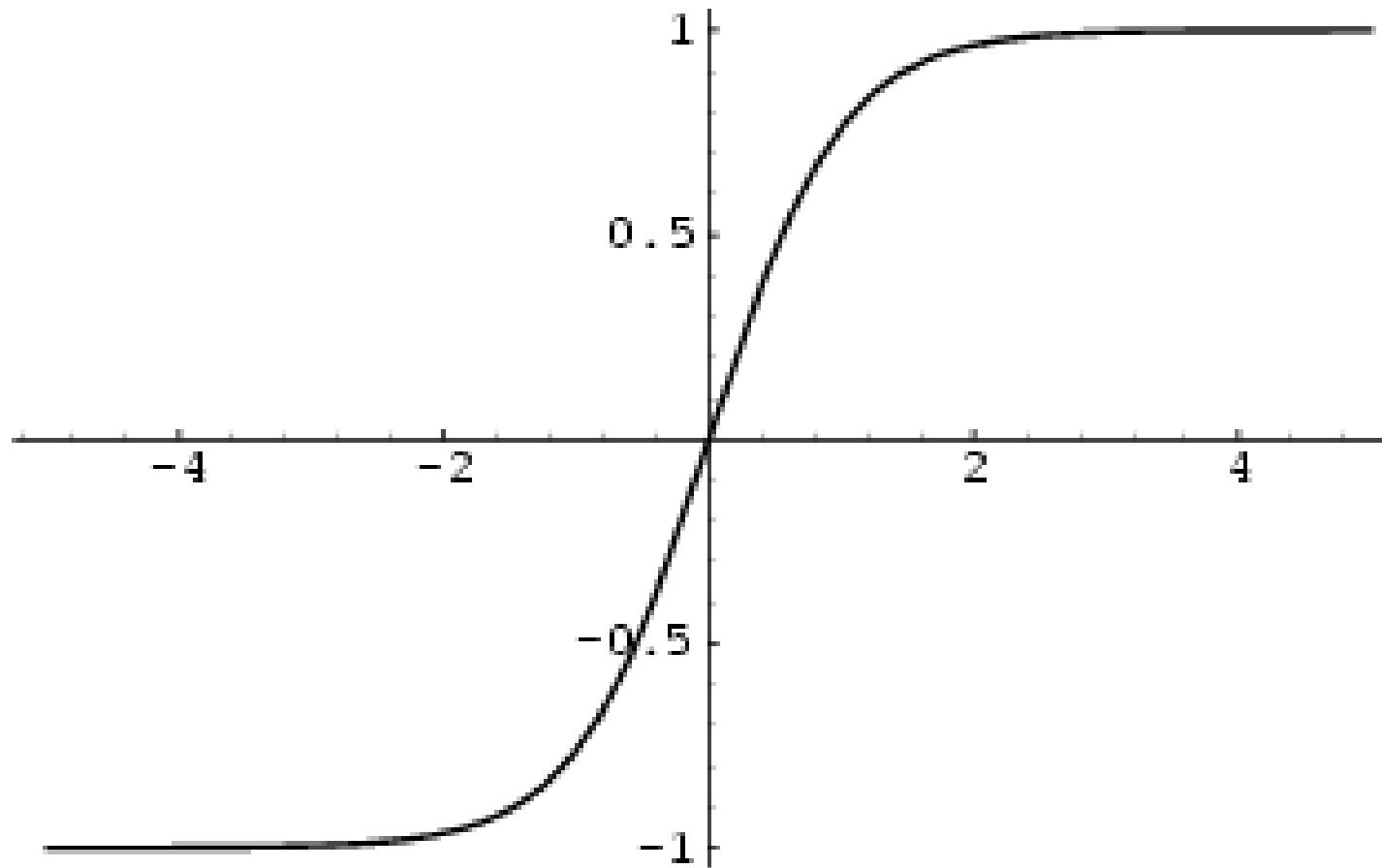
WHAT IS A DEEP NEURAL NETWORK?





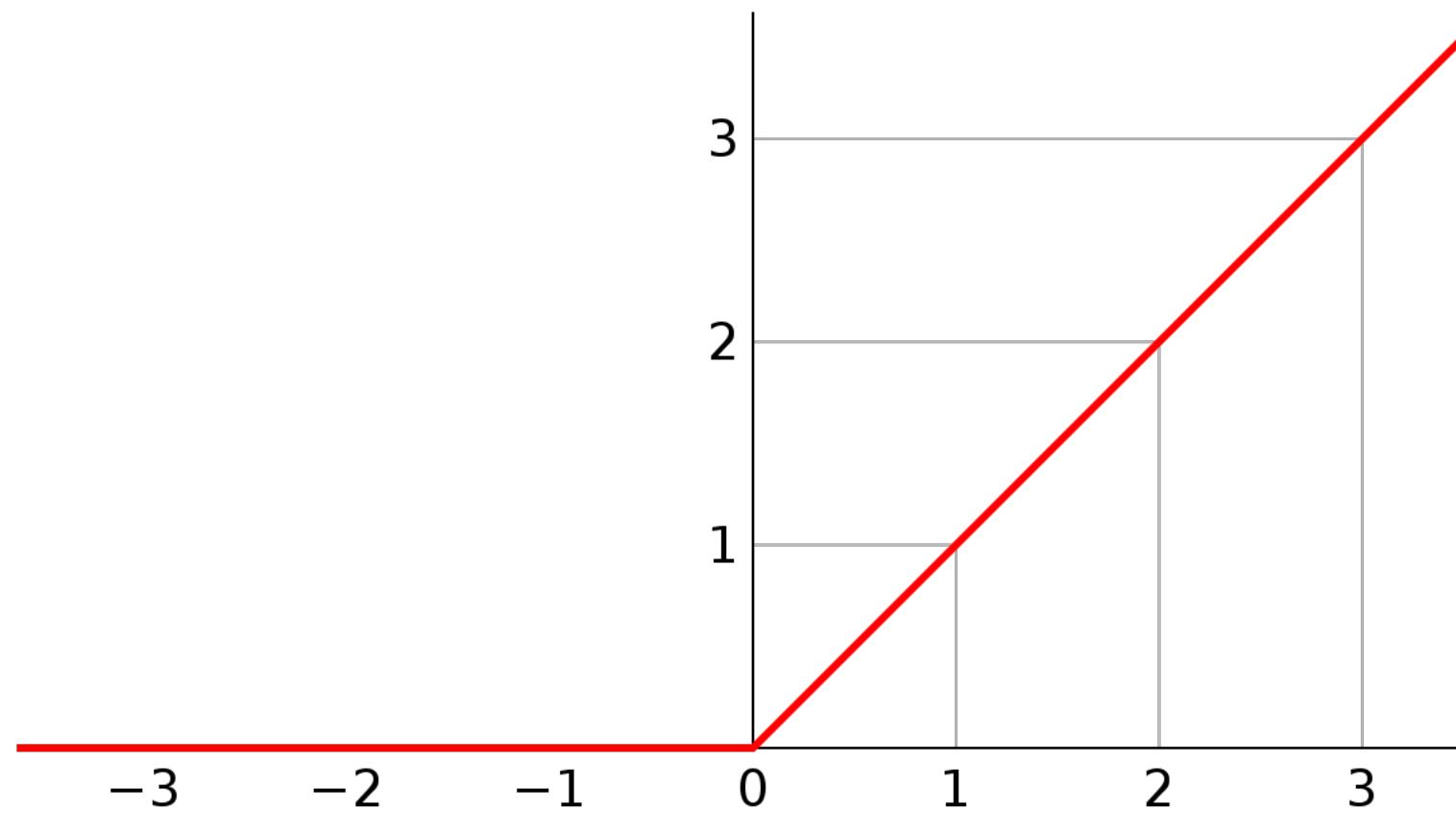
$$\frac{e^x}{e^x + 1}$$

SIGMOID SQUASHING FUNCTION

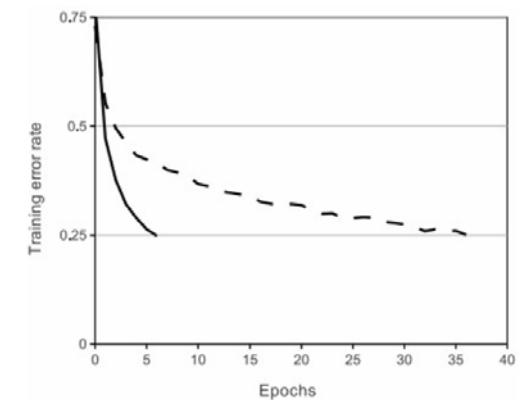


$$\frac{1 - e^{-2x}}{1 + e^{-2x}}$$

TANH SQUASHING FUNCTION

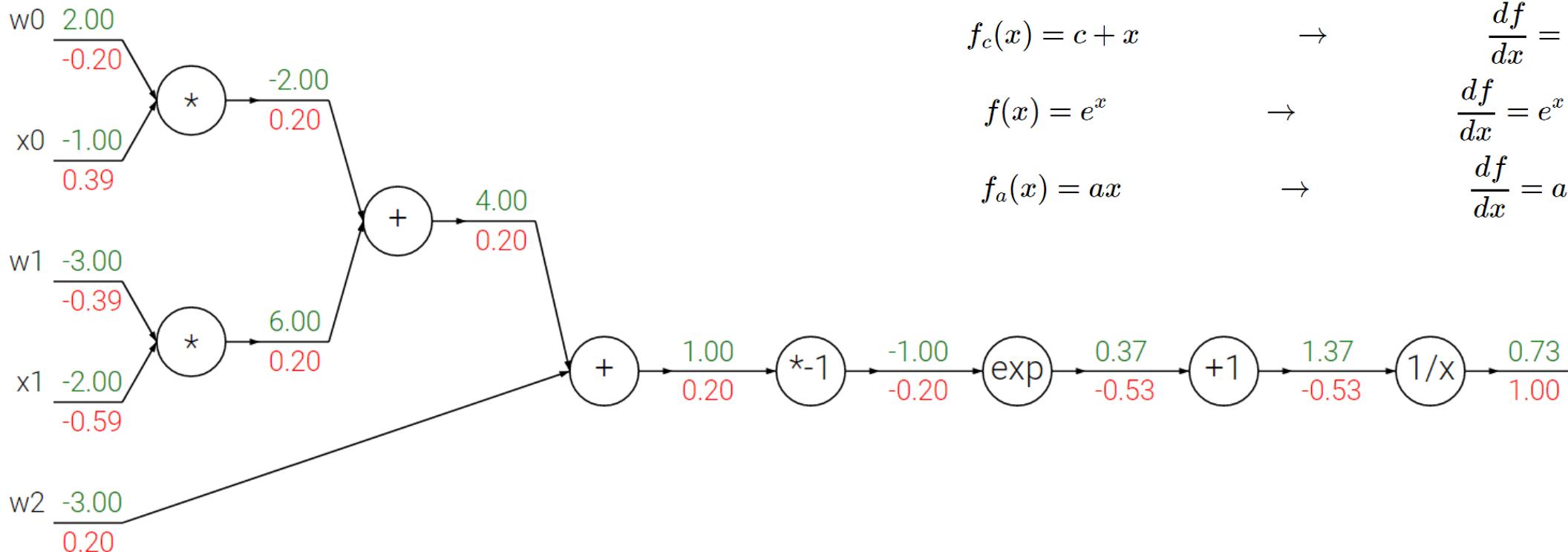


$$f(x) = x^+ = \max(0, x)$$



RELU SQUASHING FUNCTION

BACKPROPAGATION



Example circuit for a 2D neuron with a sigmoid activation function. The inputs are $[x_0, x_1]$ and the (learnable) weights of the neuron are $[w_0, w_1, w_2]$. As we will see later, the neuron computes a dot product with the input and then its activation is softly squashed by the sigmoid function to be in range from 0 to 1.

$$\frac{\partial E}{\partial w_{jk}}$$

$$\Delta w_{jk} = \eta * [x_j * (o_k - t_k) * o_k * (1 - o_k)]$$

learning
rate

error e_k

derivative of output activation φ'_k

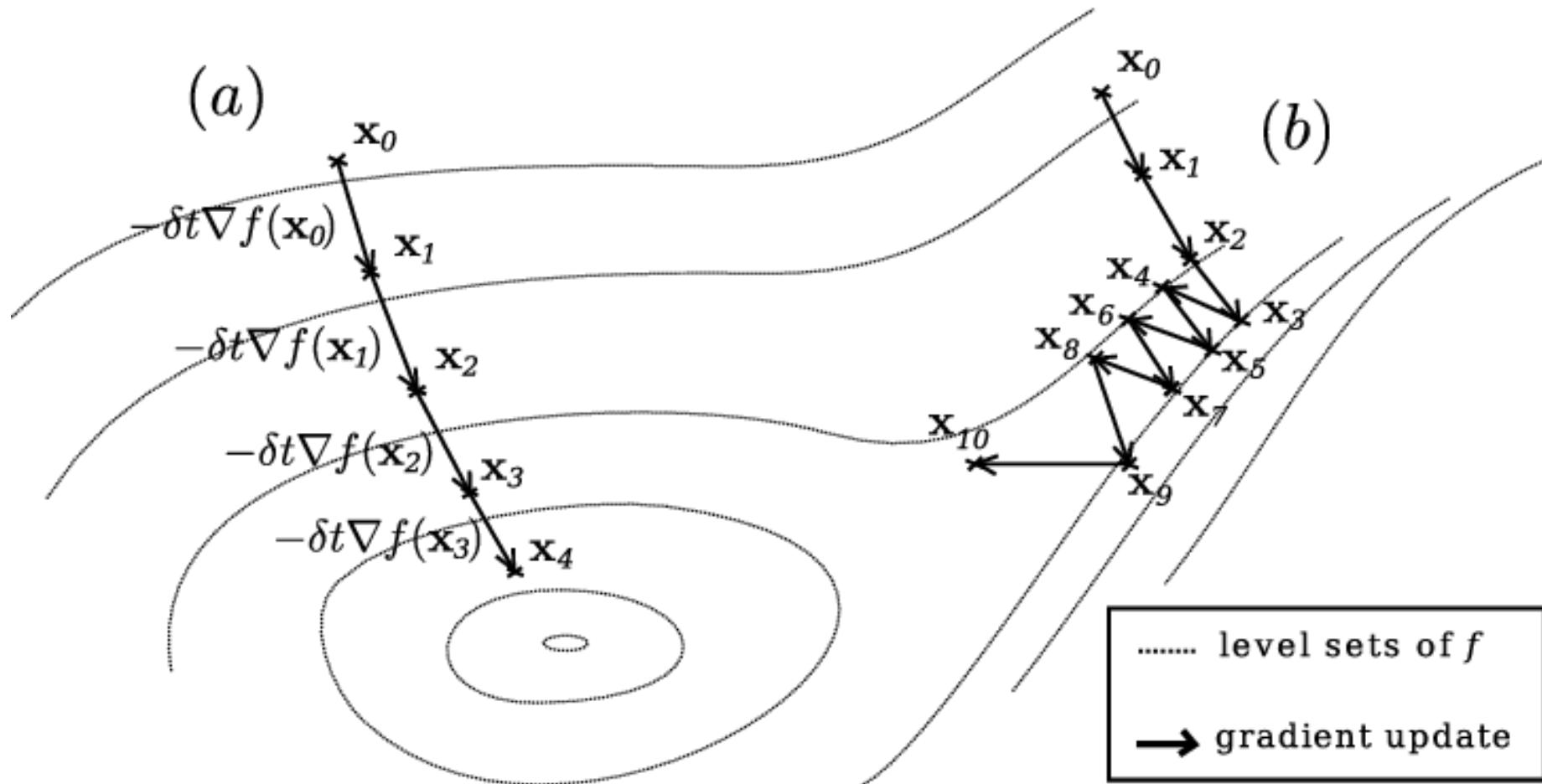
signal δ_k

WEIGHT UPDATE

```
loop maxEpochs times
    for-each training item
        get target values
        compute output values
        compute the gradient of each weight
        use gradient to compute delta for each weight
        update each weight using its delta
    end-for
end-loop
```

