

Computer Vision

Imagenet, Coco and google open images datasets are 3 most popular image datasets for computer vision. These datasets provides millions of hand annotated images with classification labels, bounding boxes for object detections and image segmentation masks. Data collection is the most difficult part in any supervised machine learning problem and availability of such datasets makes training CNNs very easy.

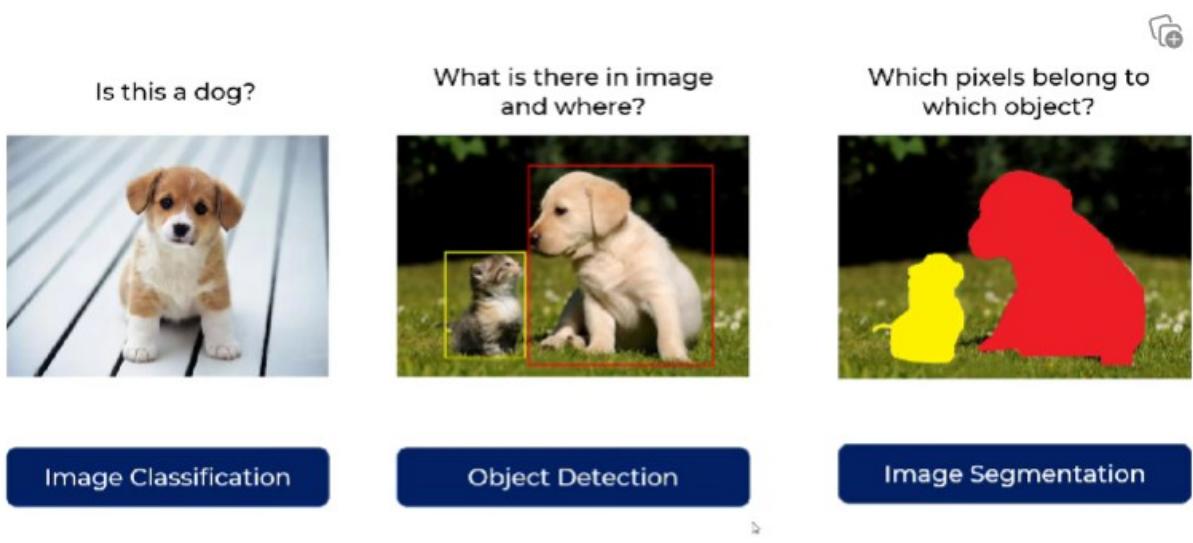


Image Classification

Process of assigning labels is called **annotation**



Object Detection

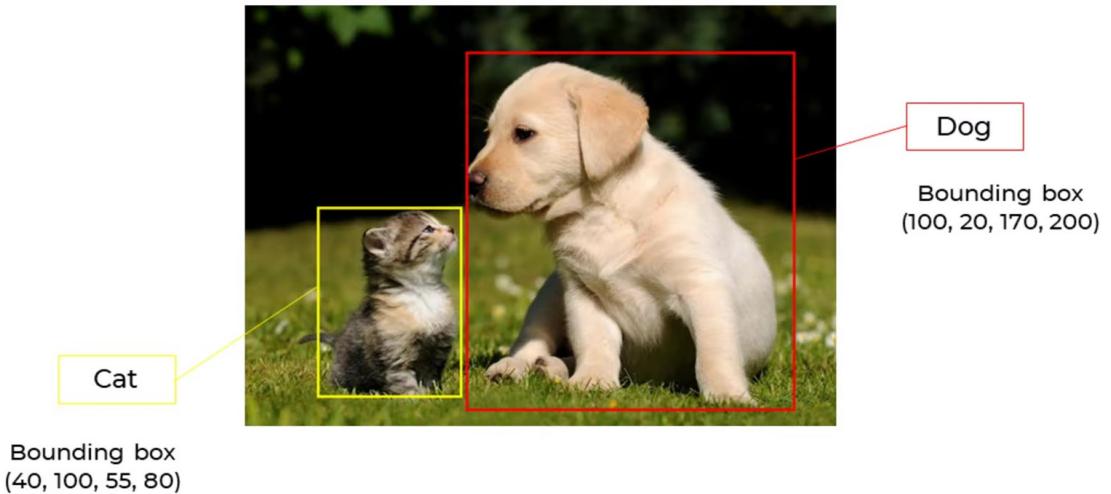
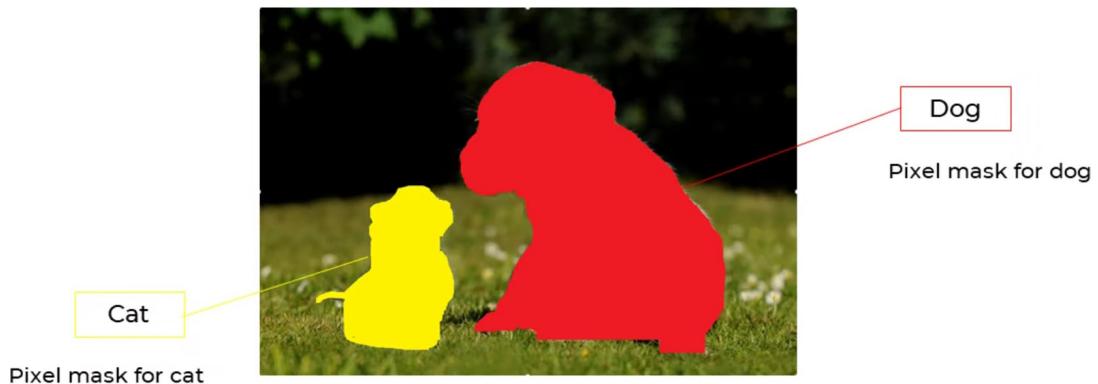


Image Segmentation



COCO 2020 Keypoint Detection Task

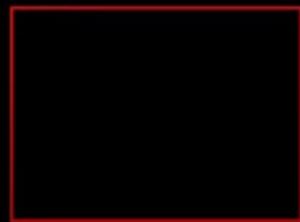


Sliding window object detection is a technique that allows you to detect objects in a picture. This technique is not very efficient as it is very compute intensive. Recently new techniques has been discovered that tried to improve performance such as R CNN, Fast R CNN, Faster R CNN etc. YOLO (You only look once) is a state of the art most modern technique that outperforms all other previous techniques such as sliding window object detection, R CNN, Fast and Faster R CNN etc. We will cover YOLO in future videos.



Picture in picture

Try different sizes iteratively



Disadvantages



Sliding Window
Object Detection

→ R CNN

→ Fast R CNN

→ Faster R CNN

→ YOLO

YOLO(You Only Look Once)

Image Classification

Is this a dog or a person?



Neural Network Output
Dog = 1
Person = 0

Object Localization

Where exactly is the dog in this image?



Neural Network Output
Dog = 1
Person = 0
+
Bounding Box

Object Localization



P_c 1
 B_x 50
 B_y 70
 B_w 60
 B_h 70
 C_1 1
 C_2 0

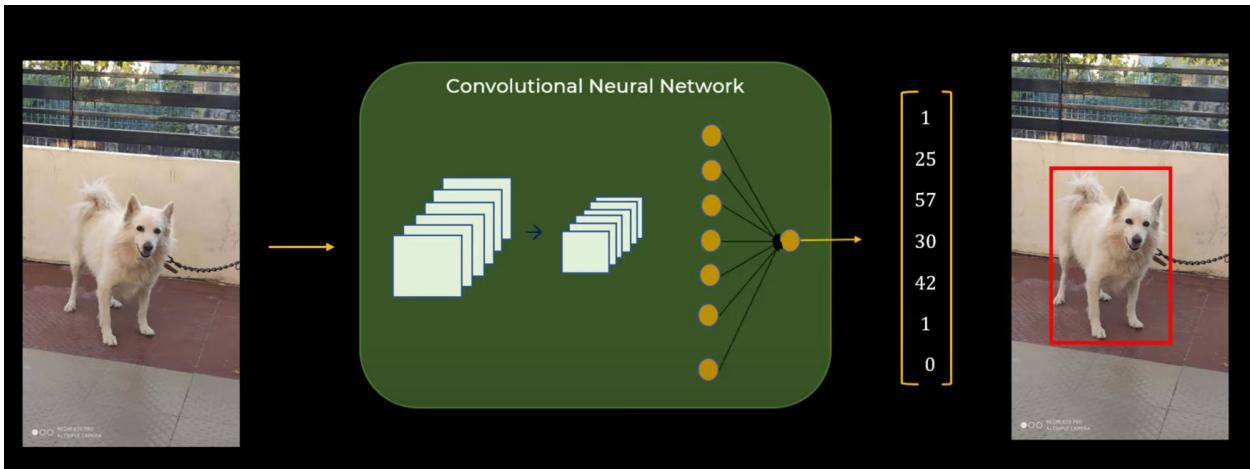
C_1 = Dog class
 C_2 = Person Class



1
30
28
28
82
0
1



0
-
-
-
-
-
-



This works ok
only for single
object. What
about multiple
objects in an
image?



Image Classification

Is this a dog or a person?