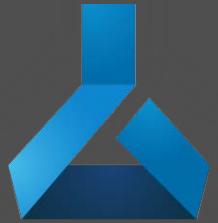
BACKPROP ALGORITHM

OPTIMIZATION/GRADIENT DESCENT

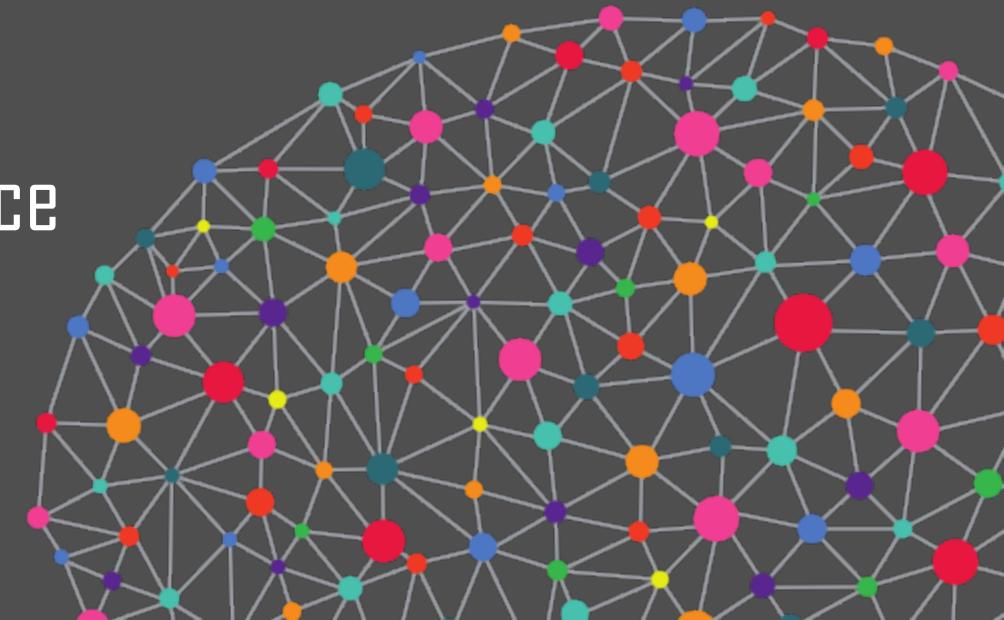


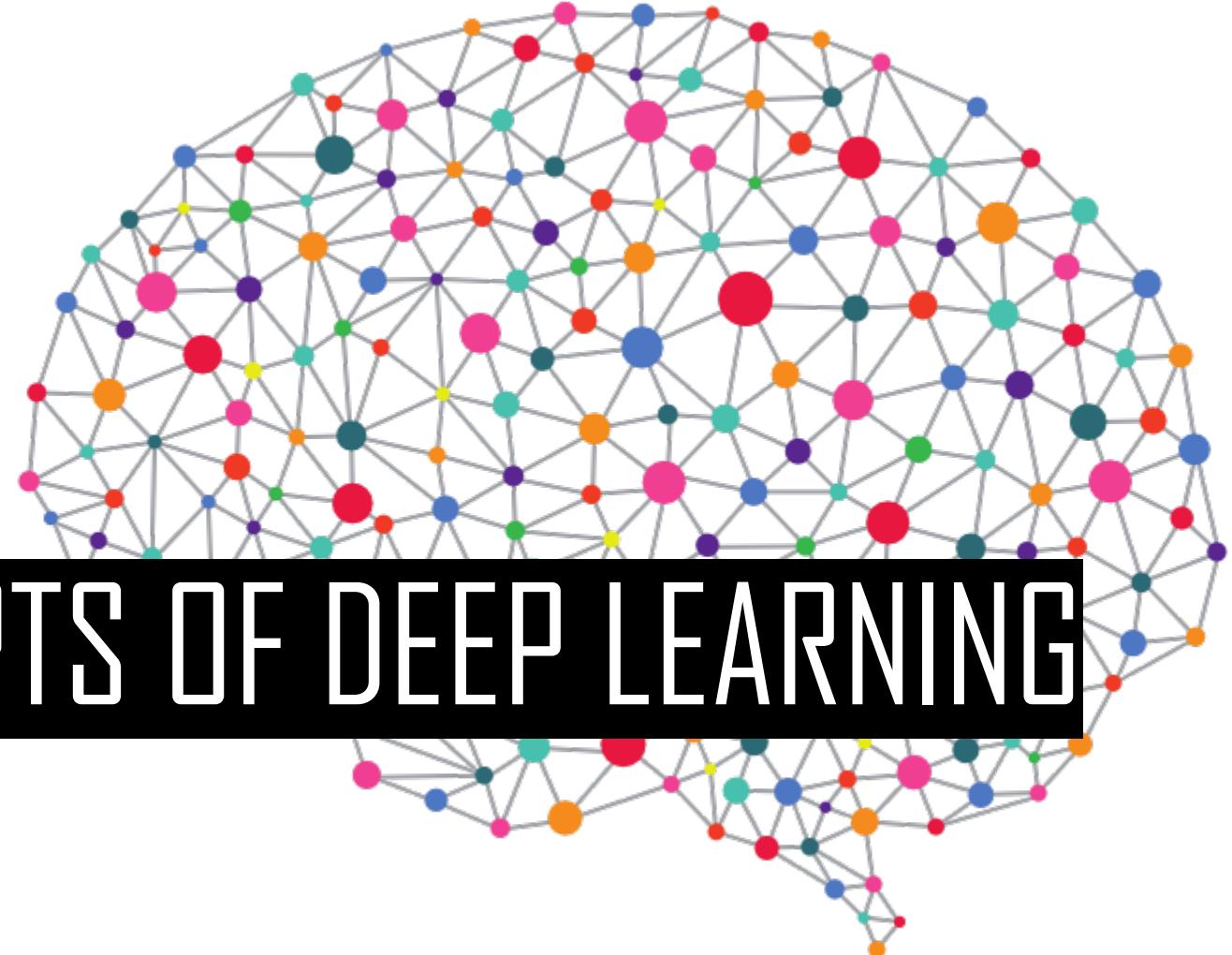
WEIGHT UPDATES IN EPOCHS / TRAINING DATA SPLIT INTO MINIBATCHES

UNIVERSAL MACHINE LEARNING PROCESS



- Define problem
- Define success
- Validation process
- Vectorise/normalize data
- Develop naïve model
- Refine model based on validation performance

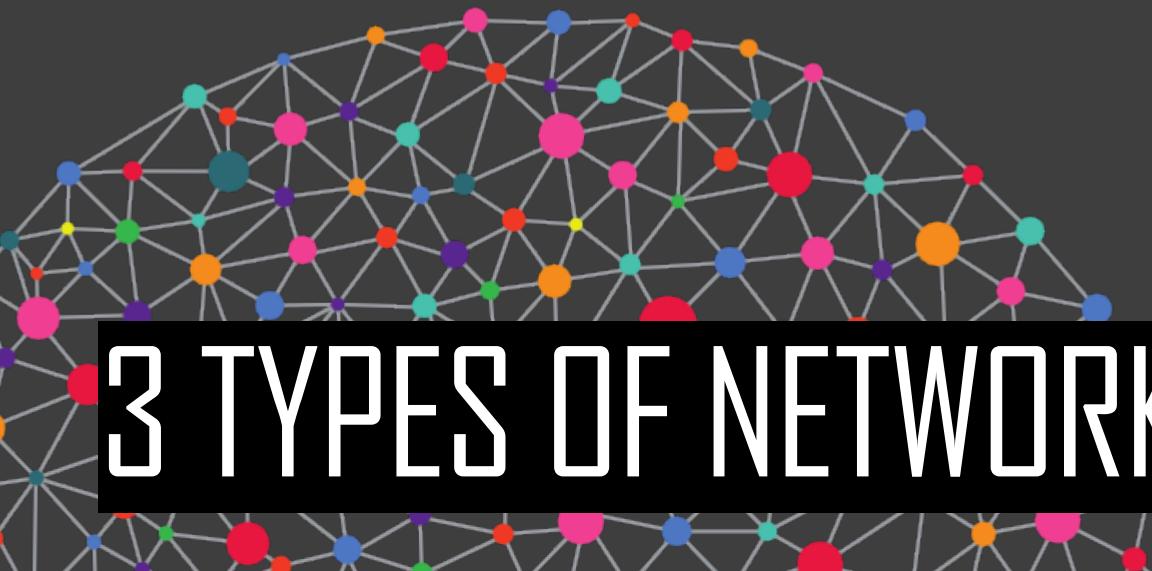




DISTILLED CONCEPTS OF DEEP LEARNING



1. DENSELY CONNECTED NETWORKS
2. CONVOLUTIONAL NEURAL NETWORKS (MODEL SPACE)
3. RECURRENT NEURAL NETWORKS (MODEL TIME)



3 TYPES OF NETWORK ARCHITECTURE



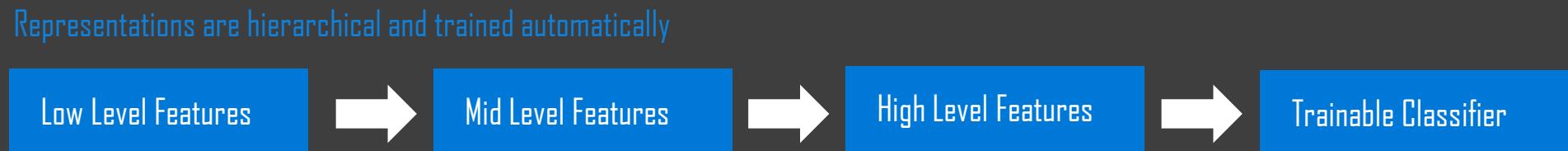
- The networks have **many levels** of depth
- Machine learns a hierarchy of representations i.e. "**representation learning**"
- **Less feature extraction** required (still need transposition, normalization etc)



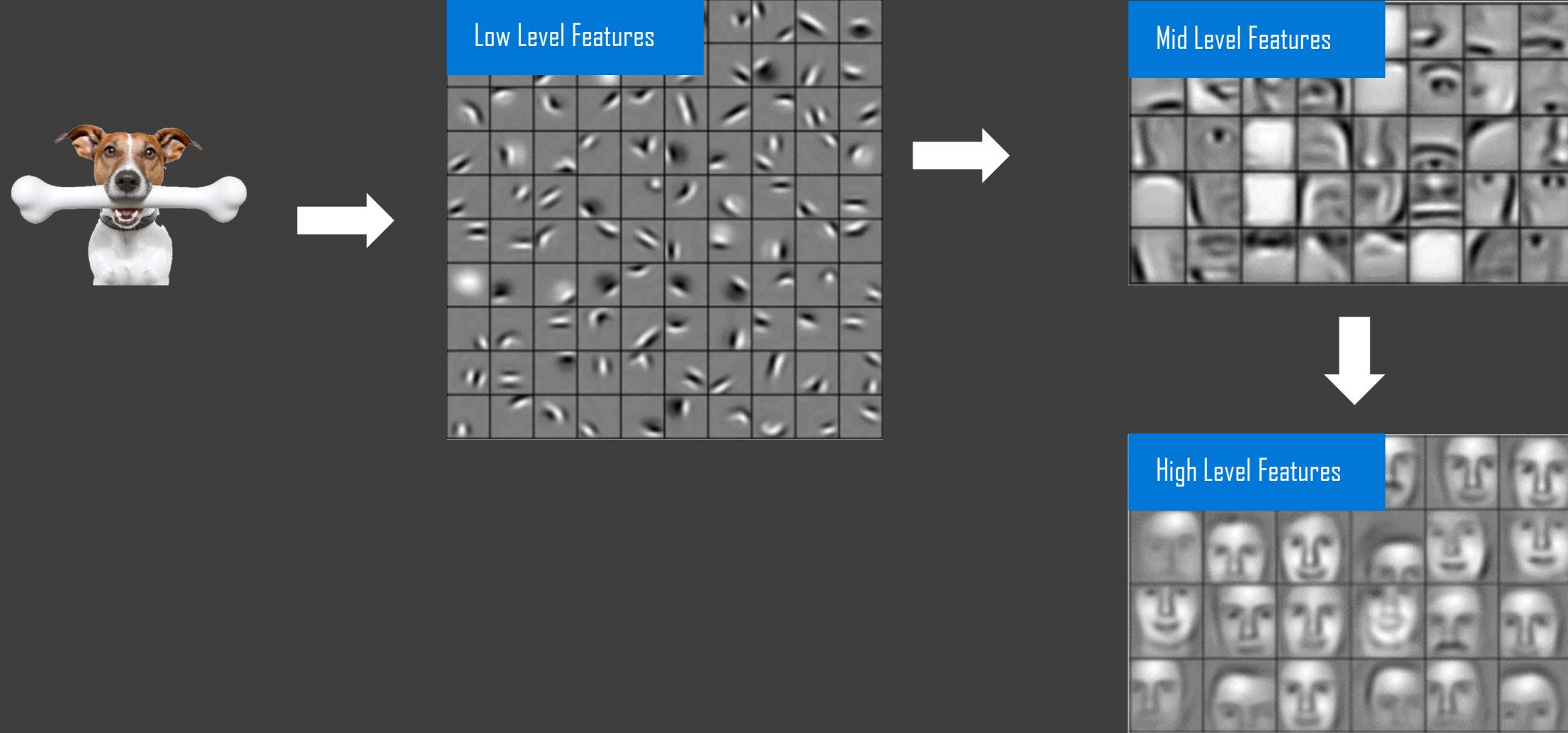
Traditional ML



Deep Learning

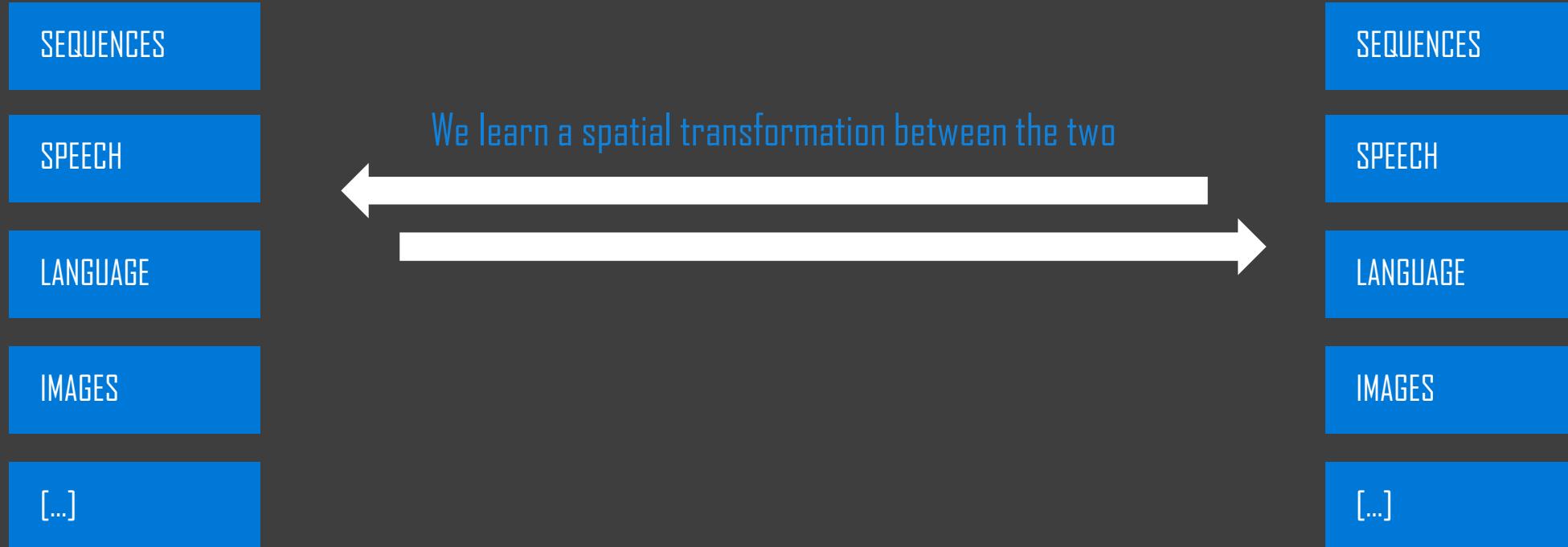


ENTIRE MACHINE IS TRAINABLE



REPRESENTATIONS ARE LEARNED AUTOMATICALLY

- Unlike other shallow ML algorithms; you can map *between data domains*
- **Generative** models are possible, not just **discriminative**



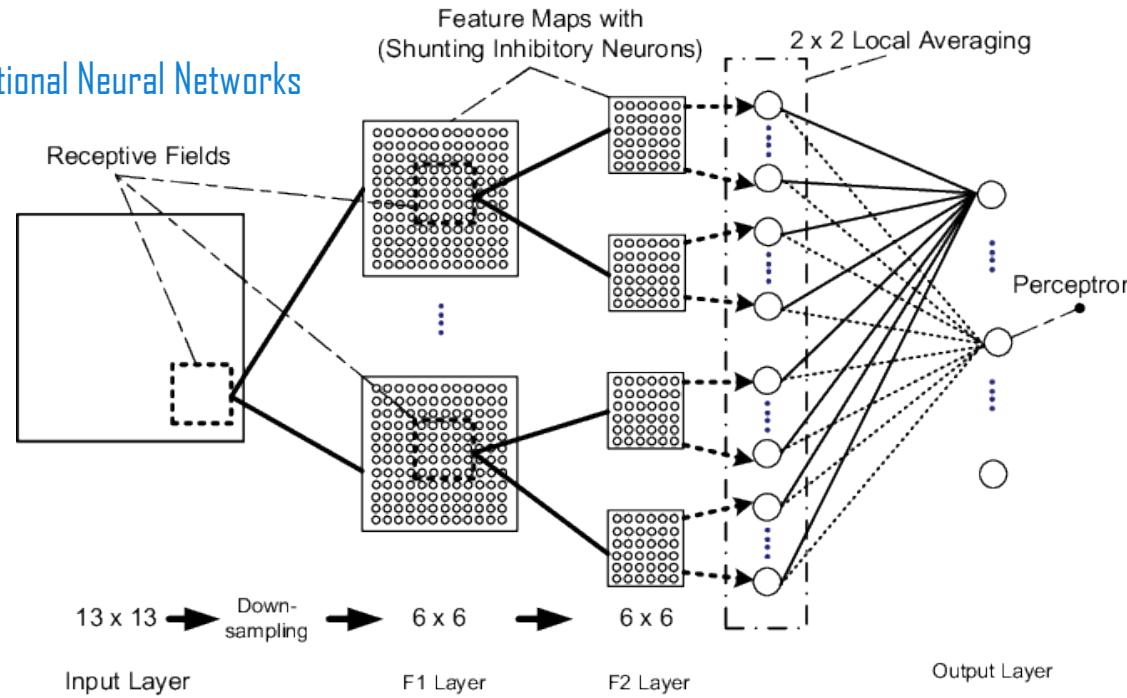
UNIVERSAL FUNCTION APPROXIMATORS



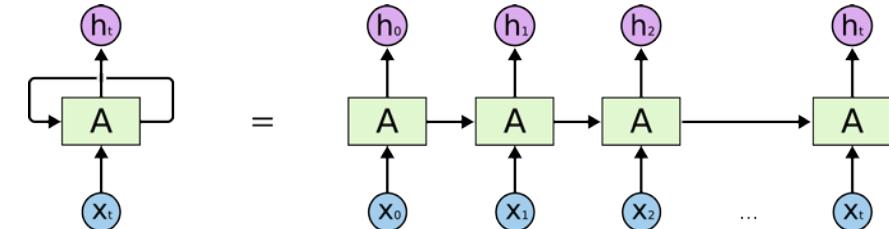
Unlike other algorithms, NNs can natively encode useful and obvious relationships in the data domain

- Local spatial dependencies (vision) i.e. CNNs
- Time dependencies (language, speech) i.e. RNNs

Convolutional Neural Networks

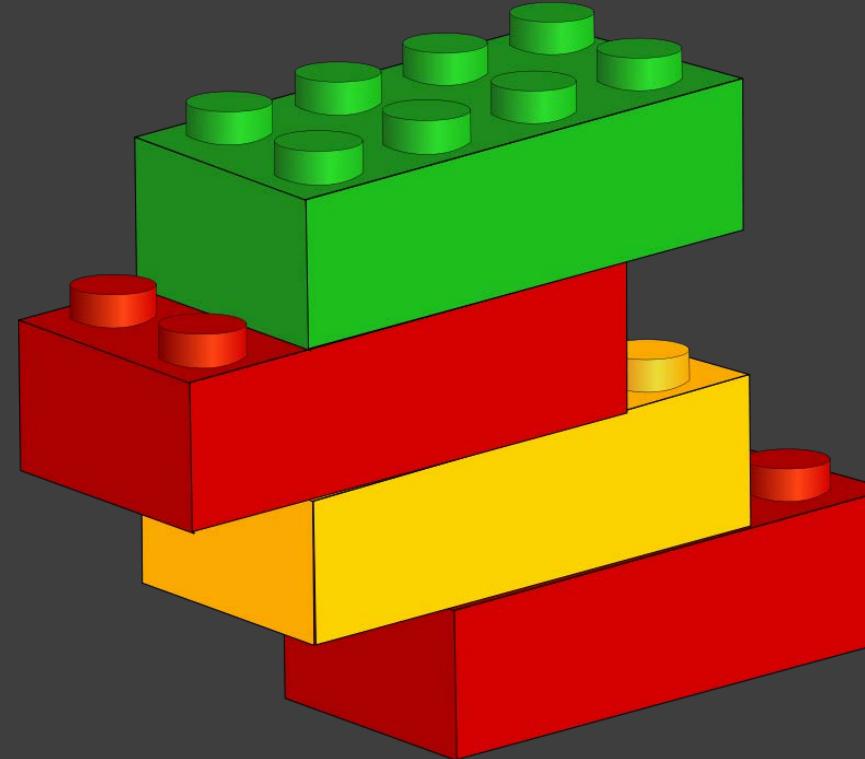


Recurrent Neural Networks



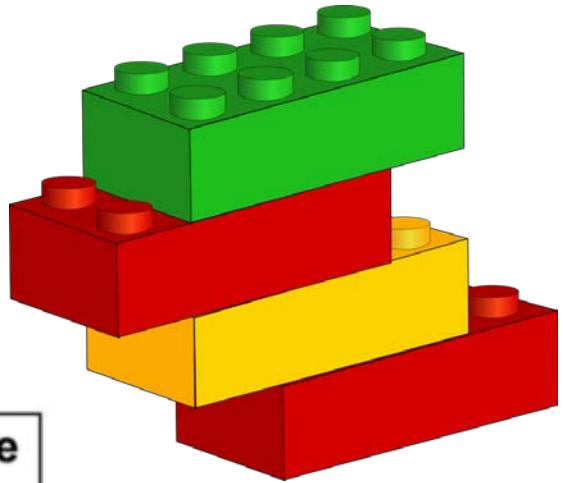
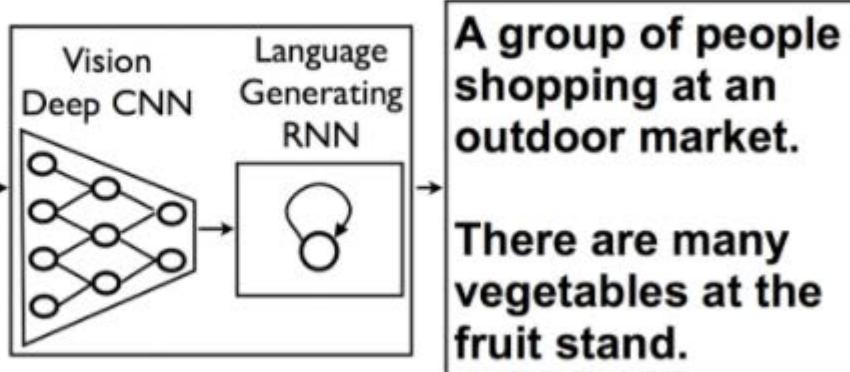
NATIVE DATA-DOMAIN FEATURES

- Composability
- ML is becoming a form of software development
- ML models are like software



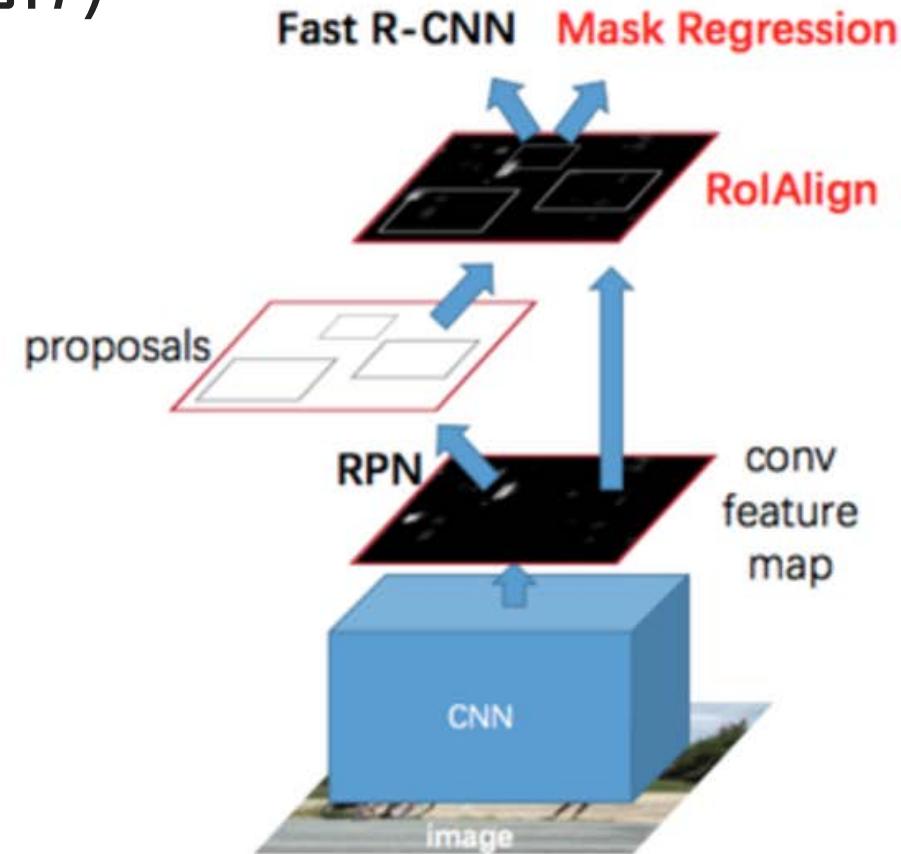
COMPOSABILITY

BUILDING PREDICTIVE ARCHITECTURES LIKE LEGO BLOCKS



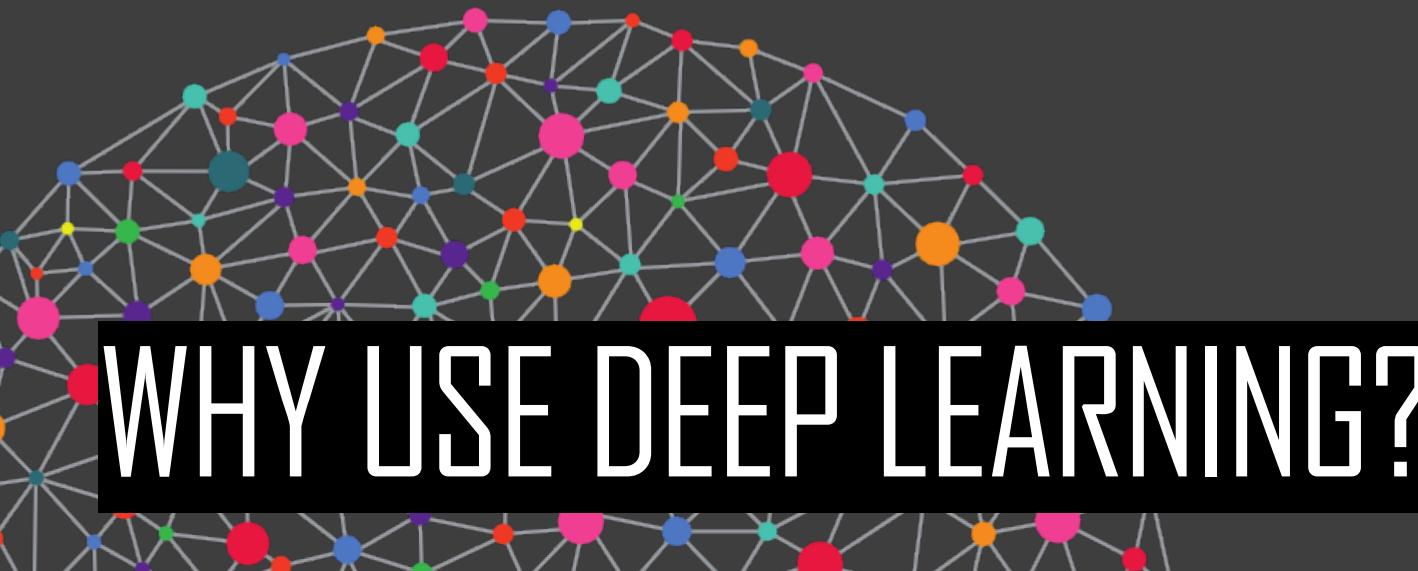
(Vinyals, et al. (2014))

MASK R-CNN ARCHITECTURE (2017)



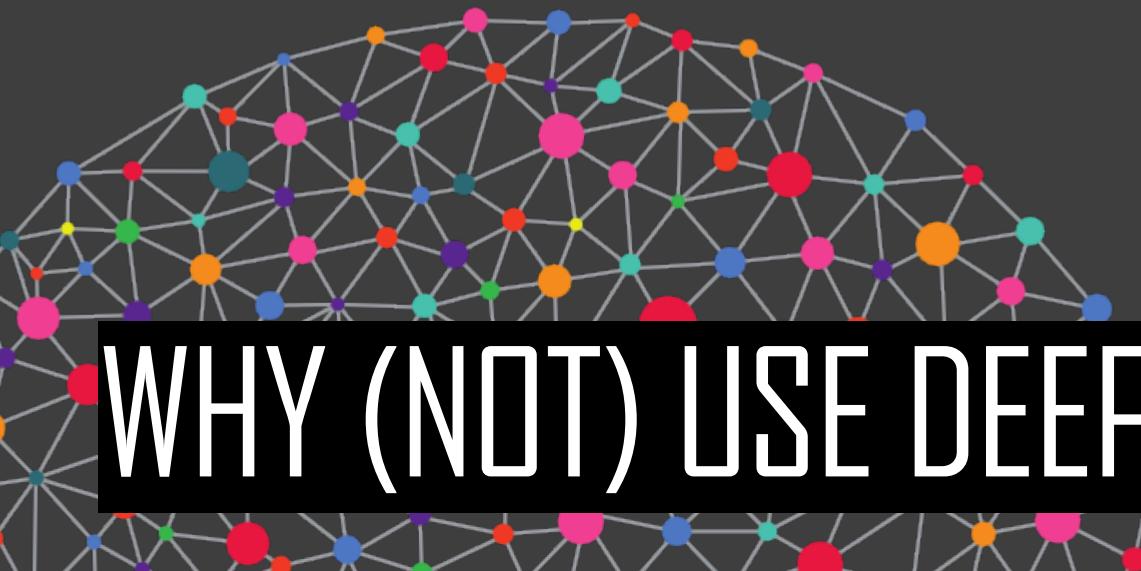
ARCHITECTURE COMPOSABILITY

1. I WANT TO MODEL TIME OR SPACE DEPENDENCIES
2. I HAVE A LOT OF DATA AND WANT A LARGE, NUANCED MODEL
3. "SOFTWARE 2.0" / NOVEL PREDICTION ARCHITECTURES
4. GENERATIVE MODELS



WHY USE DEEP LEARNING?

1. I NEED MODEL INTERPRETABILITY
2. I NEED PARSIMONY / LOW COMPUTE
3. I AM WORRIED ABOUT BIAS
4. I NEED TO UNDERSTAND MY DATA



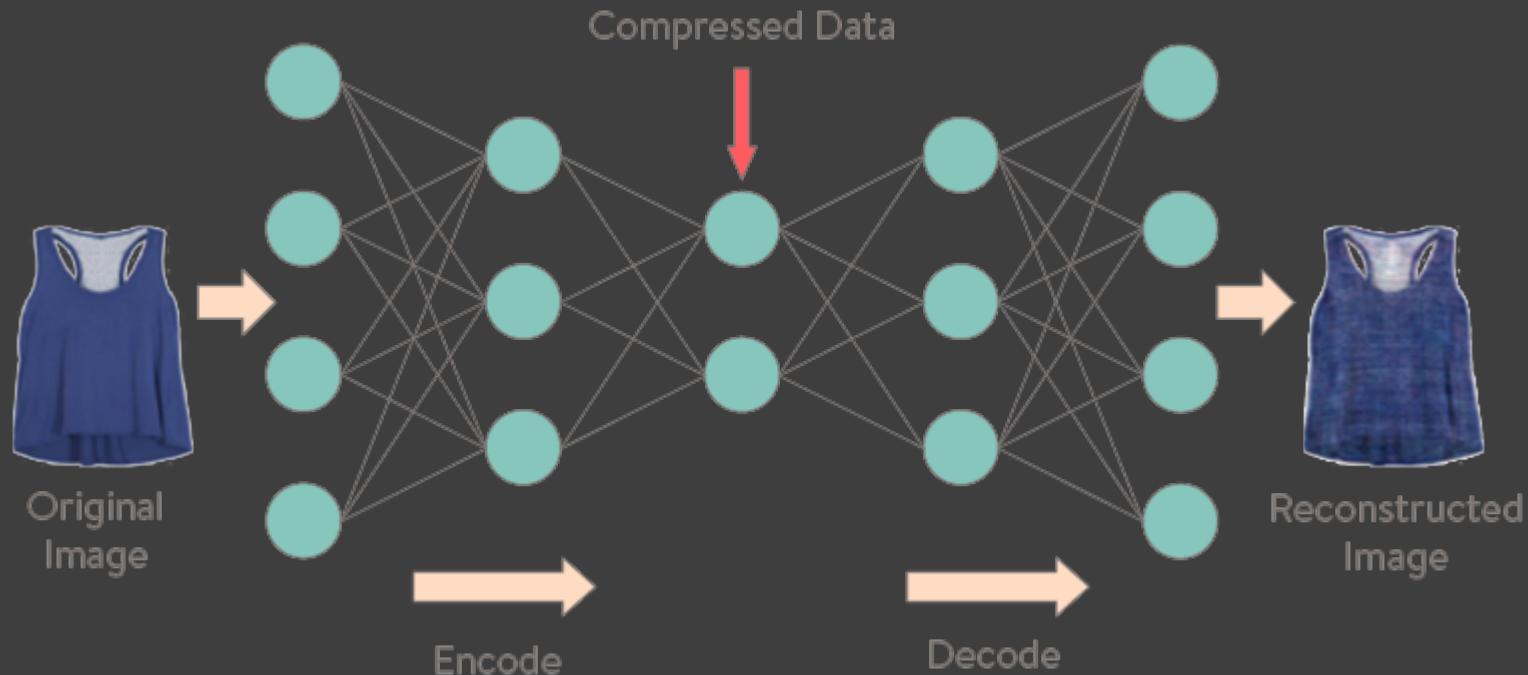
WHY (NOT) USE DEEP LEARNING?

TWO INNOVATIVE DL ARCHITECTURES



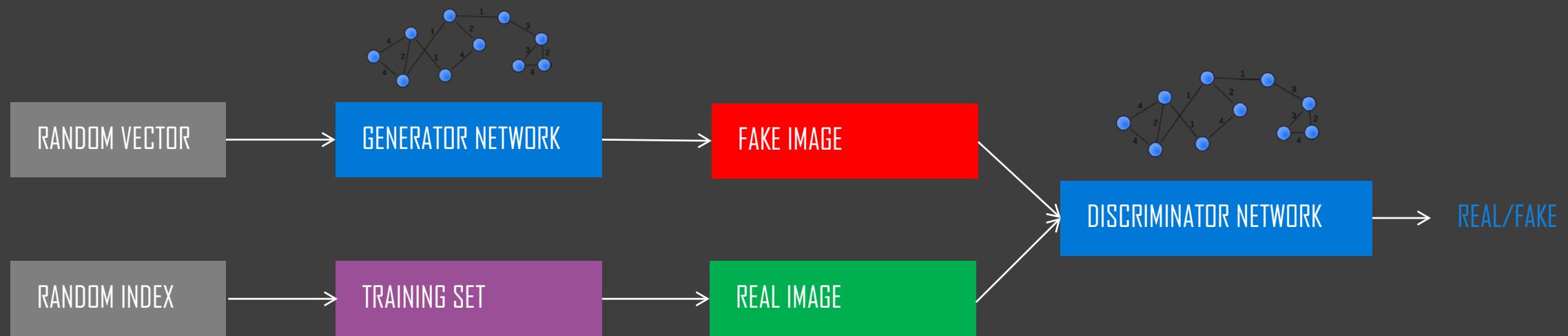


- Autoencoders are an unsupervised method to compress information into an “embedding” similar to unsupervised dimensionality reduction.