Neural Network Output

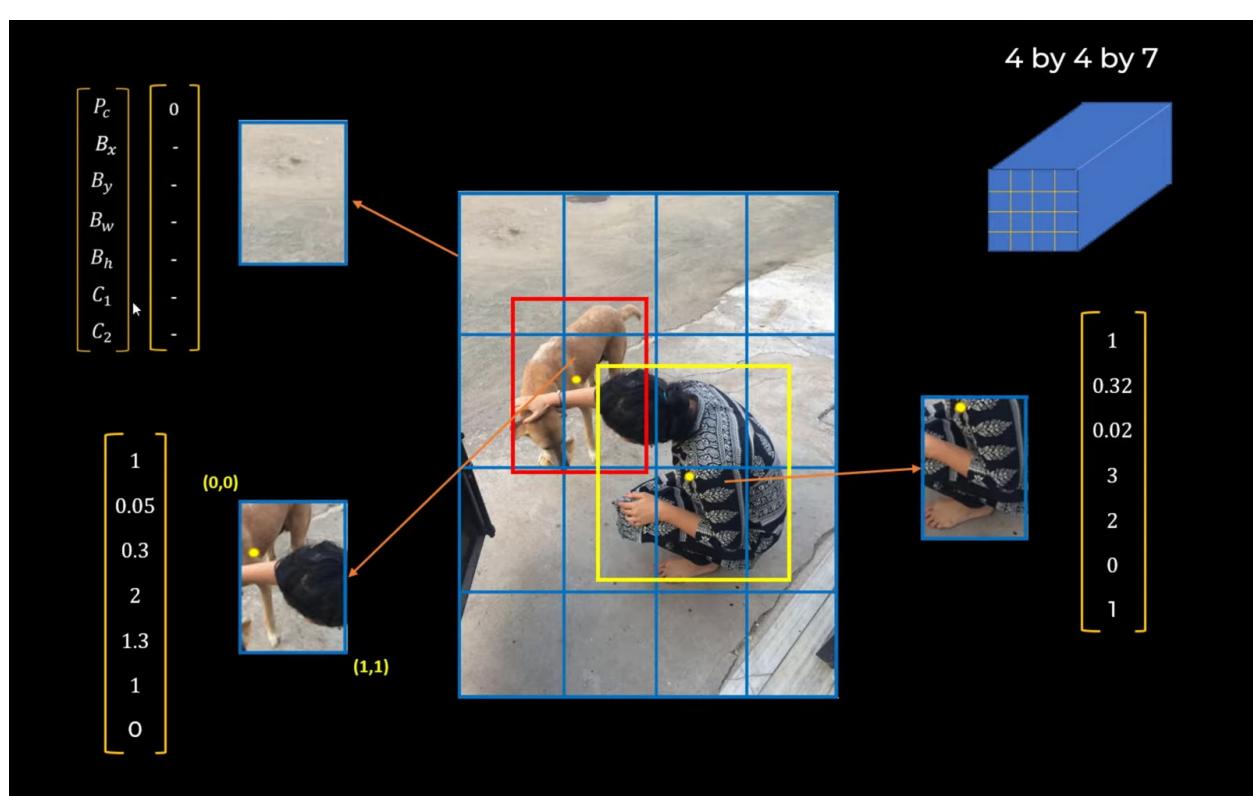
Dog = 1
Person = 0

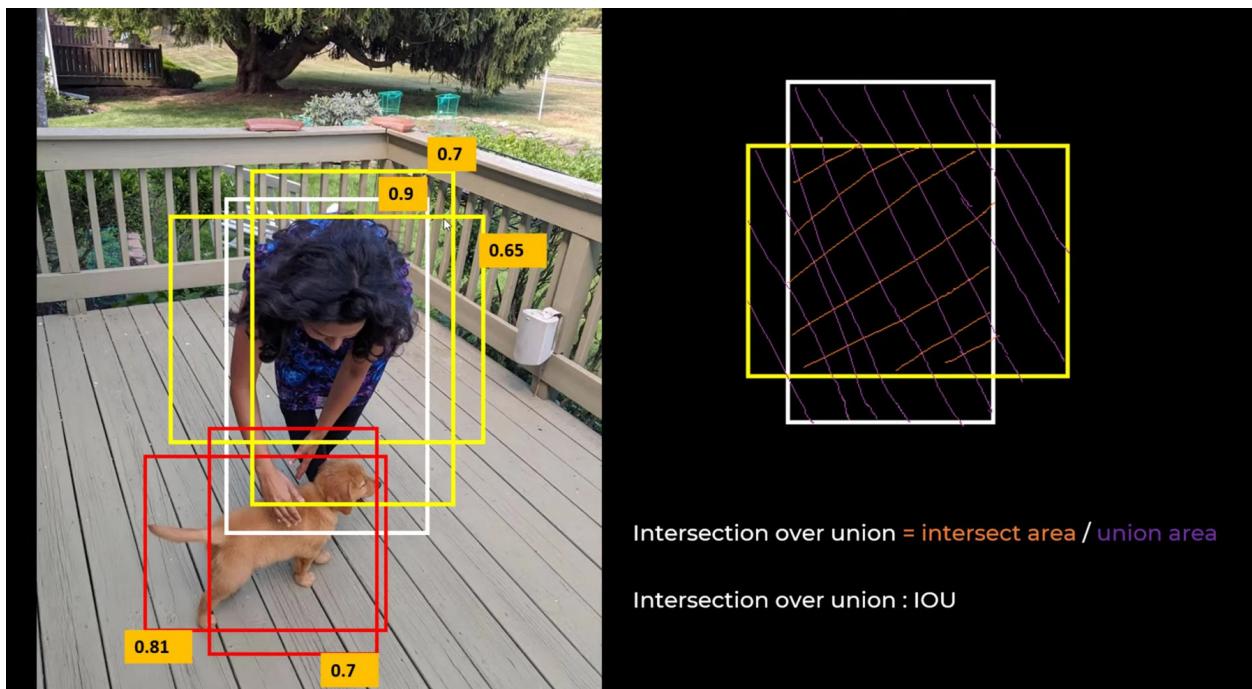
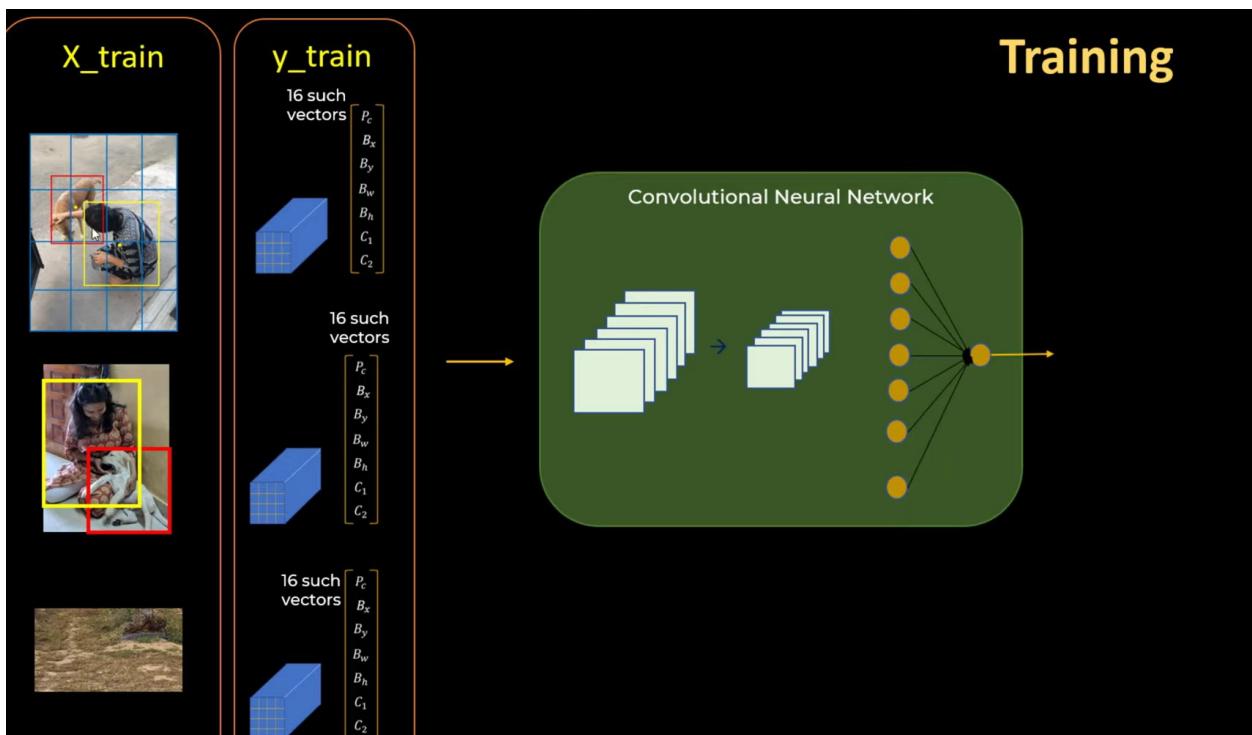
Object Localization

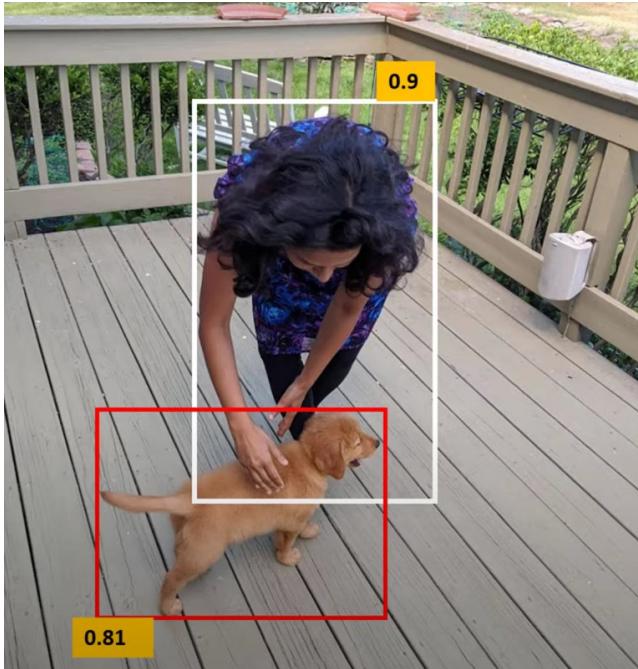
Where exactly is the dog in this image?