AUTOENCODERS

- An example; generative adversarial networks
- Game-theoretical approach to deep learning



Goodfellow et al. NIPS 2014



VERY ADVANCED NOVEL ARCHITECTURES (GANS)



IMAGES GENERATED WITH GANS

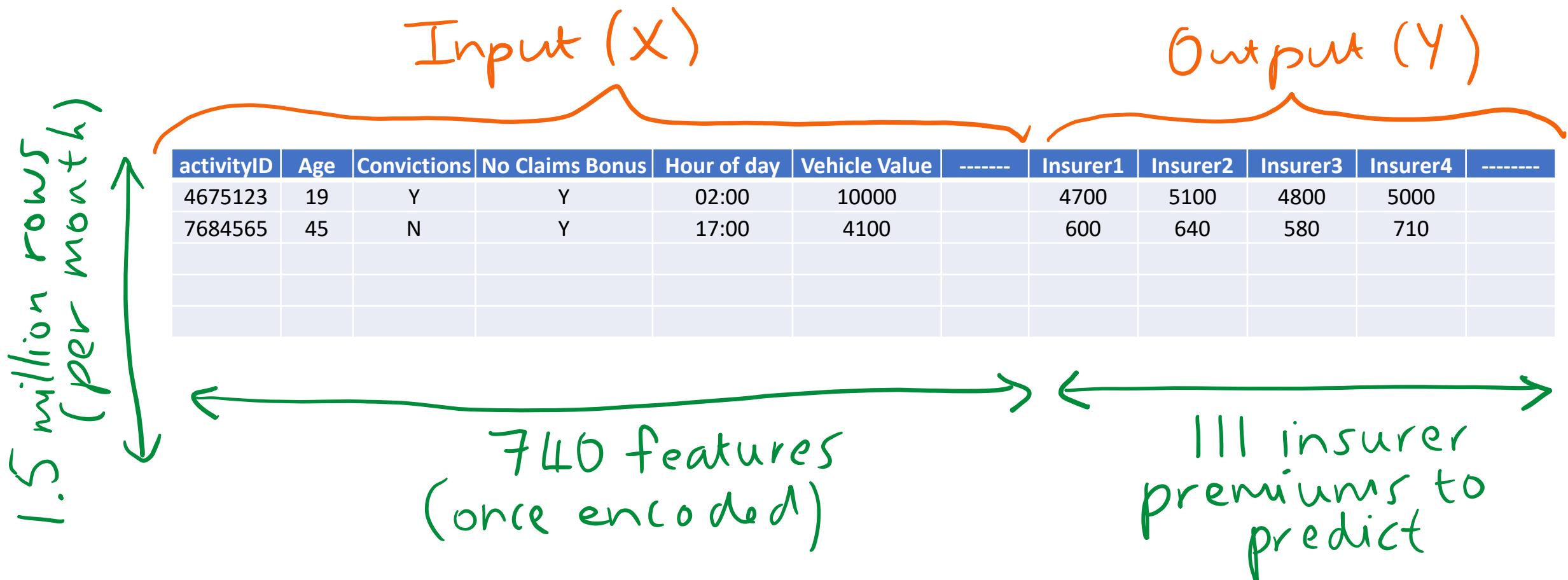
DEMO

NEURAL NETWORK PLAYGROUND



Predicting insurance premiums in a production data-science pipeline

Data

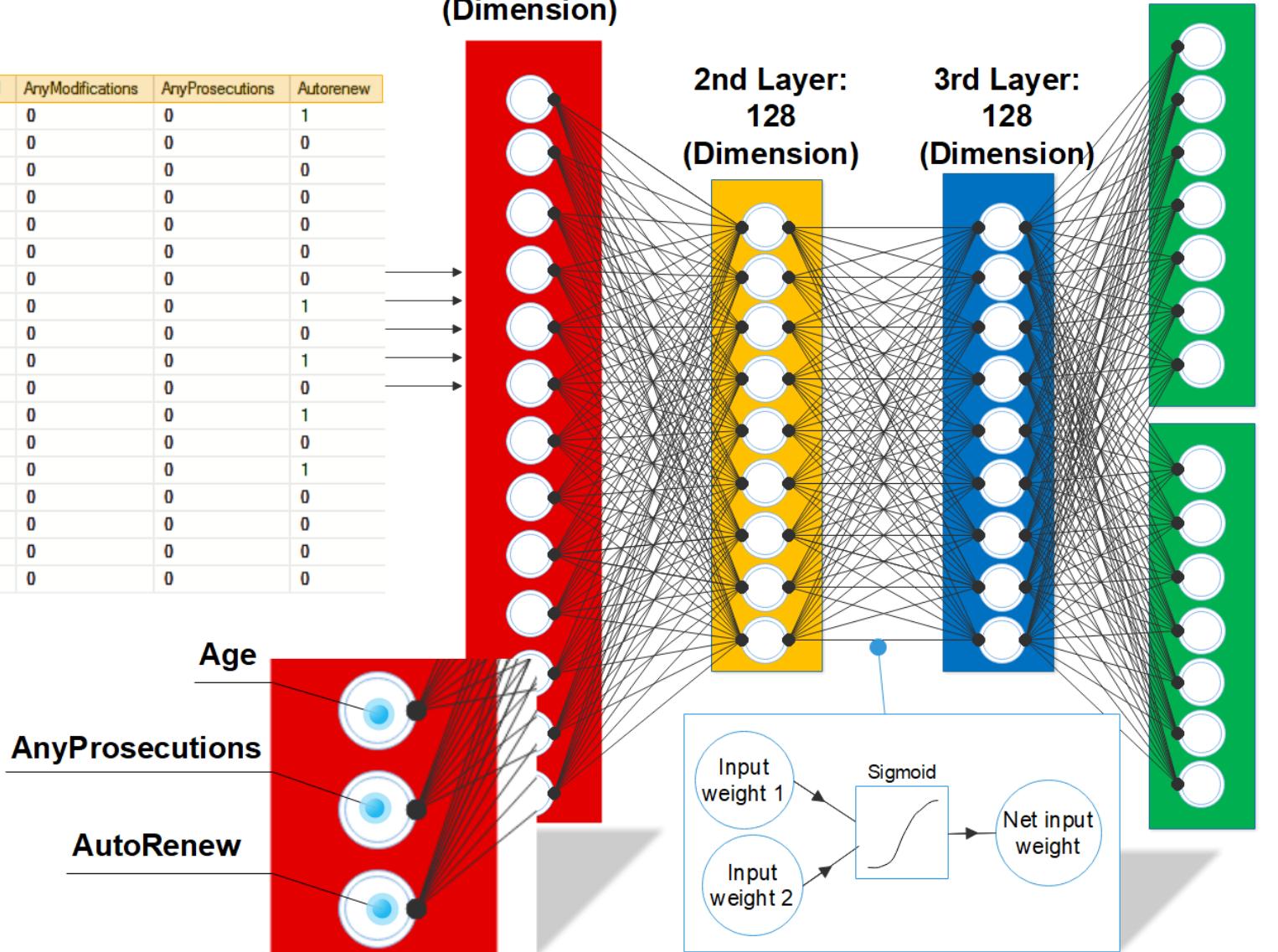


Model

	RandomId	PartitionId	Amend	AnyModifications	AnyProsecutions	Autorenew
1	9EE062C...	0	0	0	0	1
2	8BD62F1...	3	0	0	0	0
3	1FB702C...	2	0	0	0	0
4	E78B9FF...	2	0	0	0	0
5	5CCCFD4...	4	0	0	0	0
6	127C460...	3	0	0	0	0
7	0BBAB1F...	5	0	0	0	0
8	E5C4B94...	0	0	0	0	1
9	D4AE24...	8	0	0	0	0
10	ODE604A...	9	0	0	0	1
11	D1F30F6...	0	0	0	0	0
12	3177A7C...	6	0	0	0	1
13	633B8A4...	1	0	0	0	0
14	C933BF9...	10	0	0	0	1
15	107EE8C...	5	0	0	0	0
16	407A008...	8	0	0	0	0
17	348EC6E...	3	0	0	0	0
18	0EA607F...	4	0	0	0	0

Input Layer:
~700
features
(Dimension)

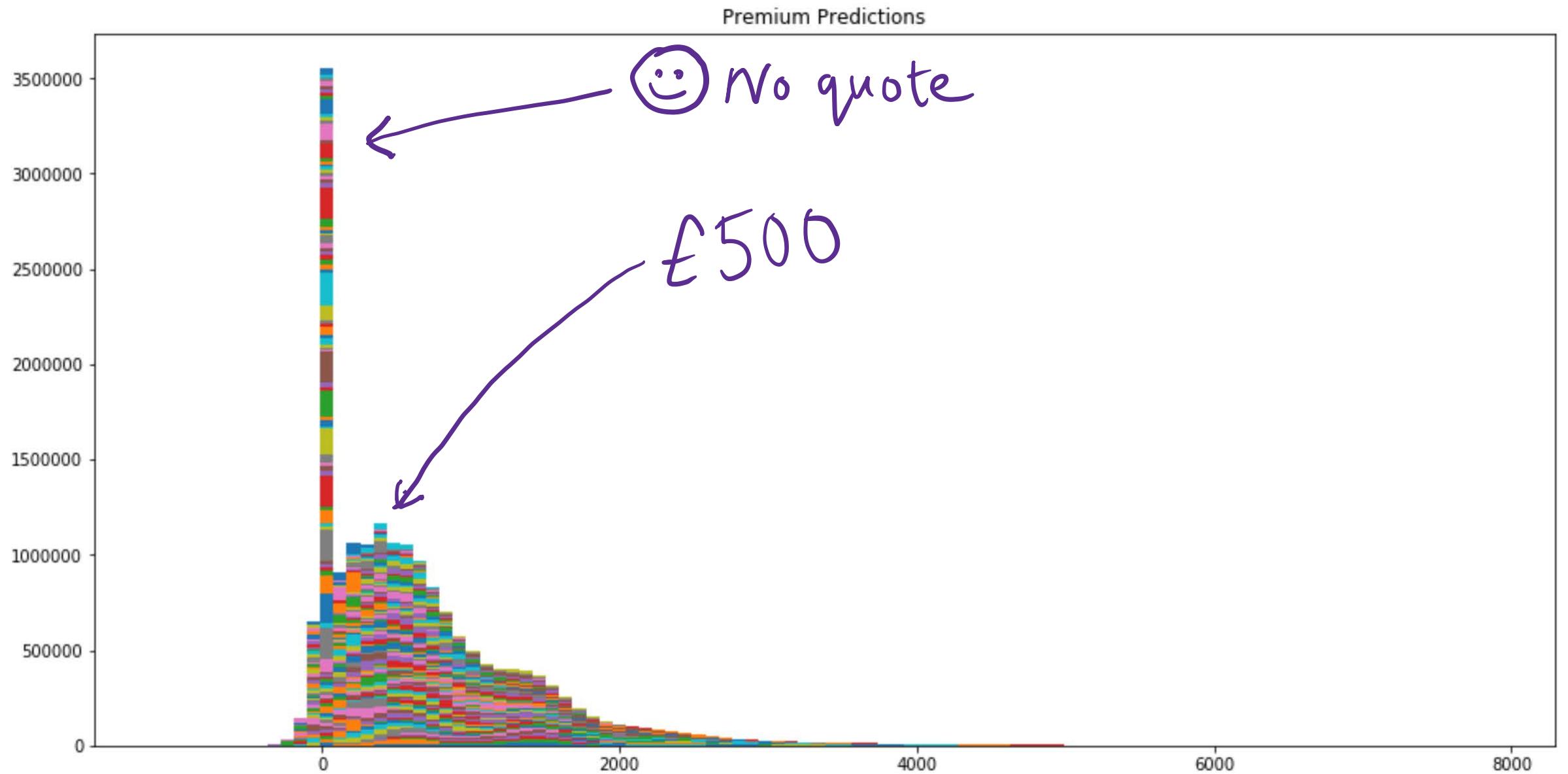
Output Layer:
110
(Dimension)