Neural Network Output
Dog = 1
Person = 0
+
Bounding Box

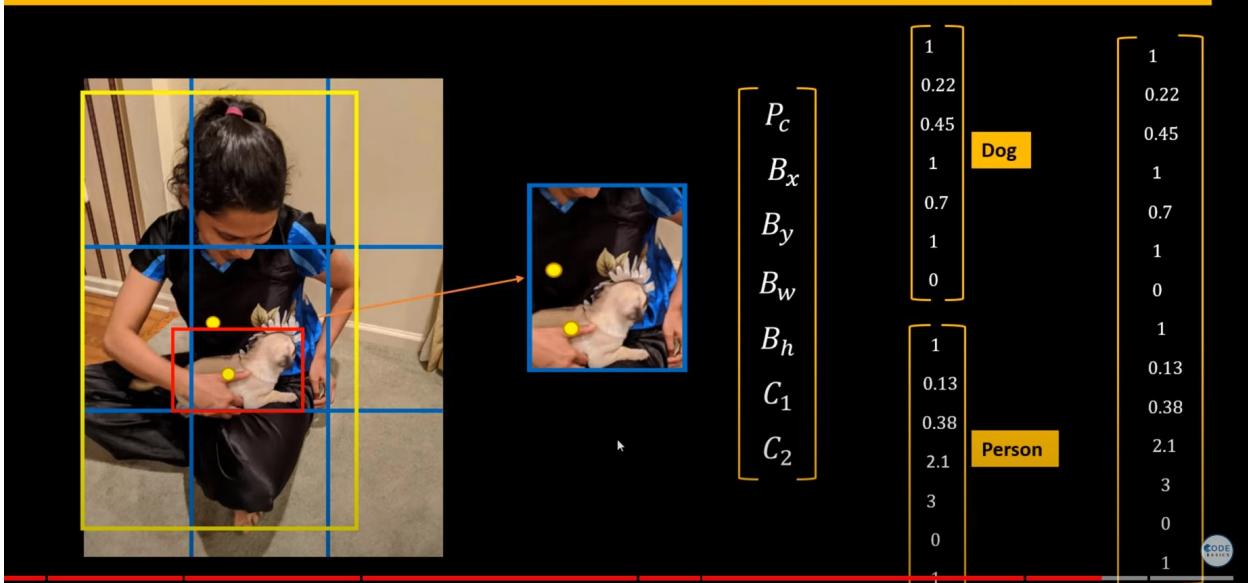




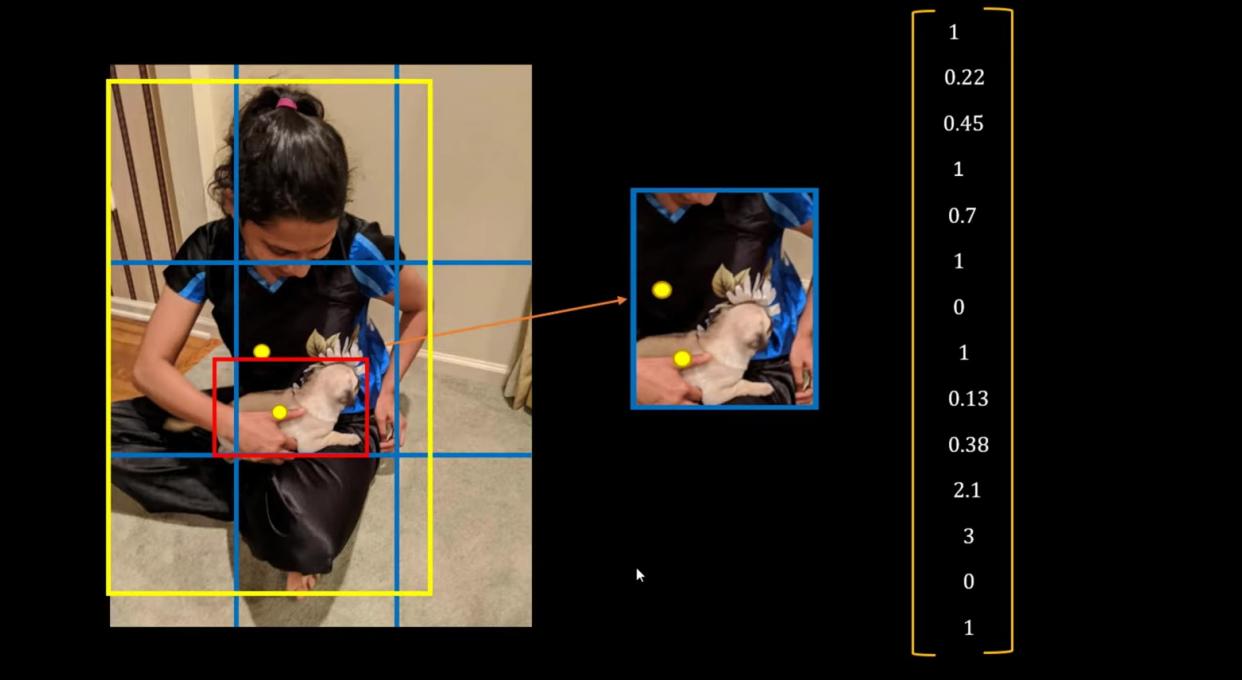


Non max suppression

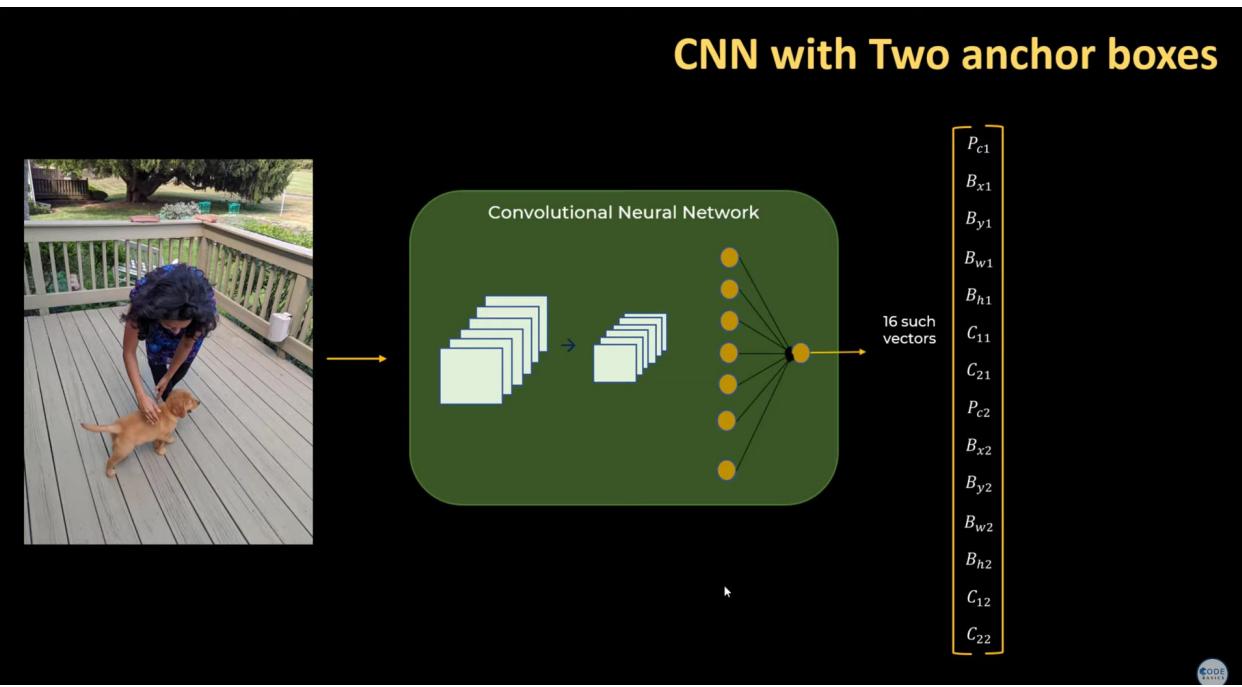
What if one grid cell has center of two objects?



This concept is called anchor boxes



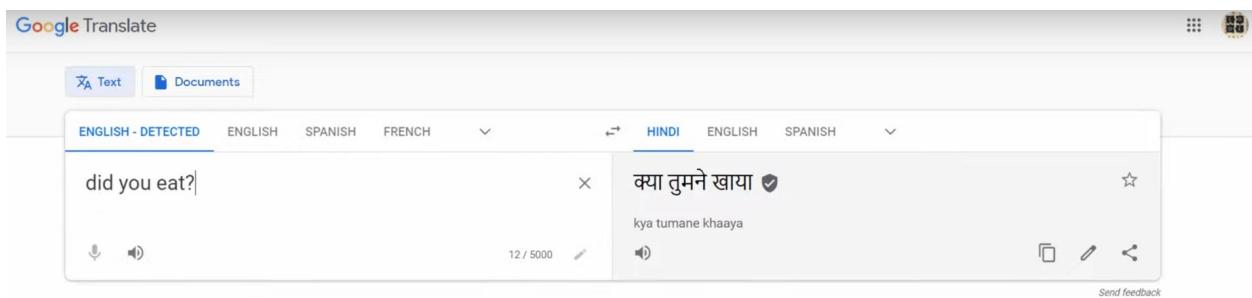
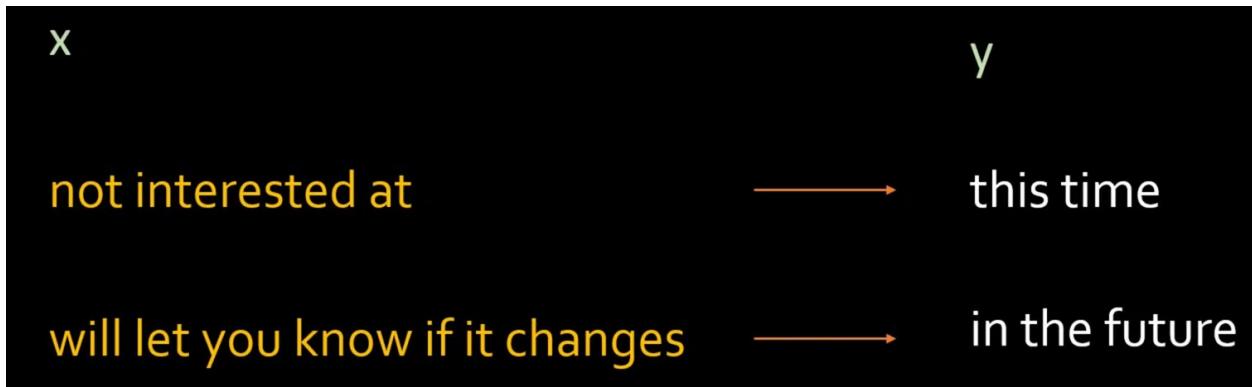
CNN with Two anchor boxes



RNN

RNN or Recurrent Neural Network are also known as sequence models that are used mainly in the field of natural language processing as well as some other areas such as speech to text

translation, video activity monitoring, etc. In this video we will understand the intuition behind RNN and see how RNN's work.