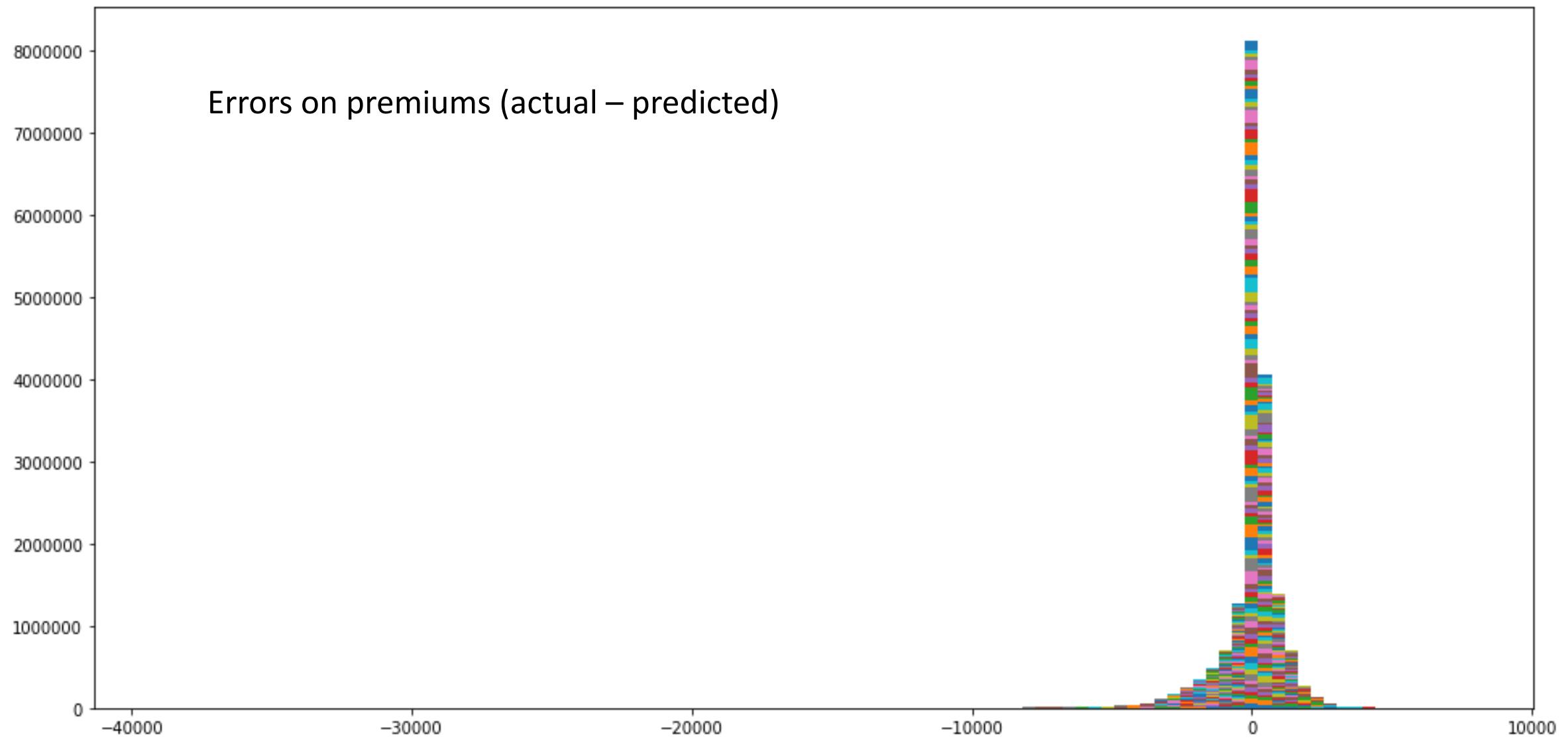
**Premium
predictions
(e.g £500)**

**Quote
predictions
(YES/NO)**

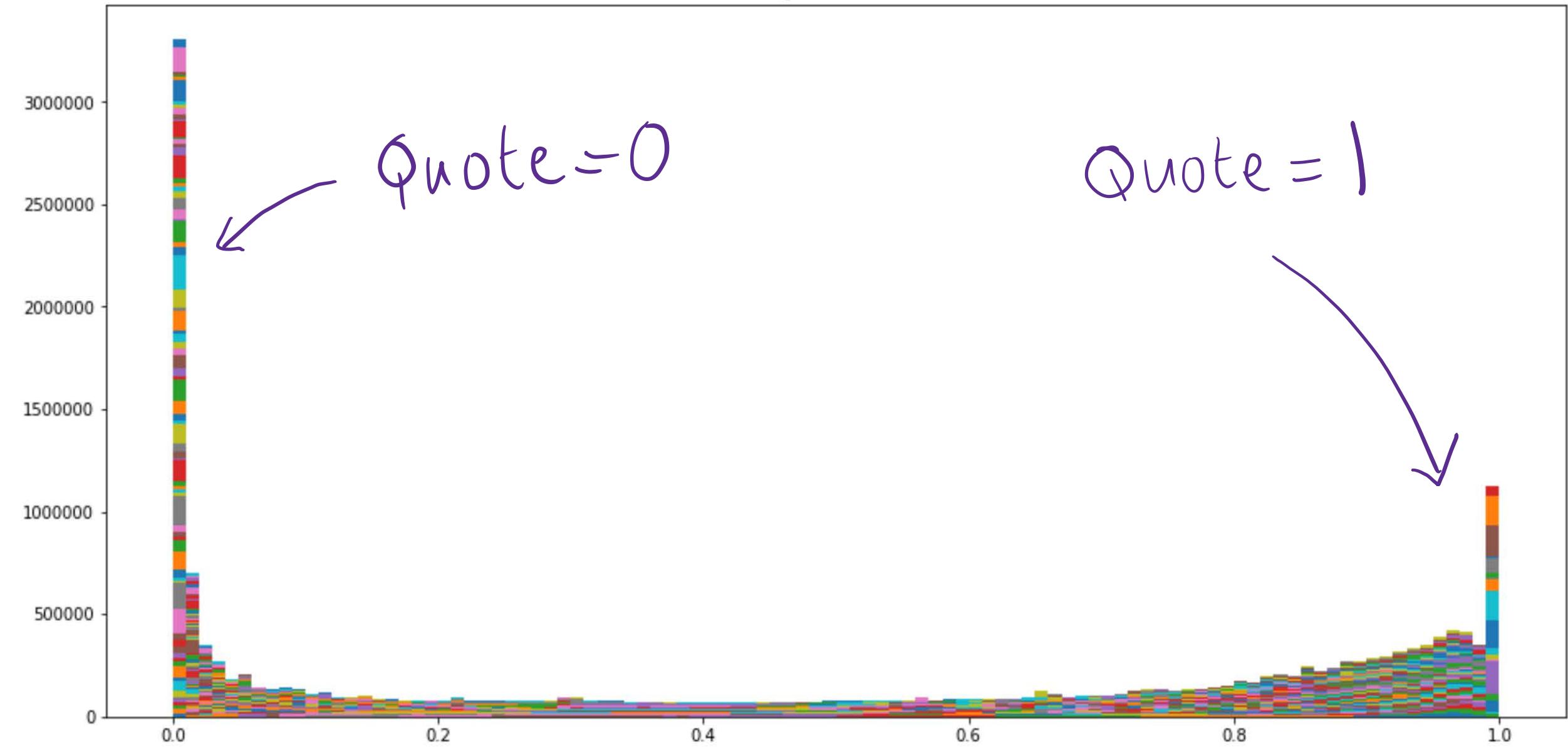
Model results on test data

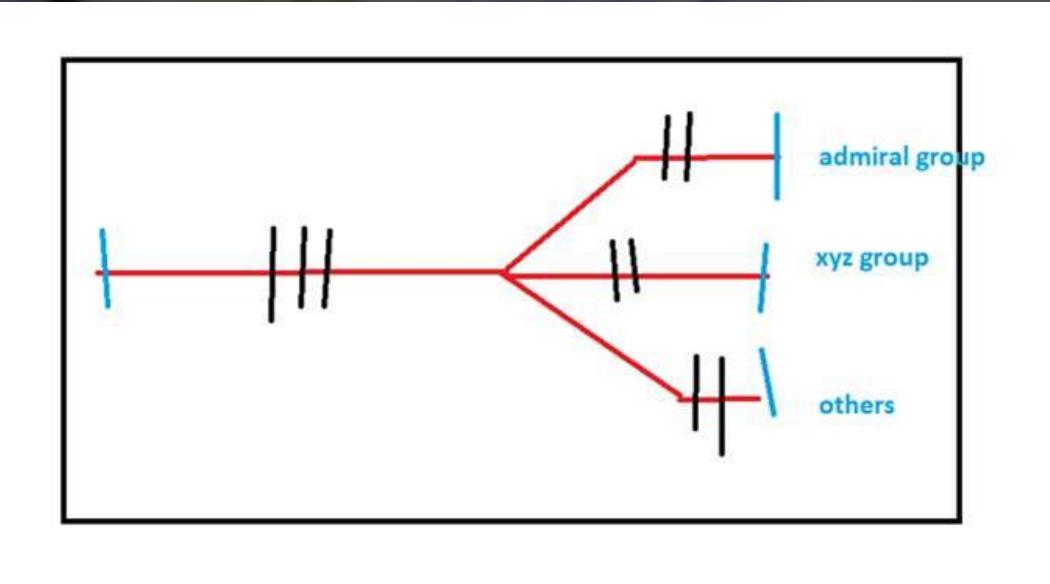
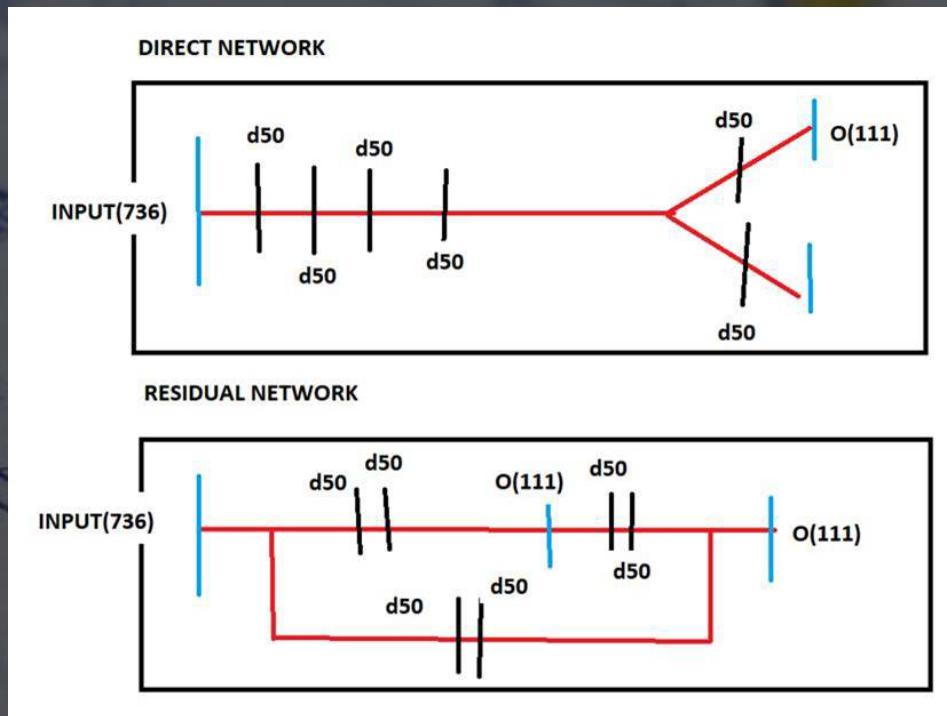


Errors on premiums (actual – predicted)

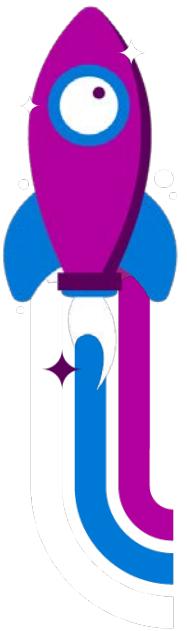


Brand Quoted Predictions





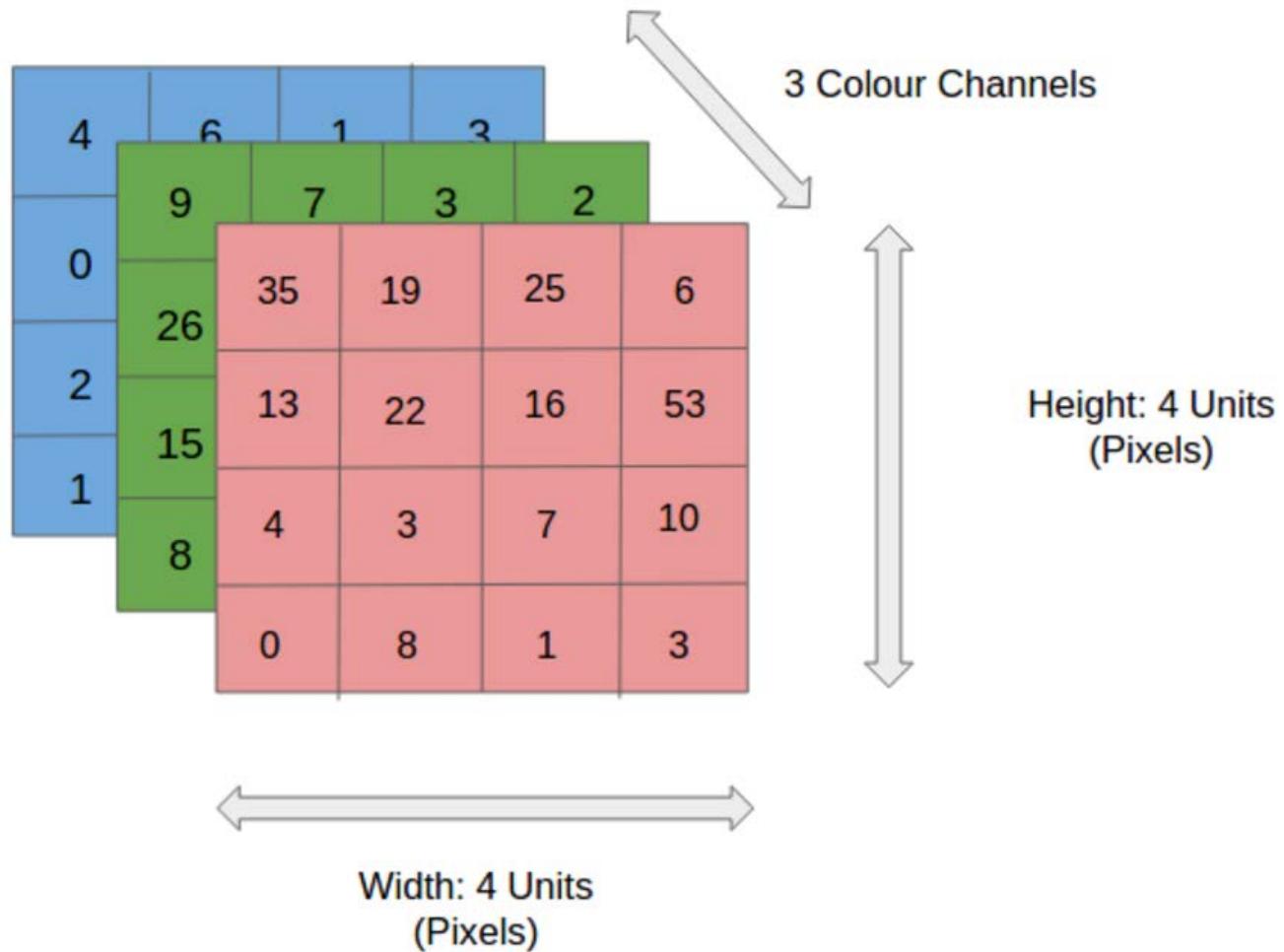
Architectures Explored



How do convolutional neural networks work?



PREPARE DATASET OF IMAGES



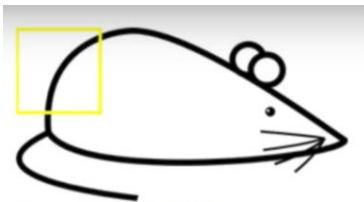
0	0	0	0	0	30	0
0	0	0	0	30	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	0	0	0	0

Pixel representation of filter



Visualization of a curve detector filter

CONVOLUTION FILTER



Visualization of the filter on the image

Visualization of the
receptive field

$$(50 \times 30) + (50 \times 30) + (50 \times 30) + (20 \times 30) + (50 \times 30) = 6600$$

0	0	0	0	0	0	30
0	0	0	0	50	50	50
0	0	0	20	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0

Pixel representation of the receptive
field

*

0	0	0	0	0	0	30	0
0	0	0	0	0	30	0	0
0	0	0	30	0	0	0	0
0	0	0	30	0	0	0	0
0	0	0	30	0	0	0	0
0	0	0	30	0	0	0	0
0	0	0	0	0	0	0	0

Pixel representation of filter

CONVOLUTION FILTER MATCH



Visualization of the filter on the image

MULTIPLY AND SUMMATION = 0

0	0	0	0	0	0	0
0	40	0	0	0	0	0
40	0	40	0	0	0	0
40	20	0	0	0	0	0
0	50	0	0	0	0	0
0	0	50	0	0	0	0
25	25	0	50	0	0	0

Pixel representation of receptive field

*

0	0	0	0	0	30	0
0	0	0	0	30	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	0	0	0	0

Pixel representation of filter

CONVOLUTION FILTER NO MATCH

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4
2	4	3

Convolved Feature

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4
2	4	3
2	3	4

Convolved Feature



CONVOLUTION

DEMO

Operation	Kernel	Image result
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	

4	6	1	3
0	8	12	9
2	3	16	100
1	46	74	27



8	12
46	100

(i)

35	19	25	6
13	22	16	63
4	3	7	10
9	8	1	3



35	63
9	10

(iii)

9	7	3	2
26	37	14	1
15	29	16	0
8	6	54	2



37	14
29	54

(ii)

35	19	25	6
13	22	16	63
4	3	7	10
9	8	1	3



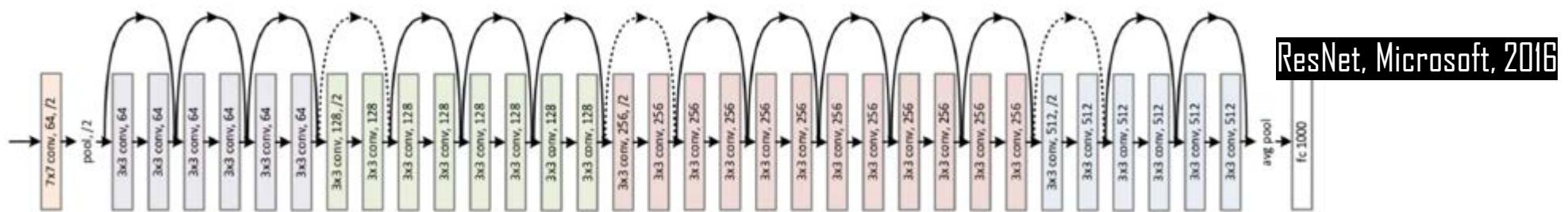
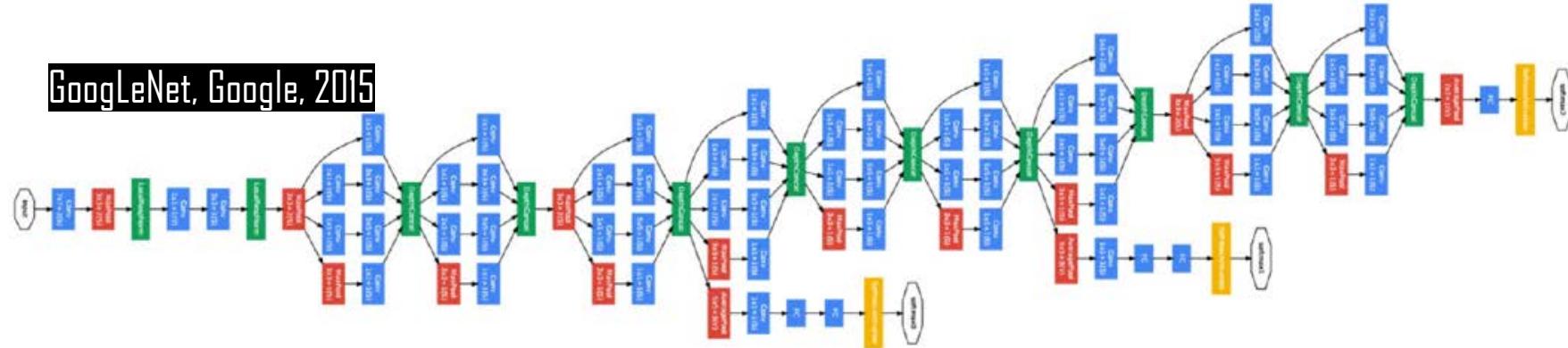
35	25	63
22	22	63
9	8	10

(iv)