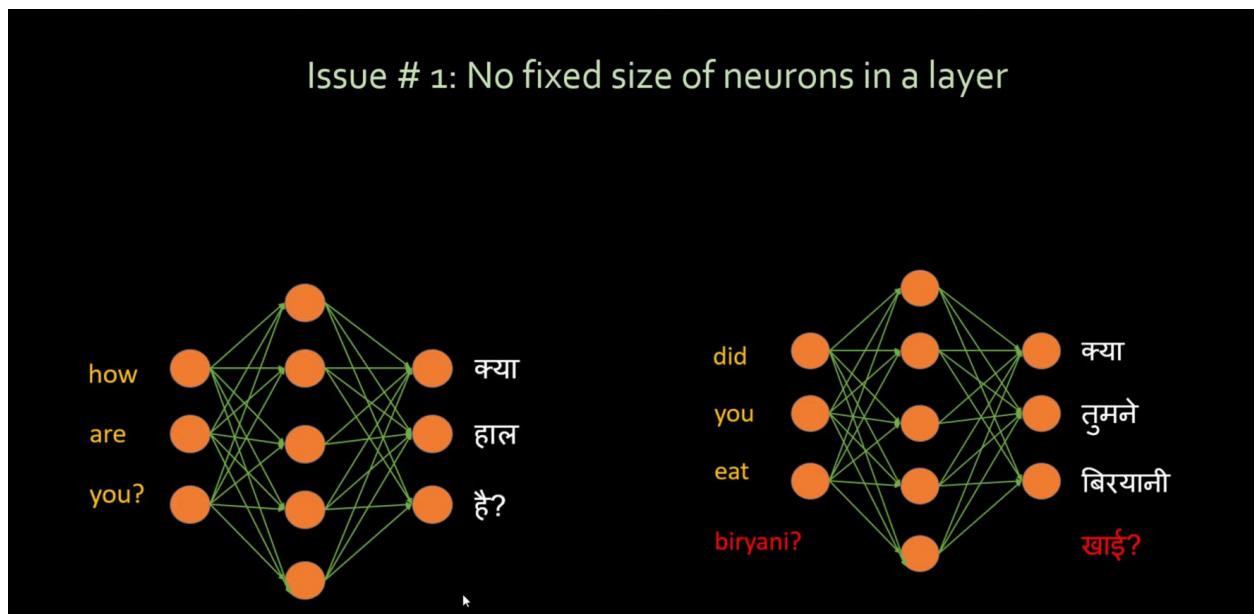


X Rudolph Smith bought 1000 shares of tesla Inc. in March 2020

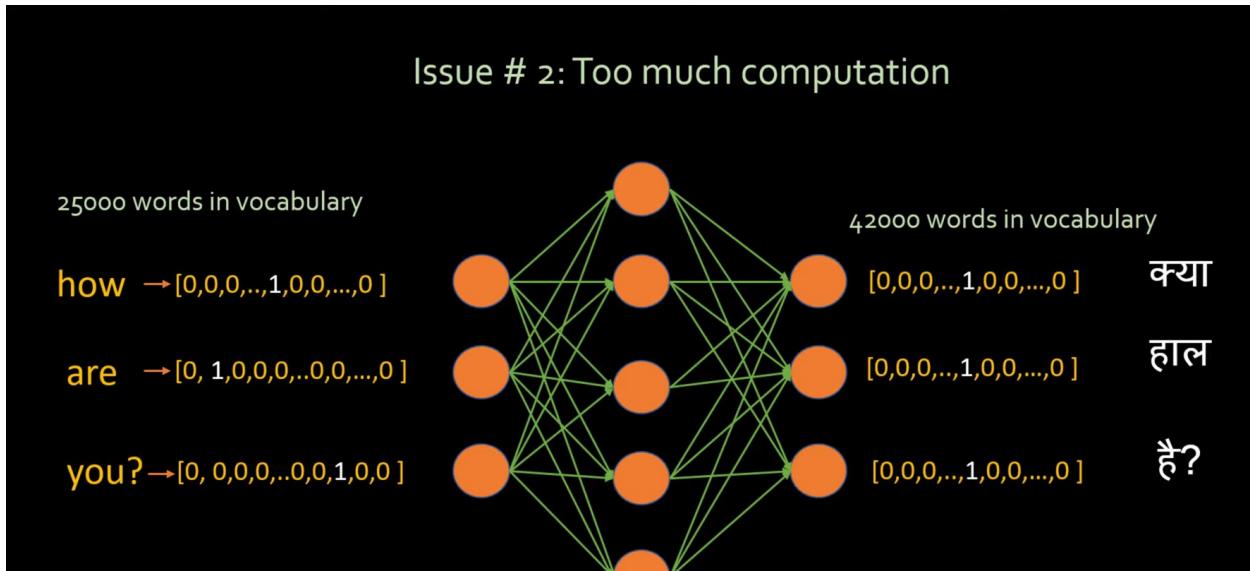
Y Person Rudolph Smith Company tesla Inc. time March 2020

NER: Named Entity Recognition

	x		y
auto complete	not interested at	→	this time
translation	how are you?	→	क्या हाल है?
NER	Rudolph Smith bought 1000 shares of tesla Inc. in March 2020	→	Person Person Company Time
Sentiment Analysis	Not only the fan was expensive, but it was broken when it arrived.	→	★☆☆☆☆



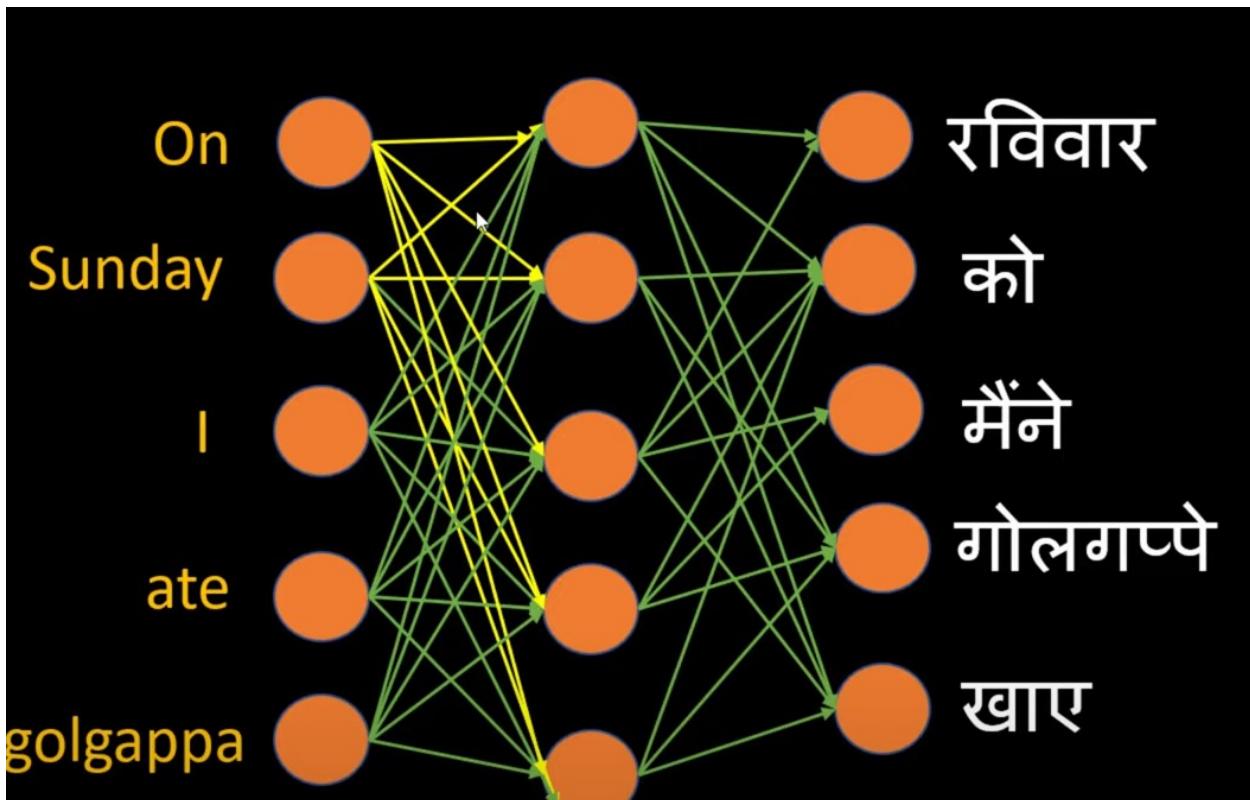
Issue # 2: Too much computation

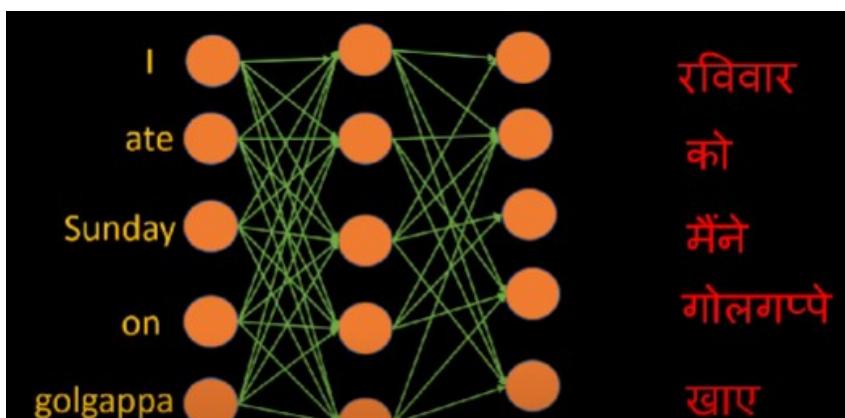
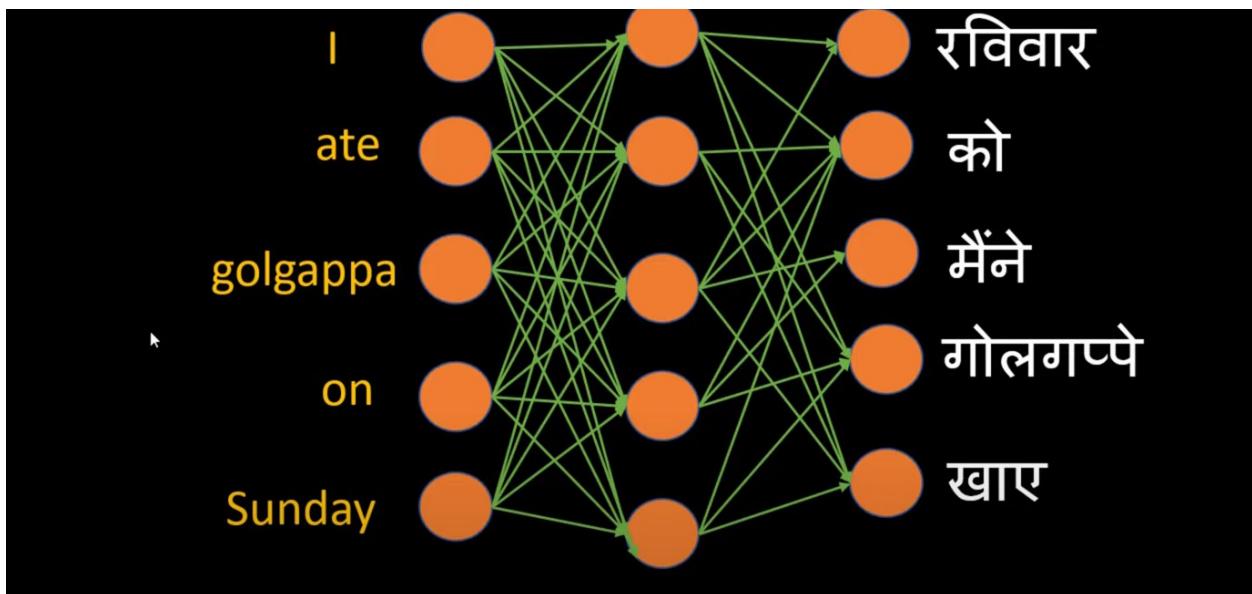


On sunday I ate golgappa

रविवार को मैंने गोलगप्पे खाए

I ate golgappa on Sunday





3 Issues using ANN for sequence problems

Variable size of
input/output
neurons

Too much
computation

No parameter
sharing

Let's once again talk about Named Entity



Named Entity Recognition

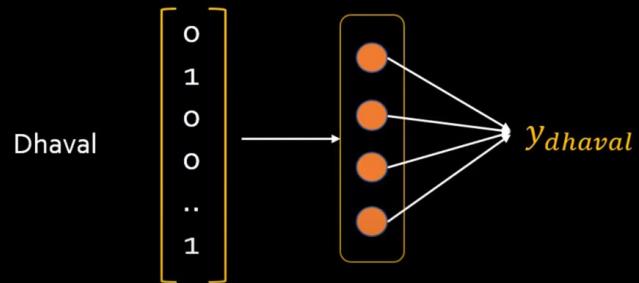


Dhaval loves baby yoda

person person
Dhaval loves baby yoda
1 0 1 1

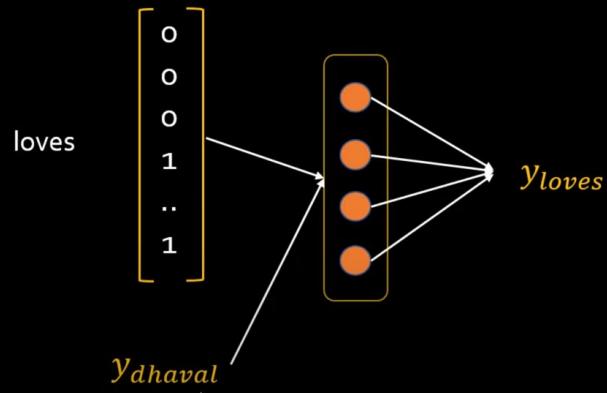
Named Entity Recognition

Dhaval loves baby yoda



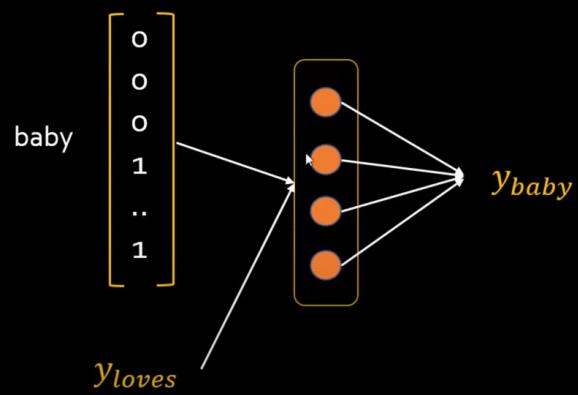
Named Entity Recognition

Dhaval loves baby yoda

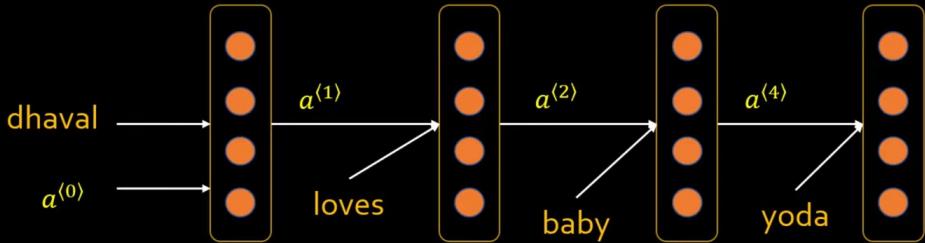


Named Entity Recognition

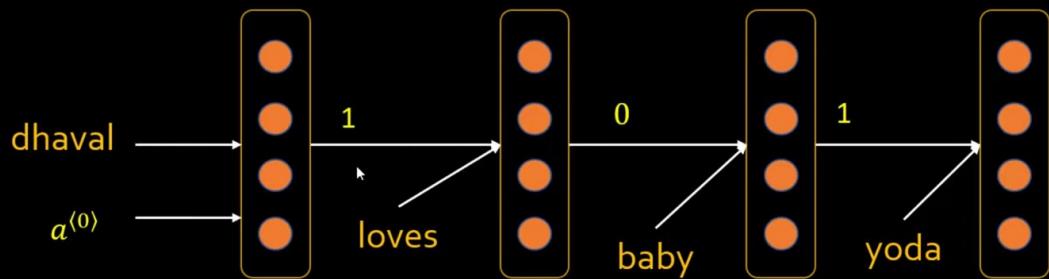
Dhaval loves baby yoda



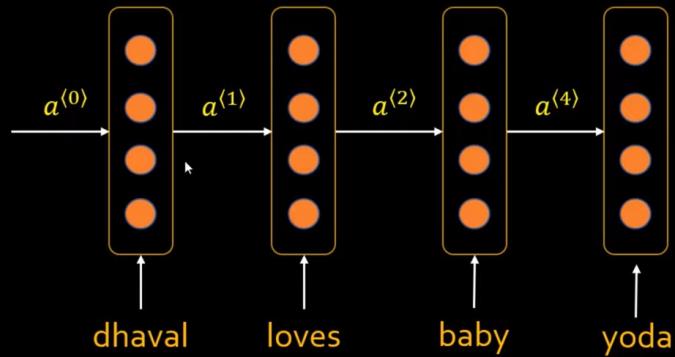
Named Entity Recognition



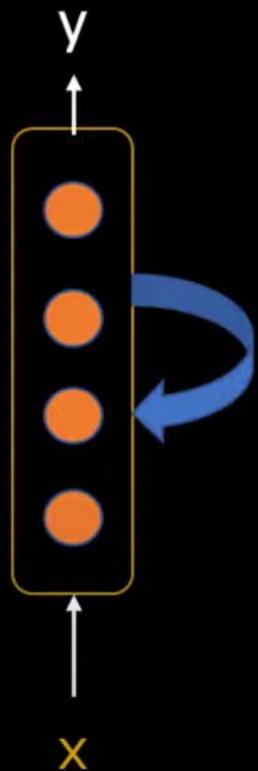
Named Entity Recognition: once network is trained



Named Entity Recognition



Generic Representation of RNN



Training : Named Entity Recognition (NER)