POOLING

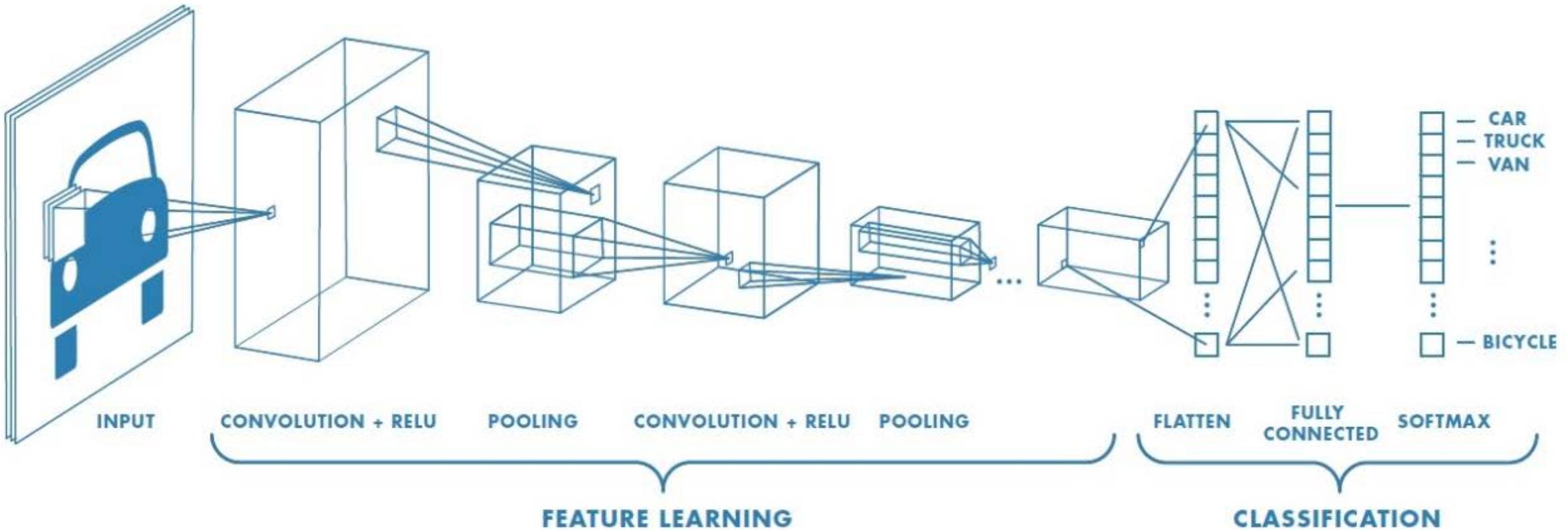


GoogLeNet, Google, 2015

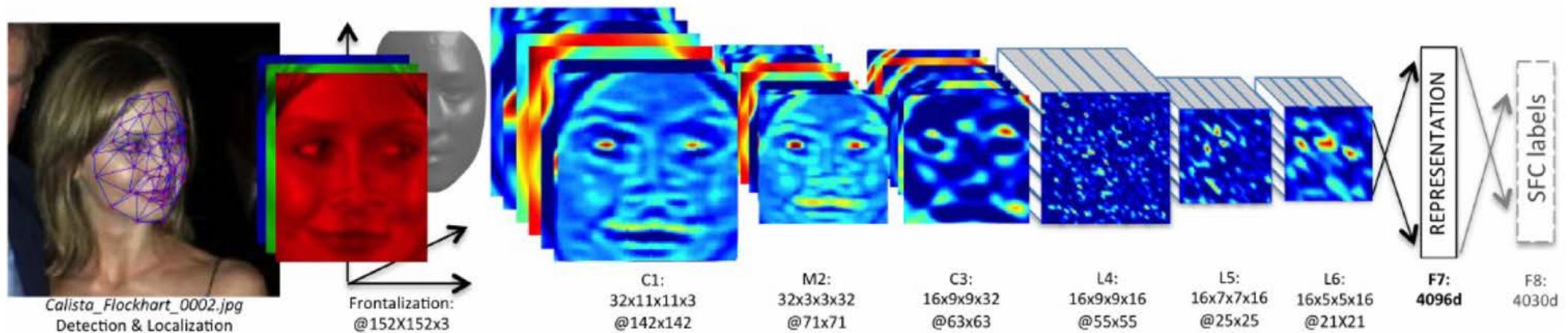


WINNING ARCHITECTURES

IMAGE CLASSIFICATION EXAMPLE



DEEP FACE

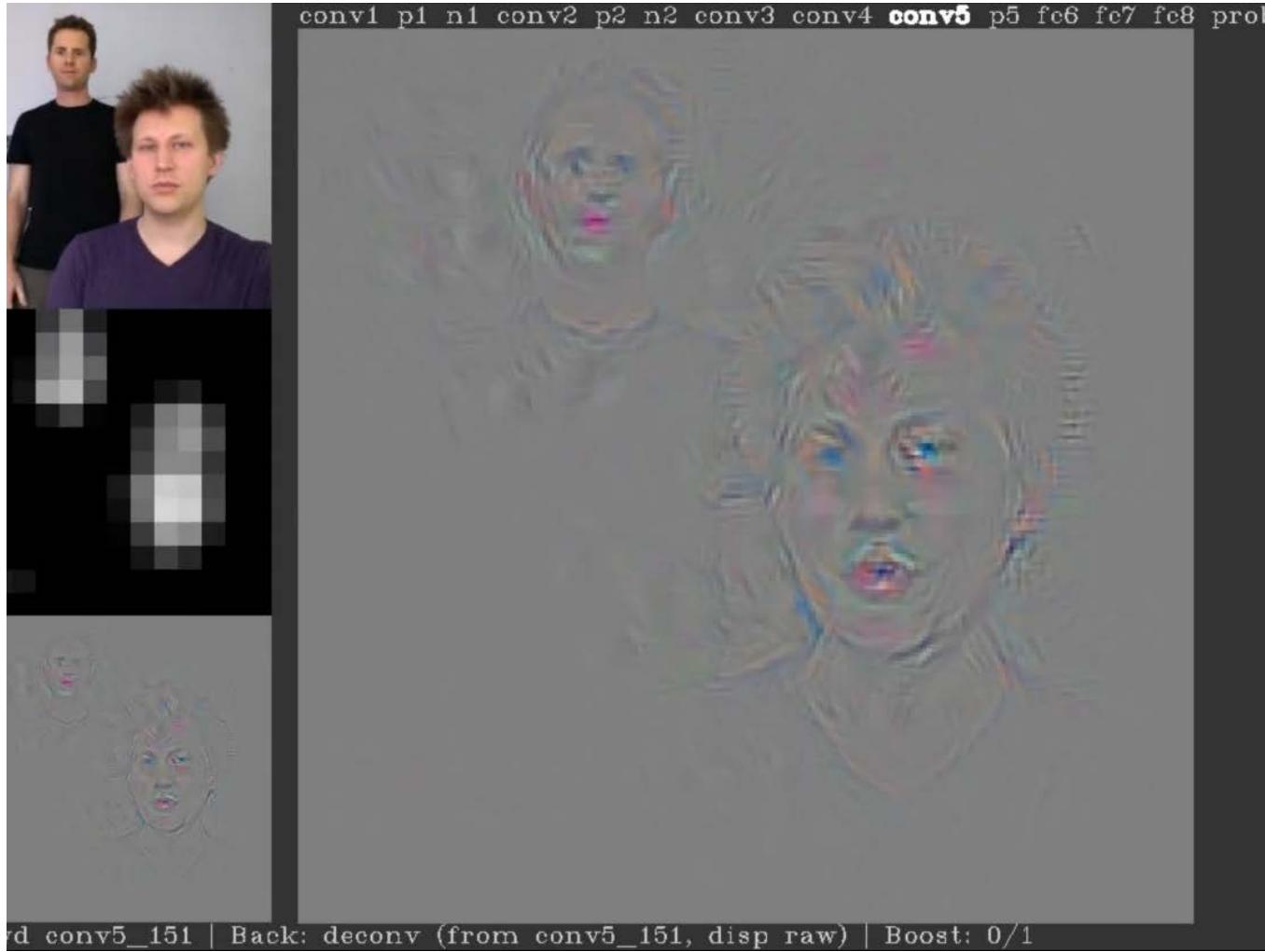


- Alignment
- CNN
- Metric Learning

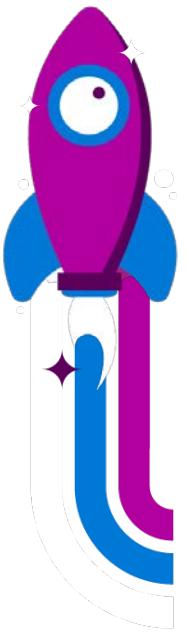
VISUALISING THE FILTERS



DEEP VISUALISATION TOOLBOX



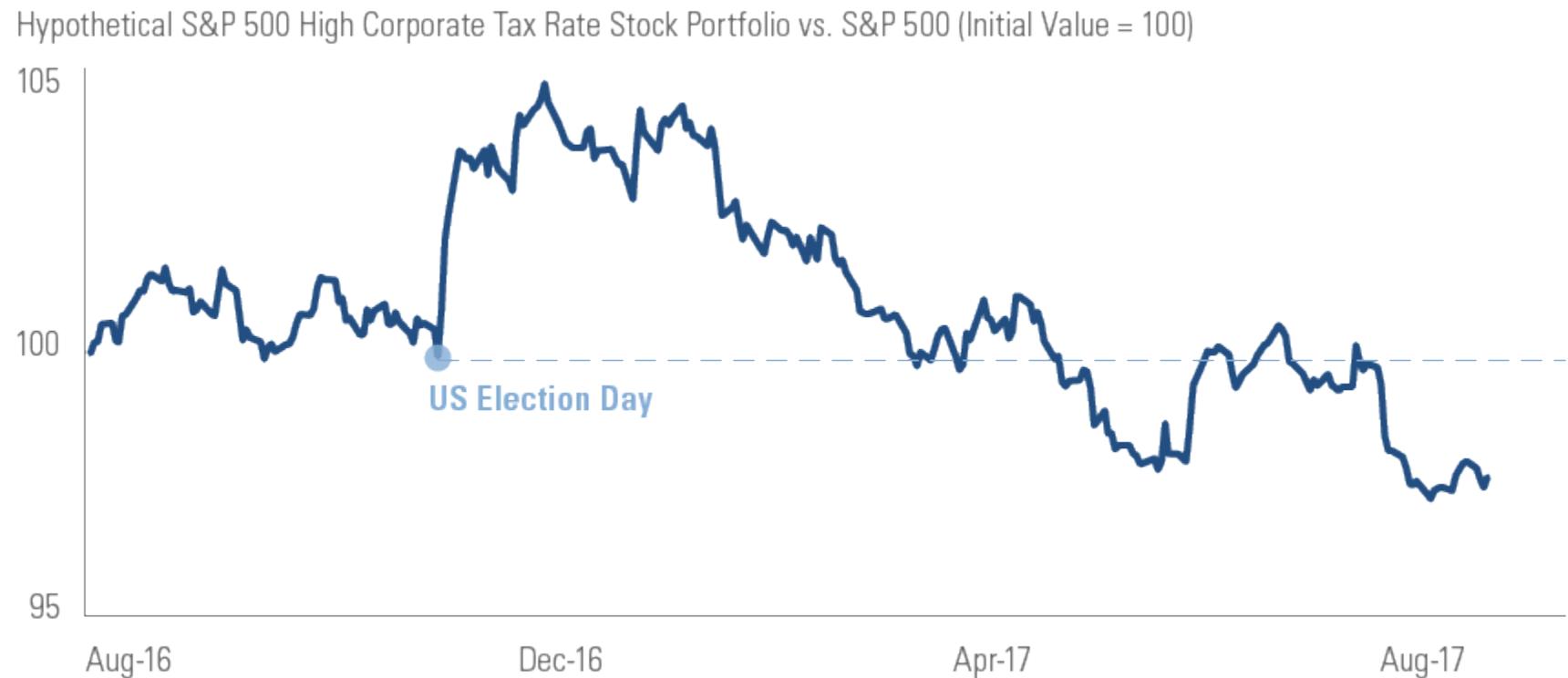
<https://www.youtube.com/watch?v=AgkfIQ4lGaM>



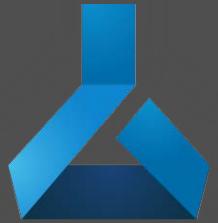
Demo of CNNs in Data Bricks/Keras



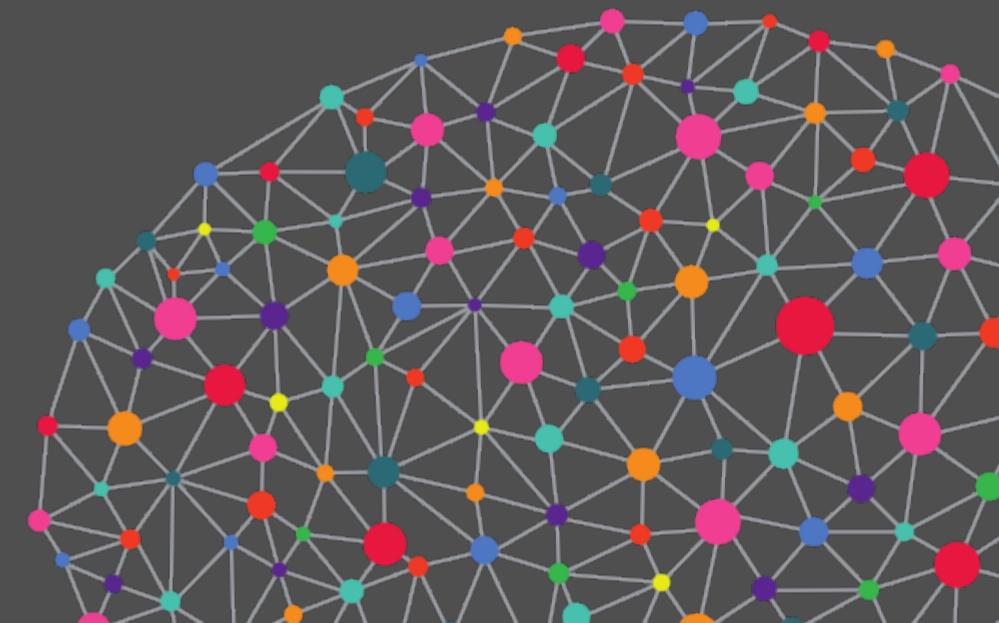
WHAT IS A SEQUENCE?



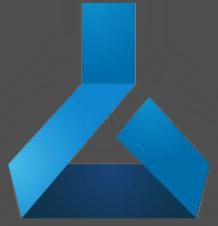
WHAT IS SEQUENCE PROCESSING



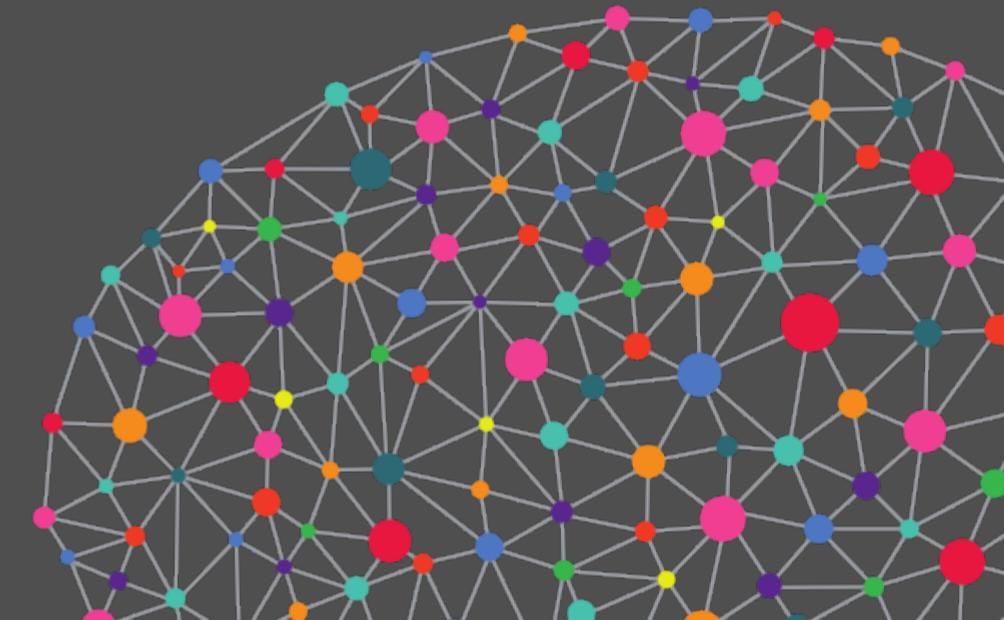
- RNNs
 - Timeseries classification
 - Anomaly detection in timeseries
 - Entity recognition
 - Revenue forecasting
 - Question + Answer
- 1d Convnets
 - Spelling correction
 - Document classification
 - Machine translation



WHAT IS SEQUENCE PROCESSING



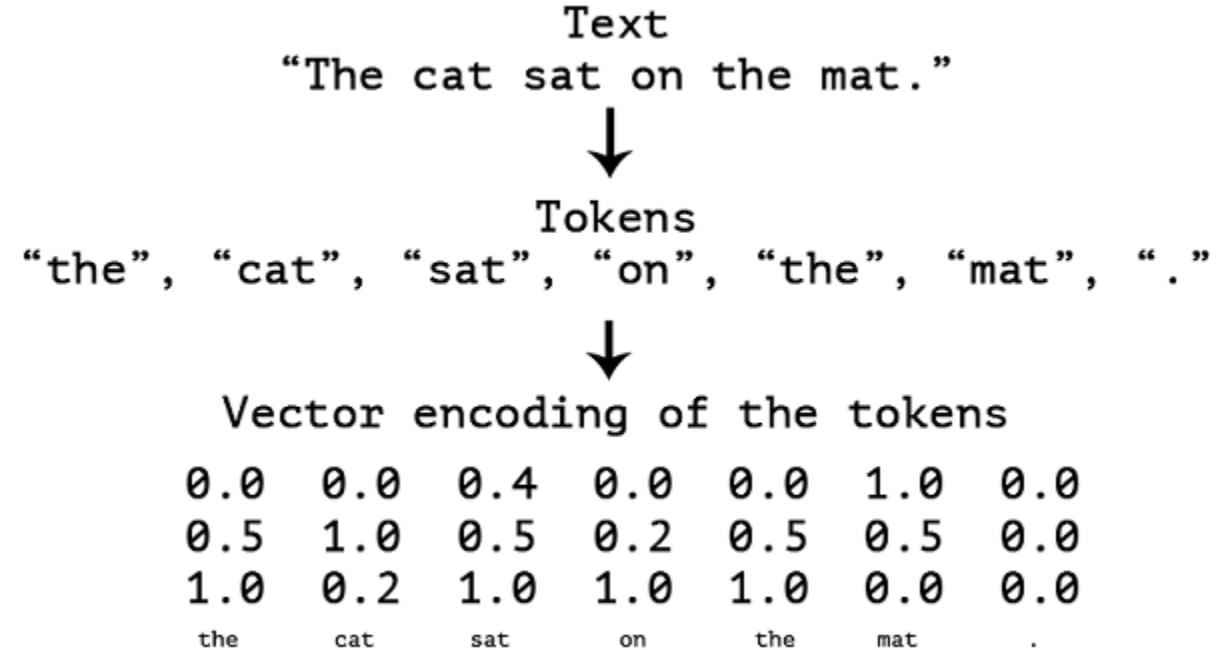
- RNNs
 - When global order matters
- 1d Convnets
 - Speed
 - Local temporal dependencies
- You can stack them!



TOKENIZATION



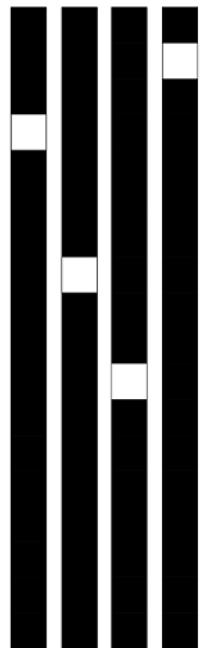
- Words
- Characters
- N-grams



N-Grams example

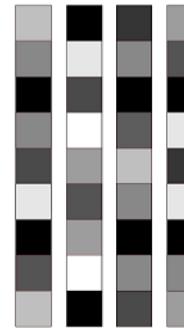
```
{"The", "The cat", "cat", "cat sat", "The cat sat",
 "sat", "sat on", "on", "cat sat on", "on the", "the",
 "sat on the", "the mat", "mat", "on the mat"}
```

WORD VECTORS VS WORD EMBEDDINGS?



One-hot word vectors:

- Sparse
- High-dimensional
- Hard-coded



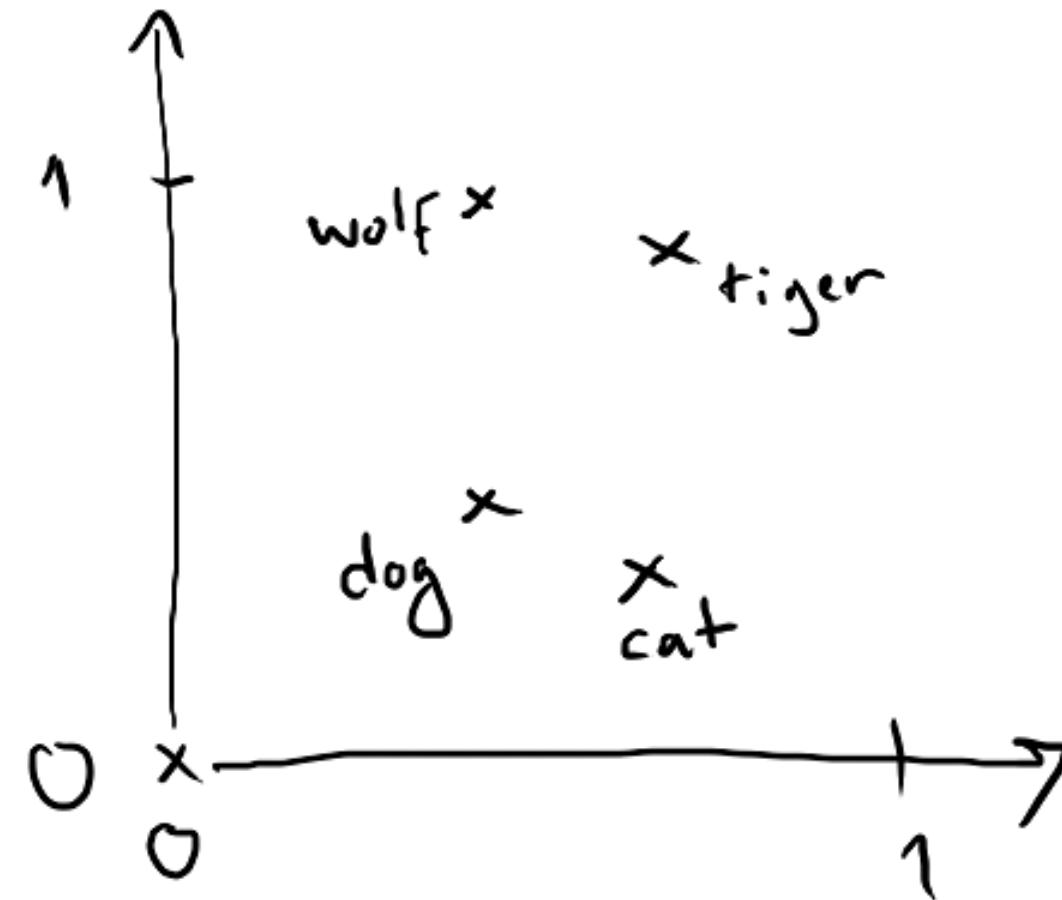
Word embeddings:

- Dense
- Lower-dimensional
- Learned from data

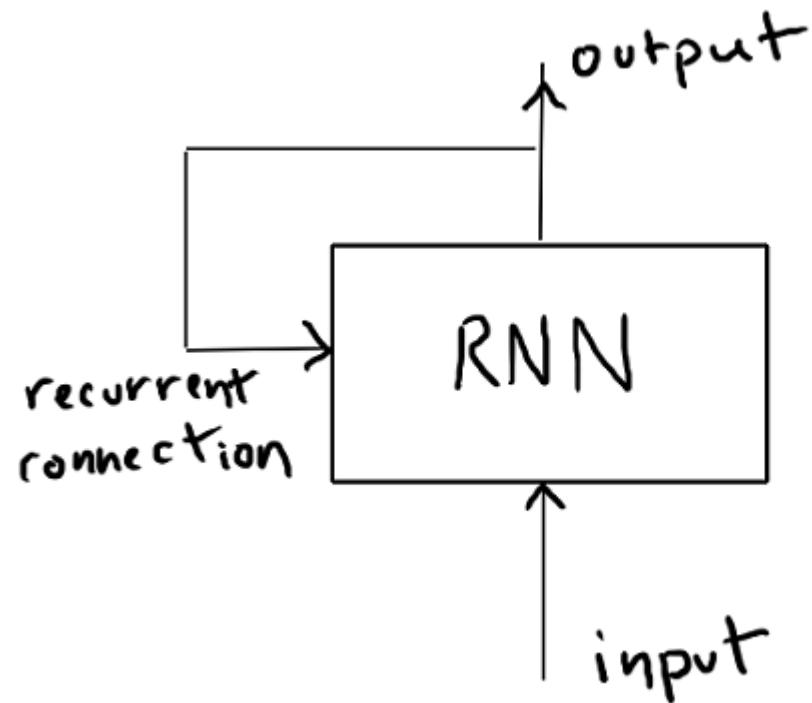
WORD EMBEDDINGS



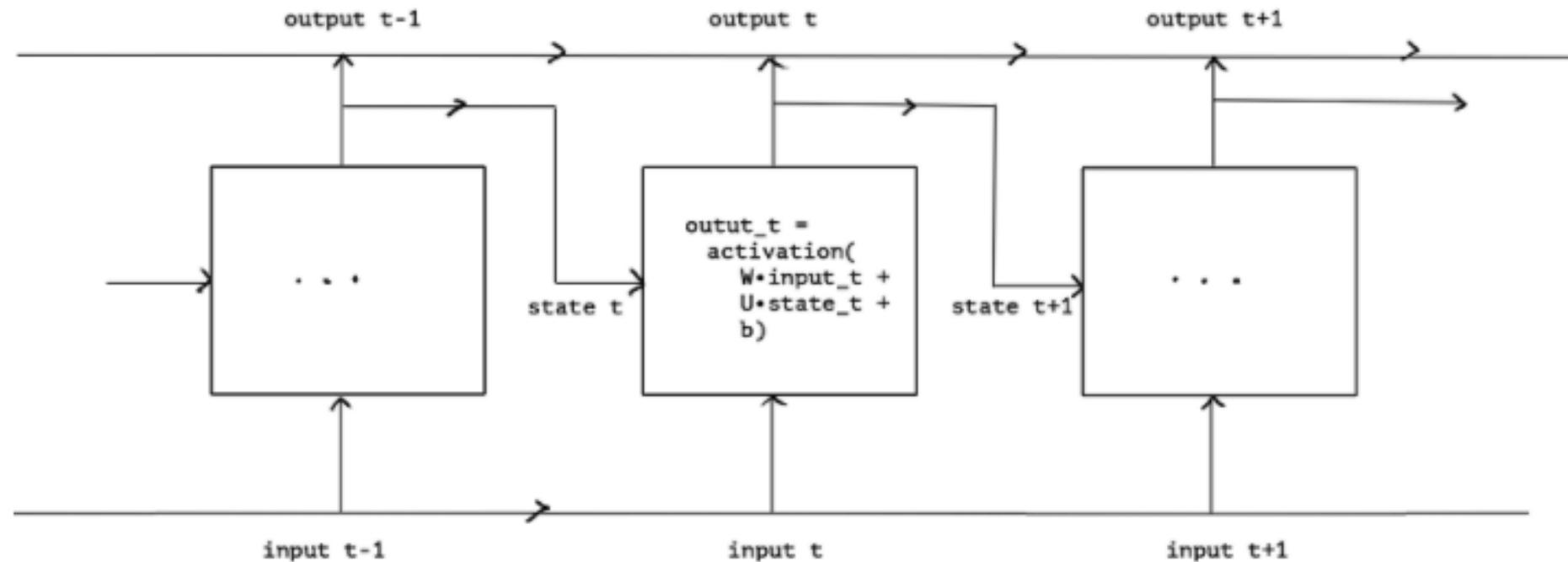
- Word2Vec
- GloVe



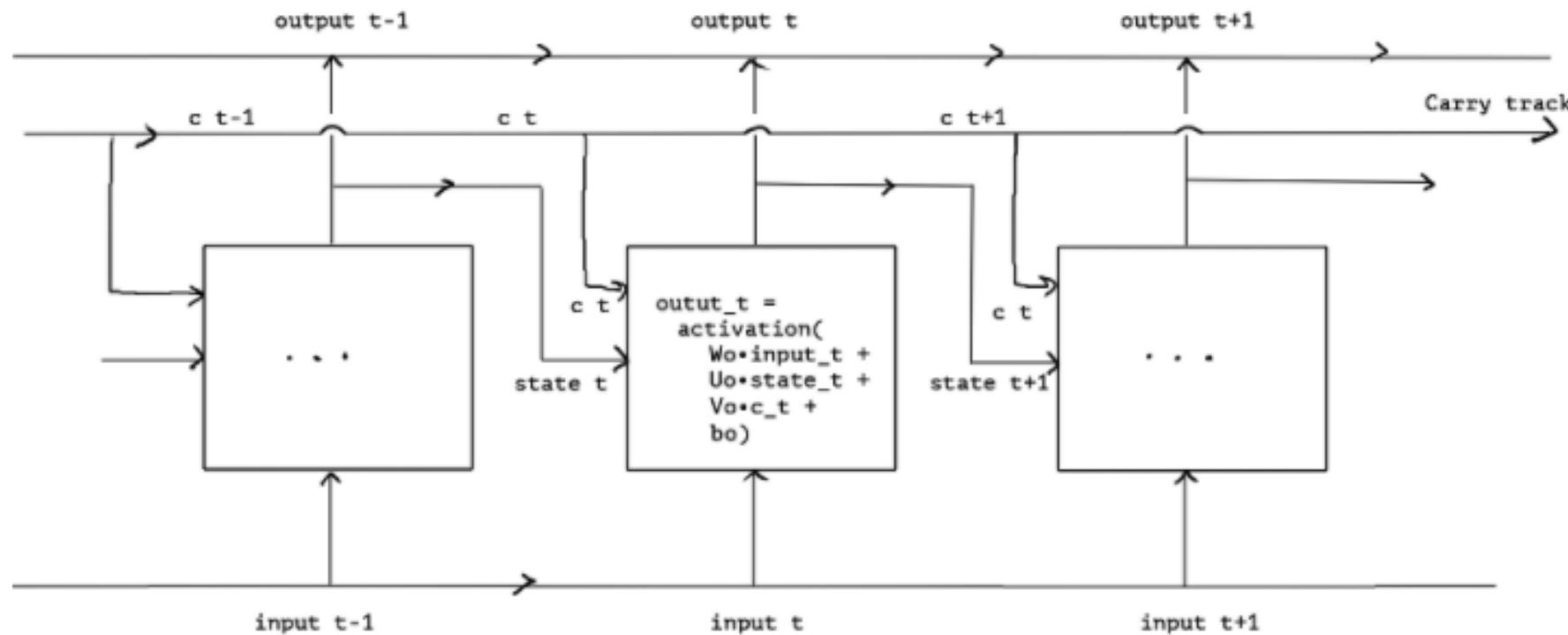
RECURRENT NEURAL NETWORKS



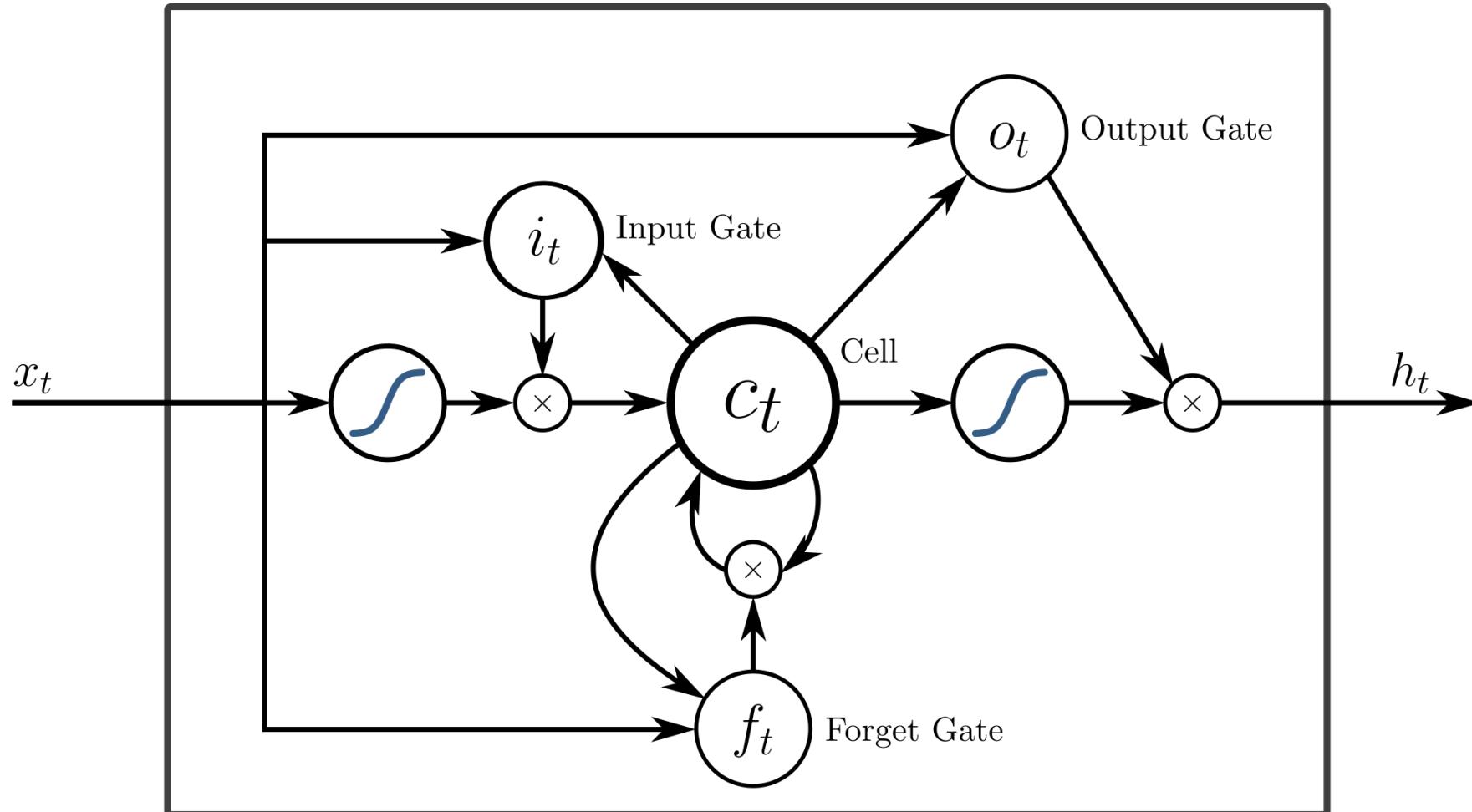
RECURRENT NEURAL NETWORKS



LSTMS(1)

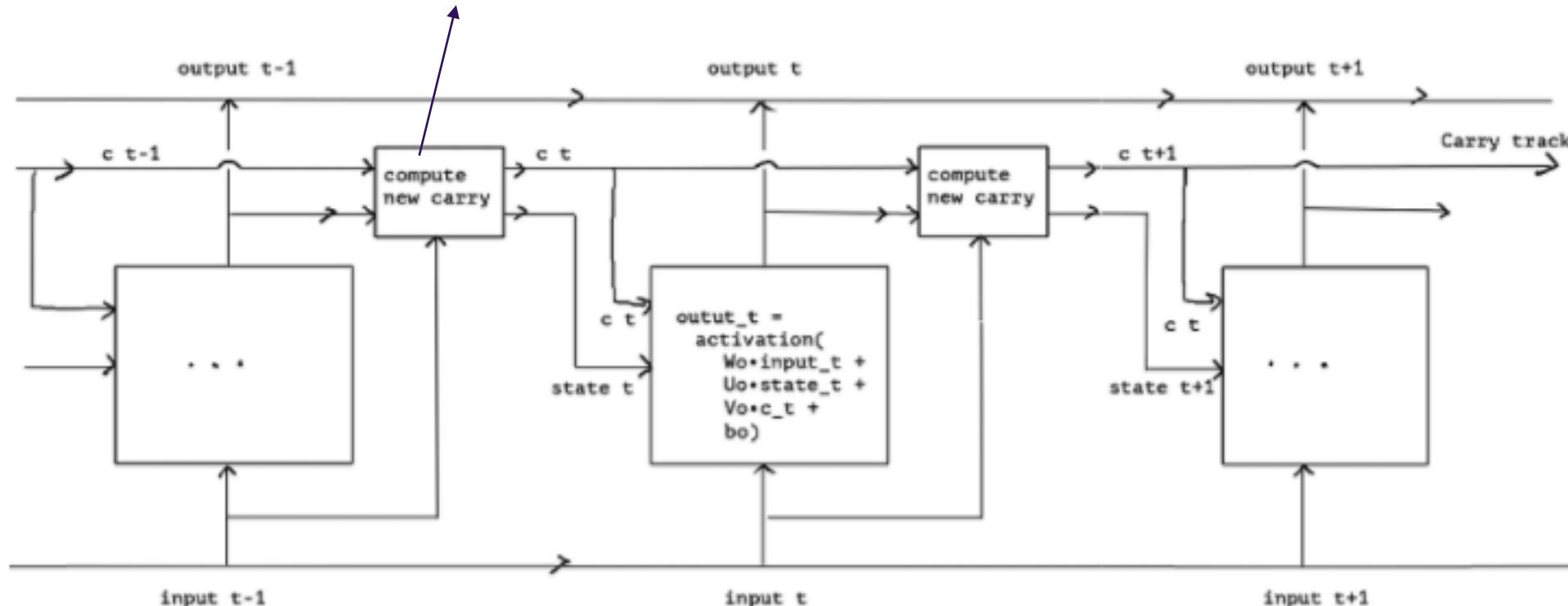


LSTMS(2)