X

y

Dhaval loves baby yoda

1 0 1 1

Bob told Ahmed that pizza is delivered

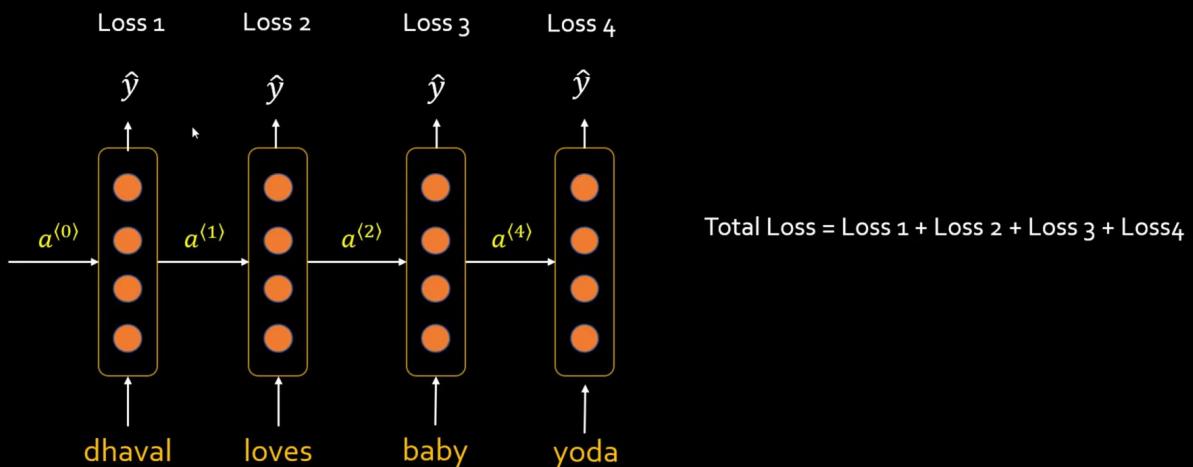
1 0 1 0 0 0 0

Ironman punched on hulk's face

1 0 0 1 1

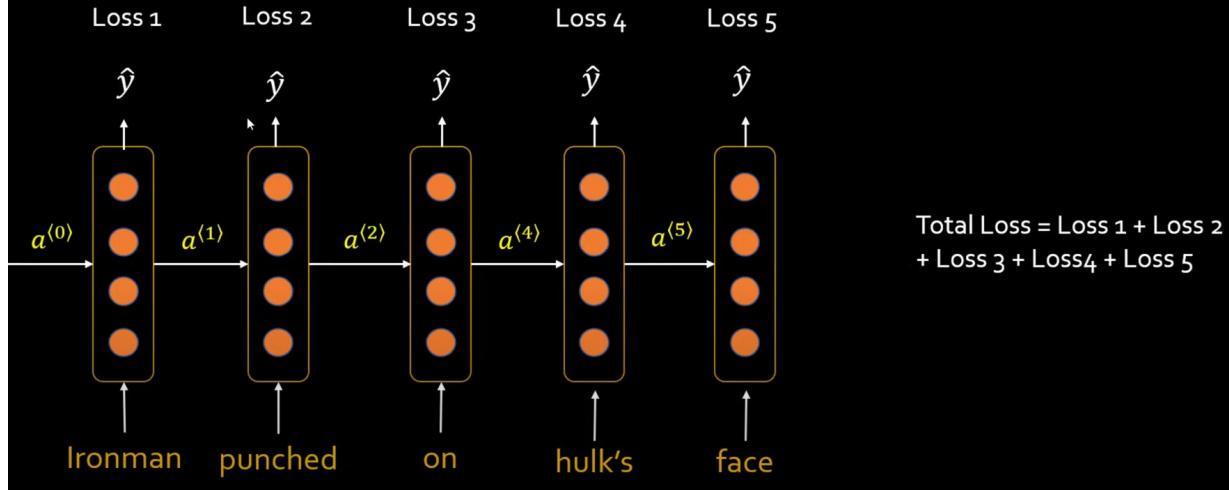
Training

Dhaval loves baby yoda \rightarrow 1 0 1 1



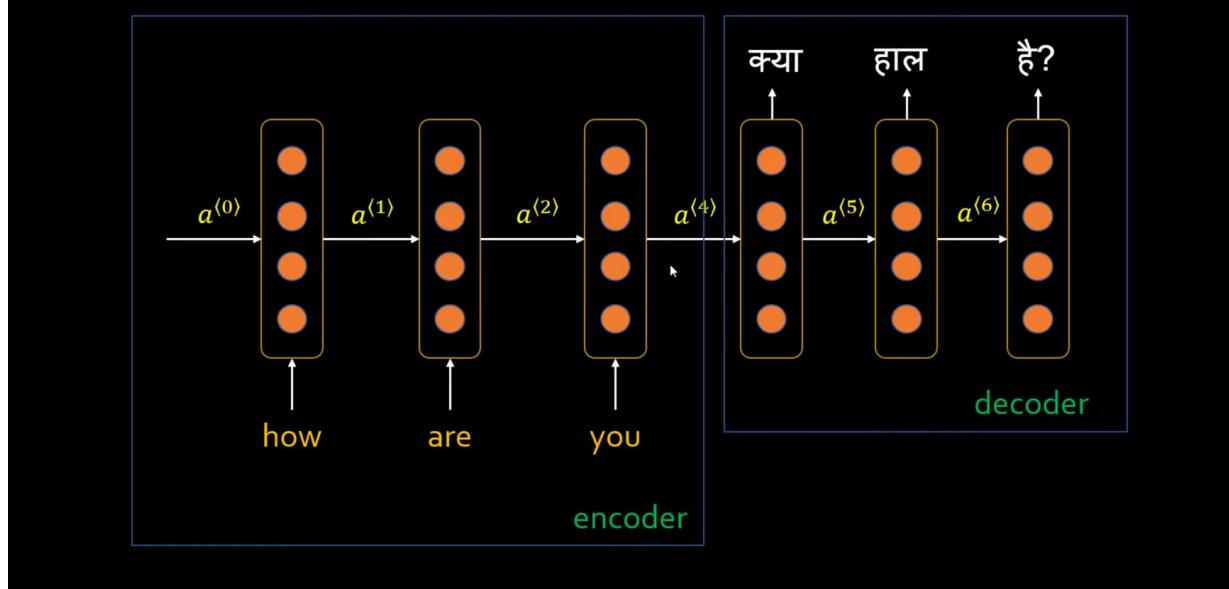
Training

Ironman punched on hulk's face $\rightarrow 10010$

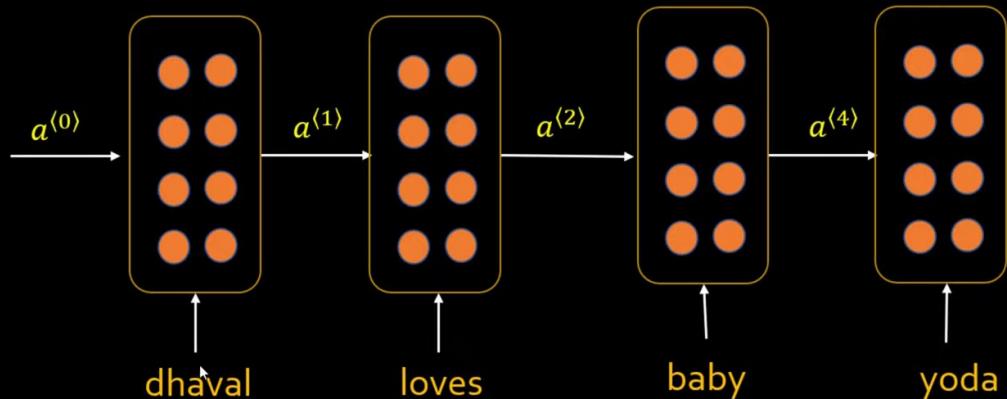


Language translation

how are you? क्या हाल है?