LSTMS(3)

$$c_{t+1} = i_t * k_t + c_t * f_t$$

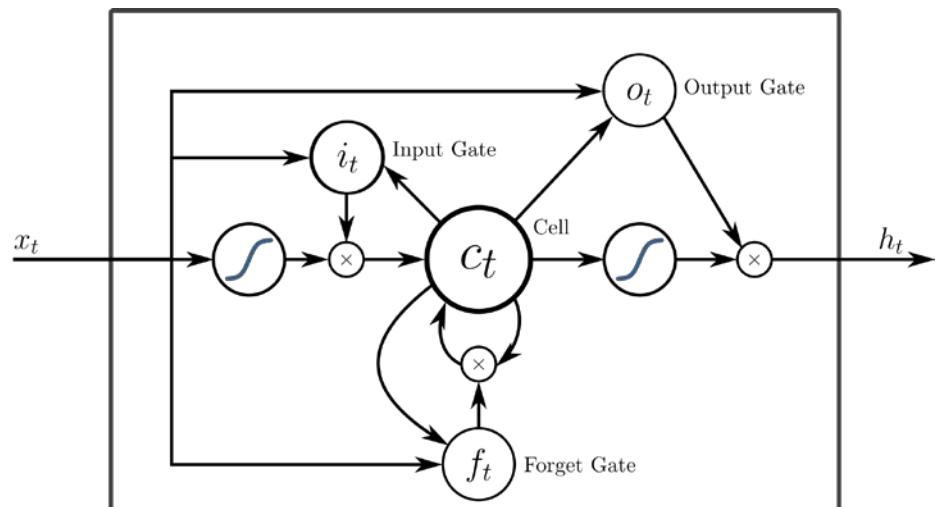


```

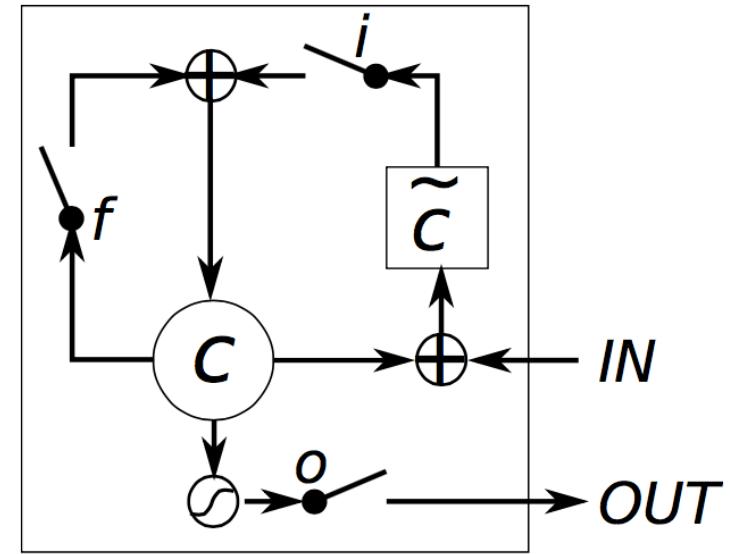
i_t = activation(dot(state_t, Ui) + dot(input_t, Wi) + bi)
f_t = activation(dot(state_t, Uf) + dot(input_t, Wf) + bf)
k_t = activation(dot(state_t, Uk) + dot(input_t, Wk) + bk)

```

LSTM VS GRU

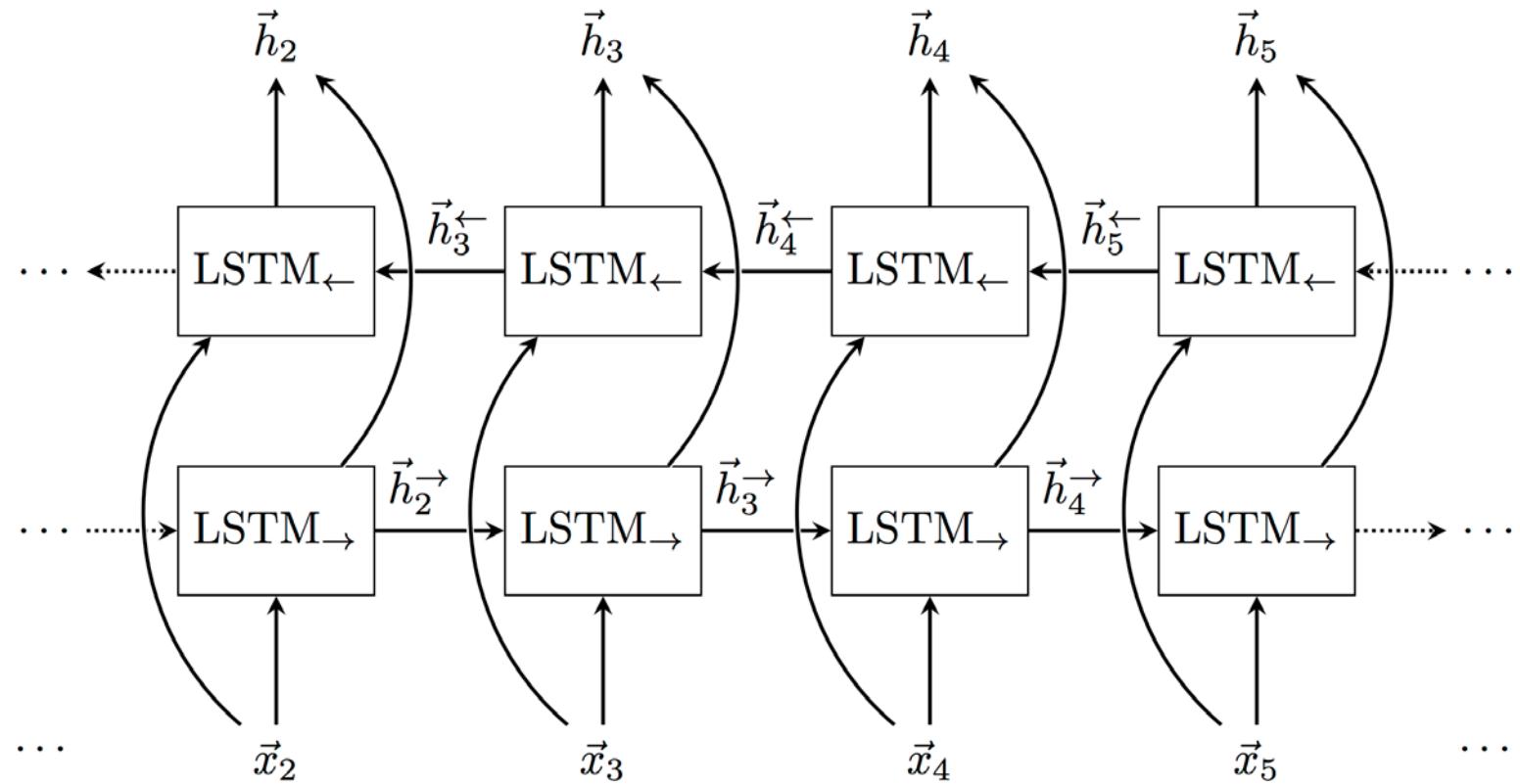


LSTM

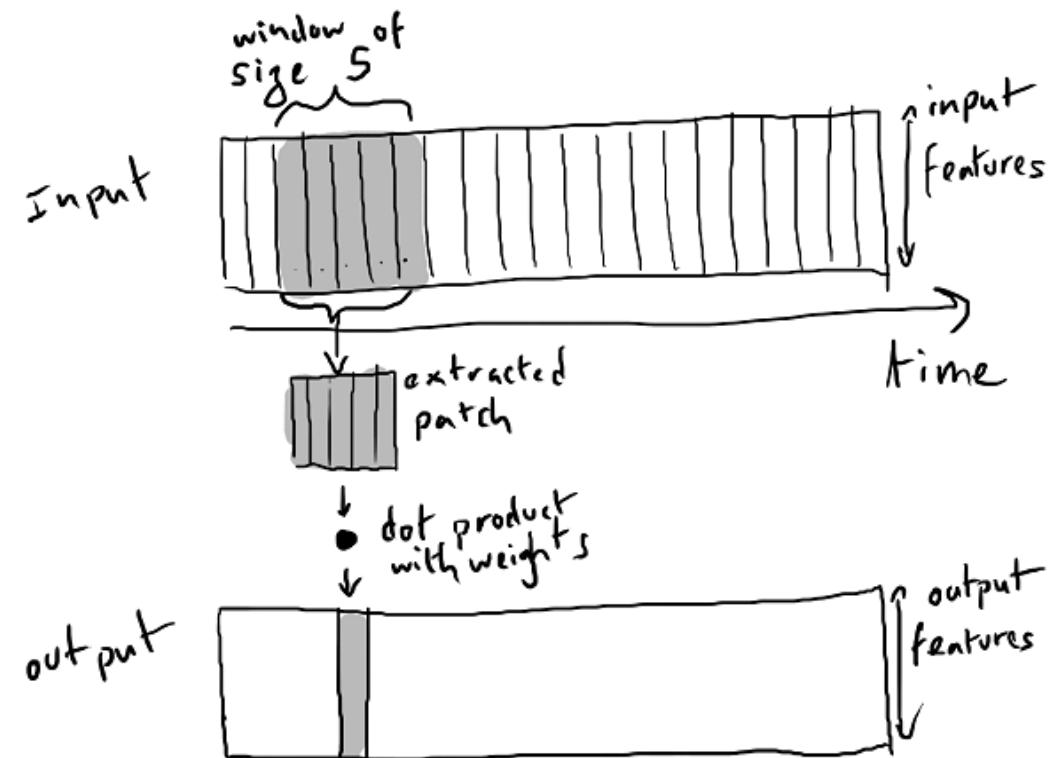


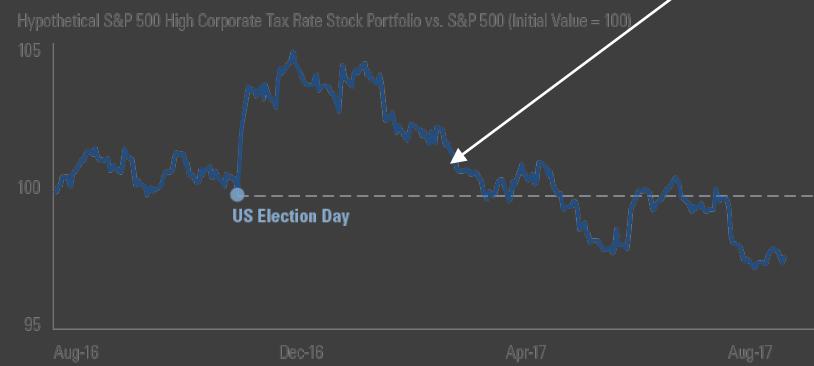
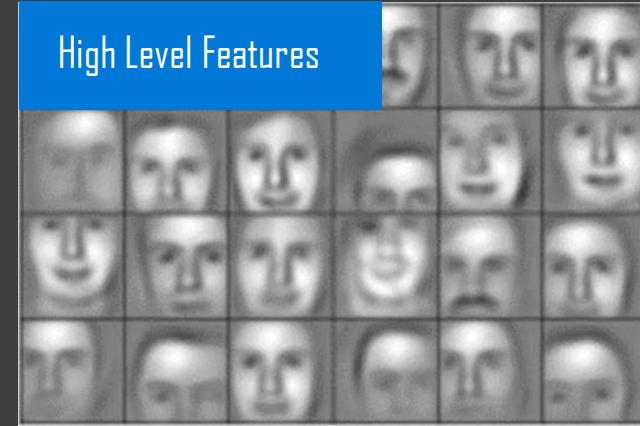
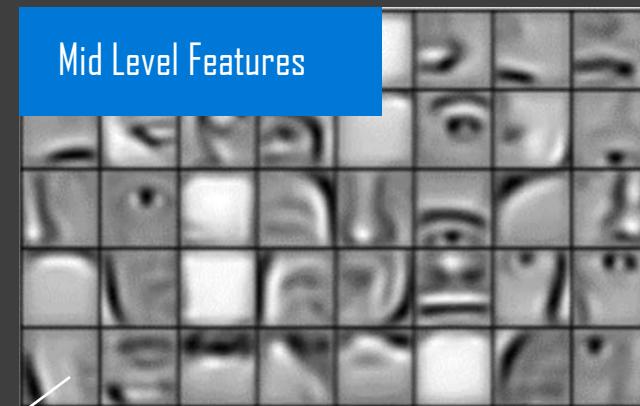
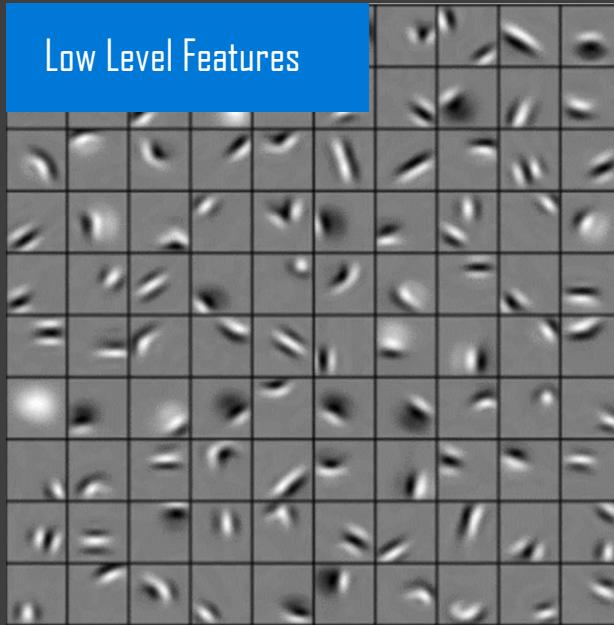
GRU

BI-DIRECTIONAL LSTMS

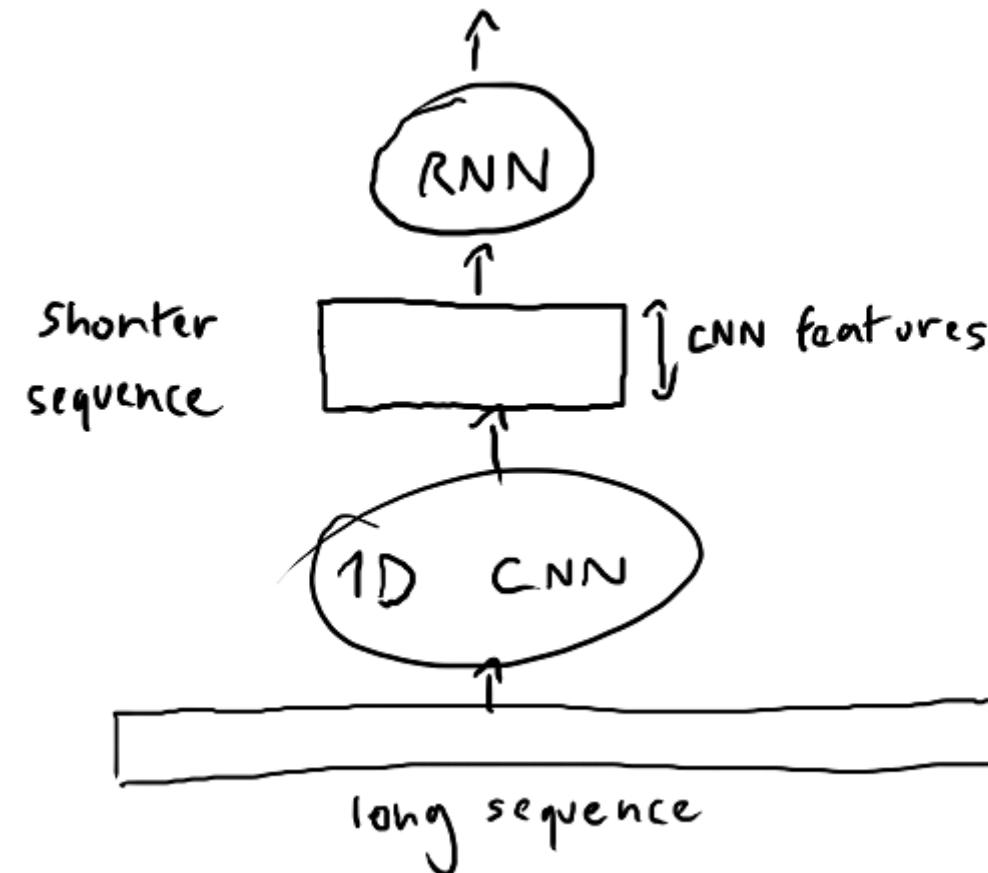


1D-CNNs

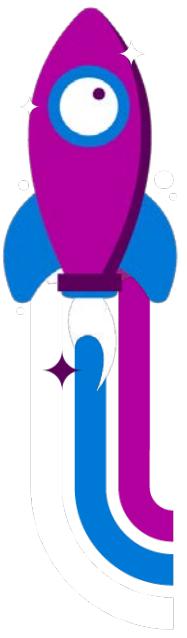




2DCNNs - SAME CONCEPT



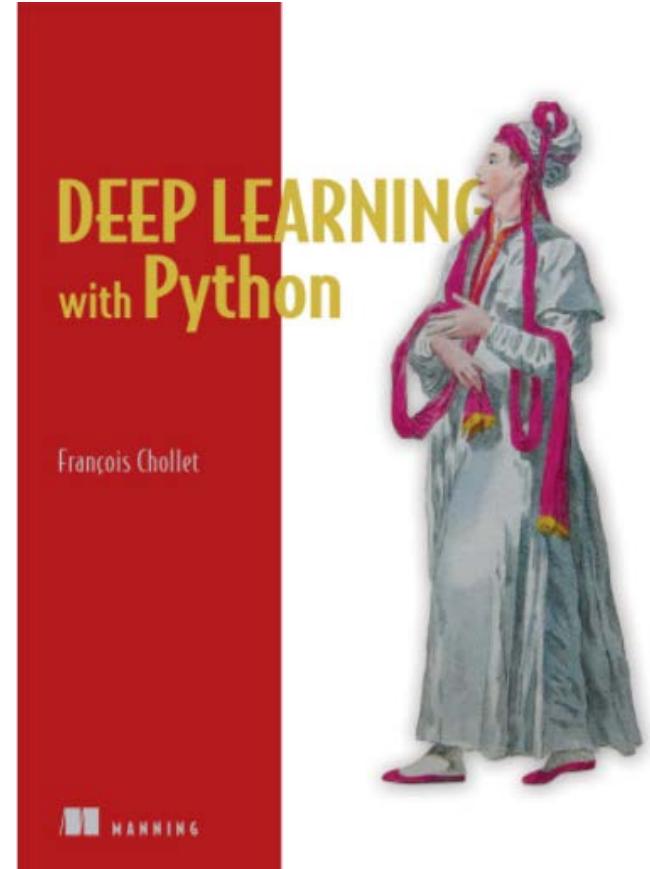
STACKING IS COOL



Demo of RNNs in Data Bricks/Keras



Francois Chollet is a living legend, please buy his book



<https://www.manning.com/books/deep-learning-with-python>



Tim Scarfe

@ecsquendor

youtube.com/machinelearningatmicrosoft