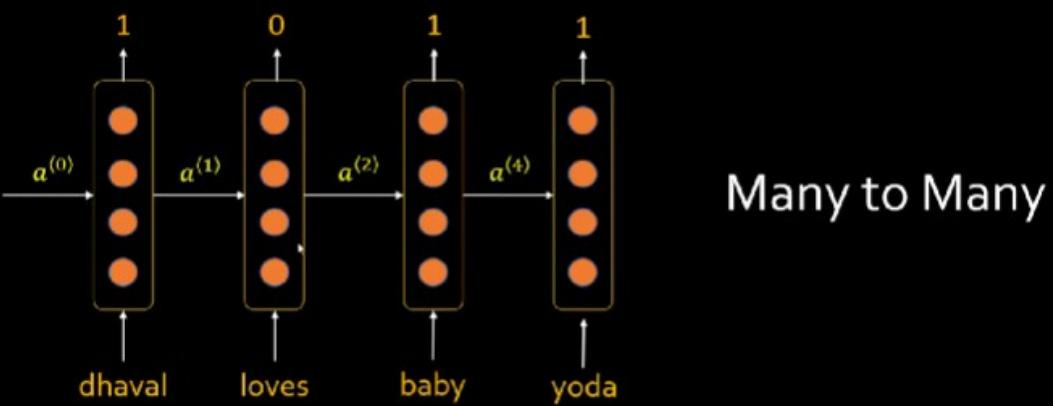


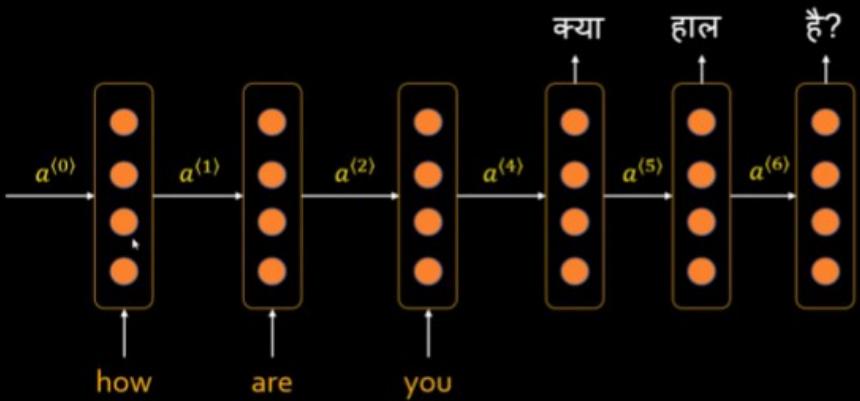
Deep RNN



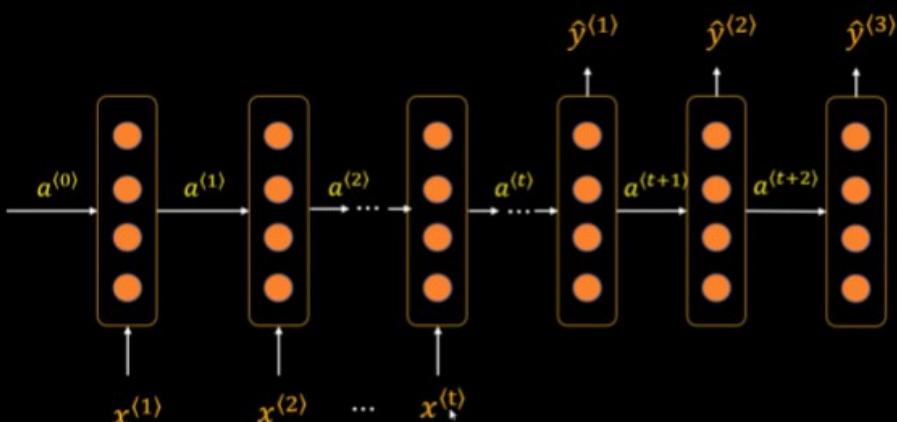
Deep RNN has multiple hidden layers.

In this video we will discuss different types of RNN types such as, 1) One to many 2) Many to many 3) Many to one





Many to Many



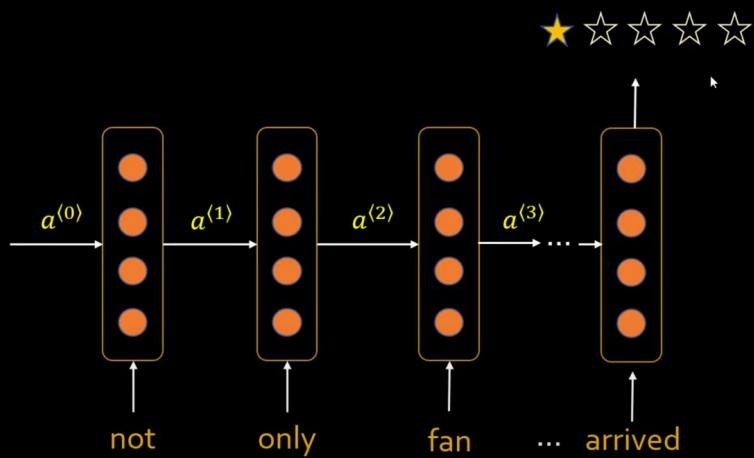
Many to Many

Sentiment Analysis

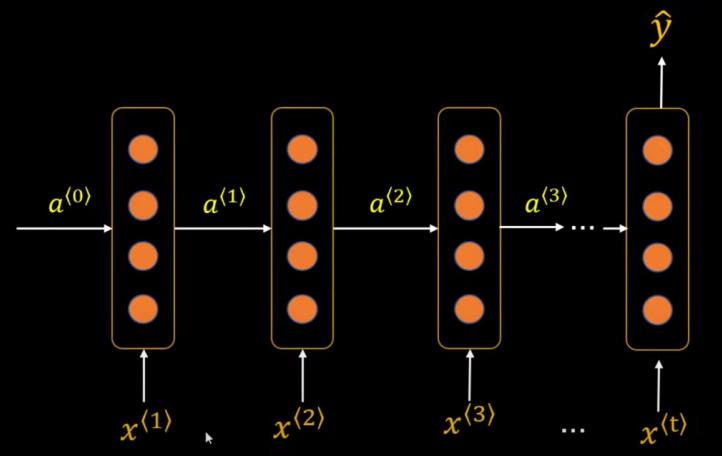
Not only the fan was expensive,
but it was broken when it arrived.



The fan works like a charm, I
wasn't expecting such a good
quality at this cheap price



Many to One

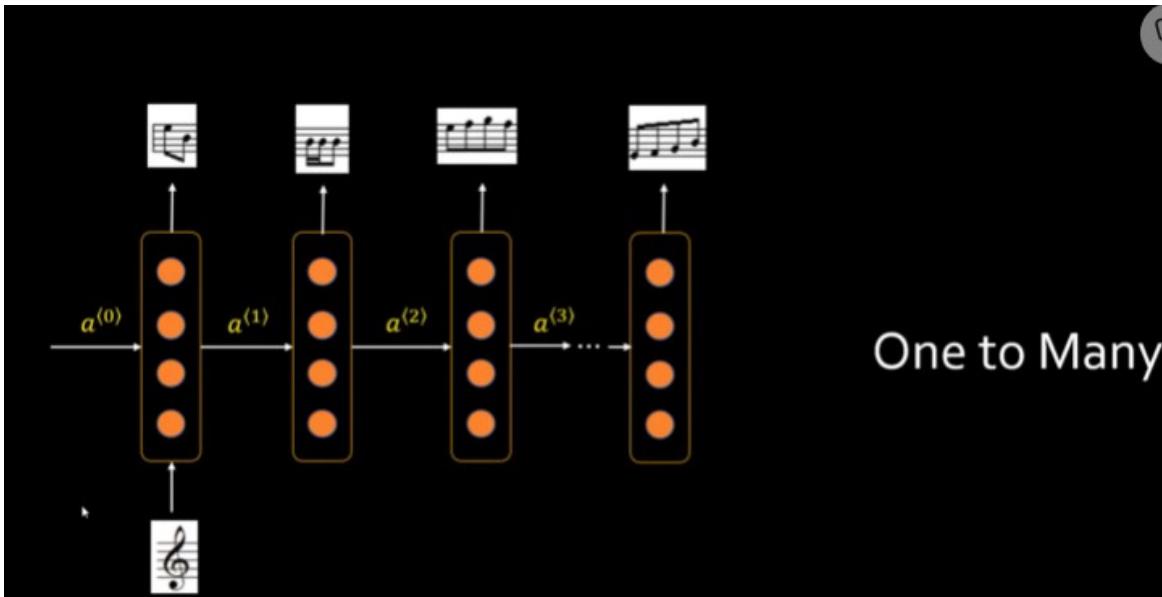


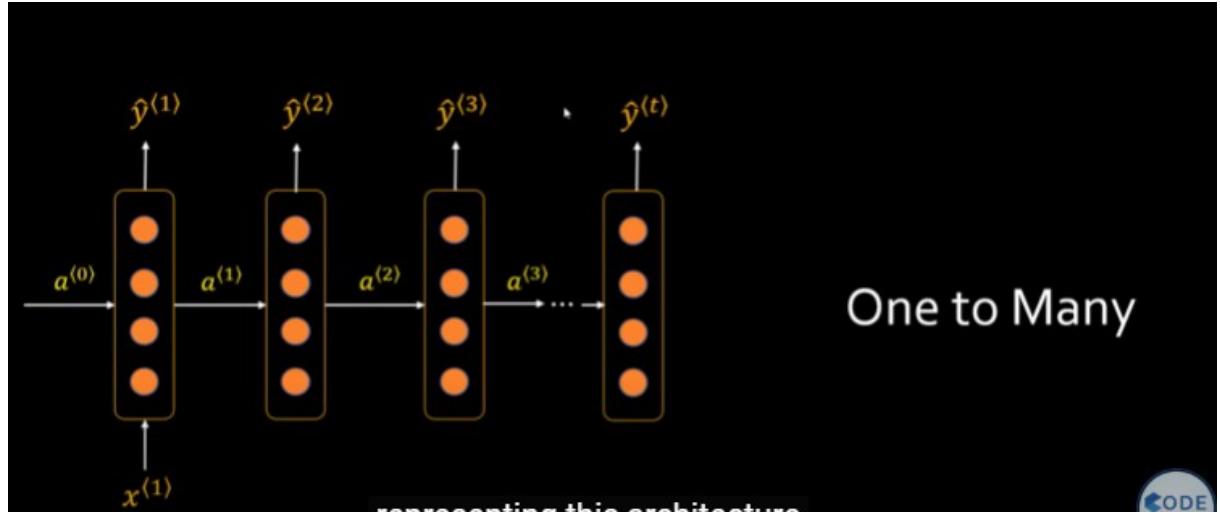
Many to One

Music Generation

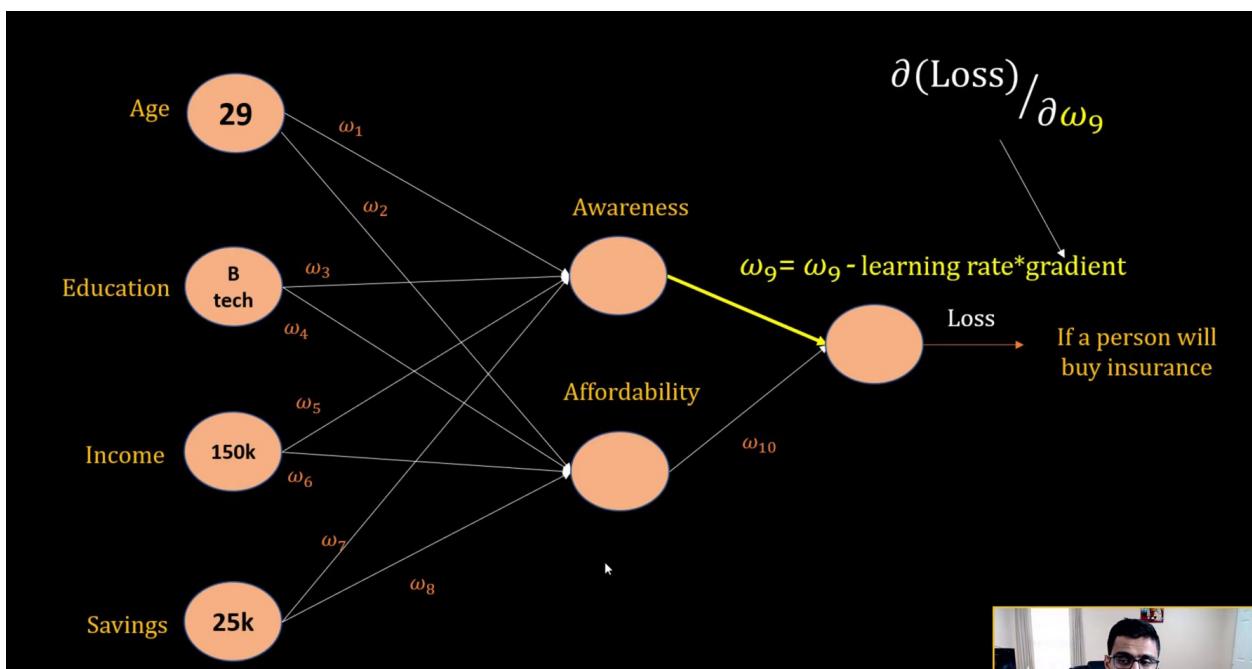


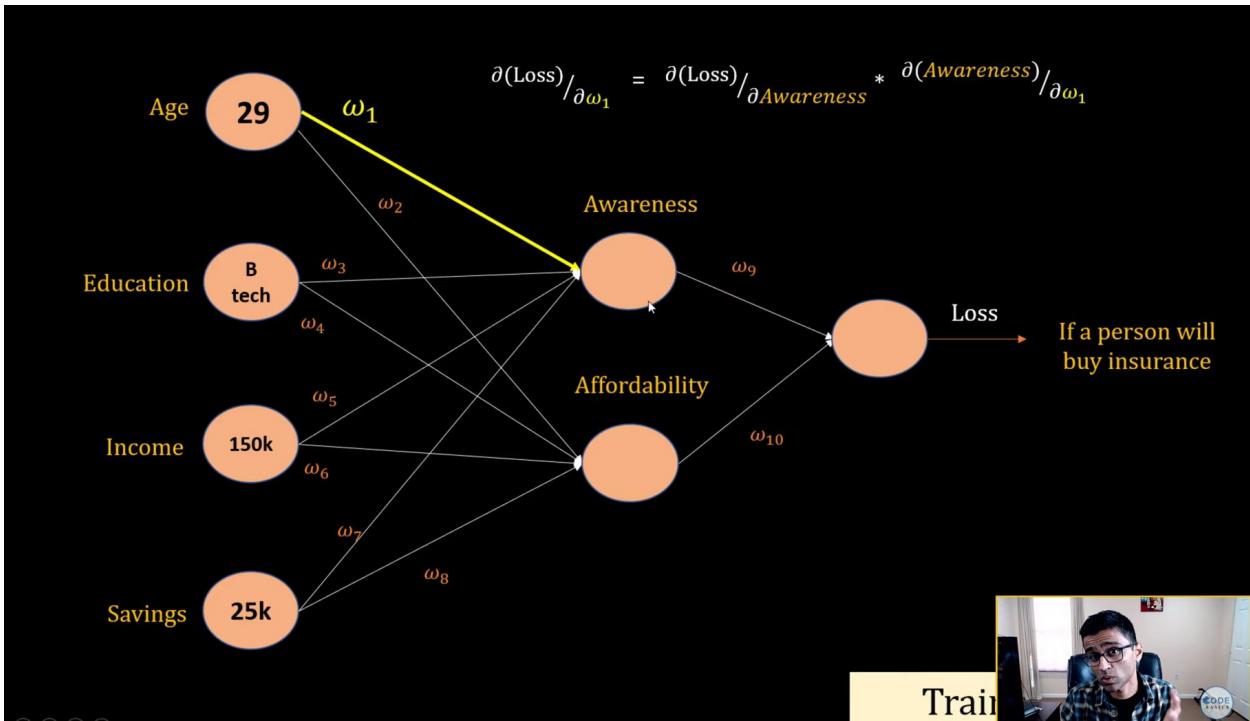
One to Many





Vanishing gradient is a common problem encountered while training a deep neural network with many layers. In case of RNN this problem is prominent as unrolling a network layer in time makes it more like a deep neural network with many layers. In this video we will discuss what vanishing and exploding gradients are in artificial neural network (ANN) and in recurrent neural network (RNN)





Train

As number of hidden layers grow, gradient becomes very small and weights will hardly change . This will hamper the learning process.

Vanishing Gradients

When individual derivatives are large, the final derivative will also become huge and weights would change drastically.

Exploding Gradients

$$\text{gradient} = d1 * d2 * d3 * d4 * \dots * dn$$

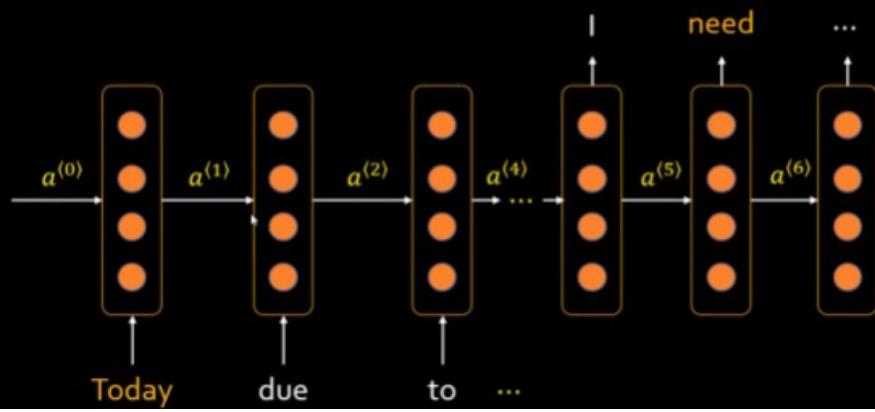
Vanishing gradient problem is more prominent in very deep neural networks.

Vanishing gradient problem in RNN

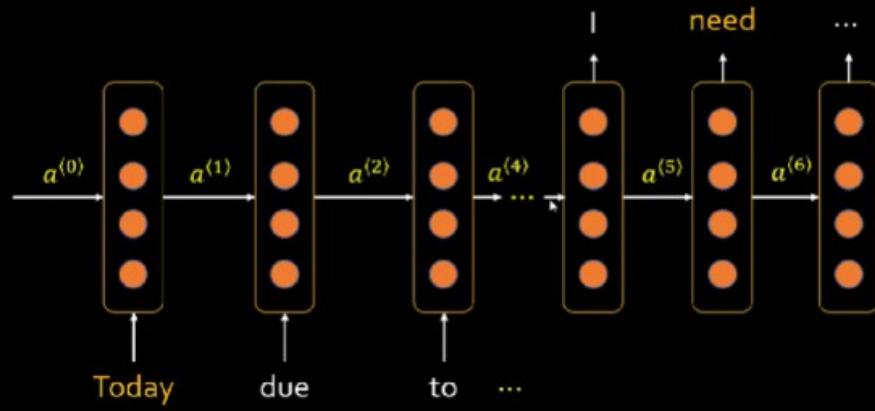
Today, due to my current job situation and family conditions, I need to take a loan.

Last year, due to my current job situation and family conditions, I had to take a loan.

Today, due to my current job situation and family conditions, I need to take a loan.



Today, due to my current job situation and family conditions, I need to take a loan.

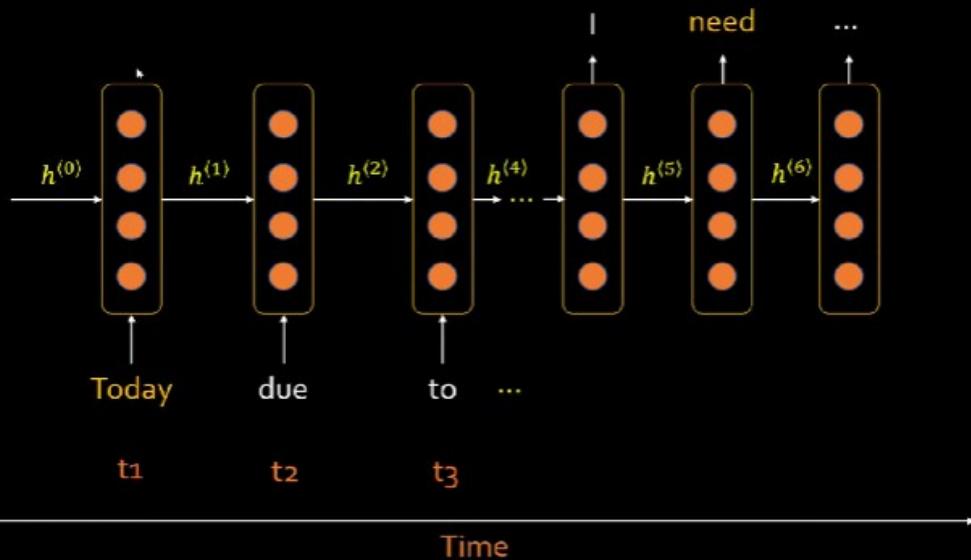


GRU

LSTM

To address the problem of vanishing effect, we use GRU and LSTM

Today, due to my current job situation and family conditions, I need to take a loan.



Due to short length memory of RNN, using word "Today/Last Year" at beginning would not be able to predict "need/had".

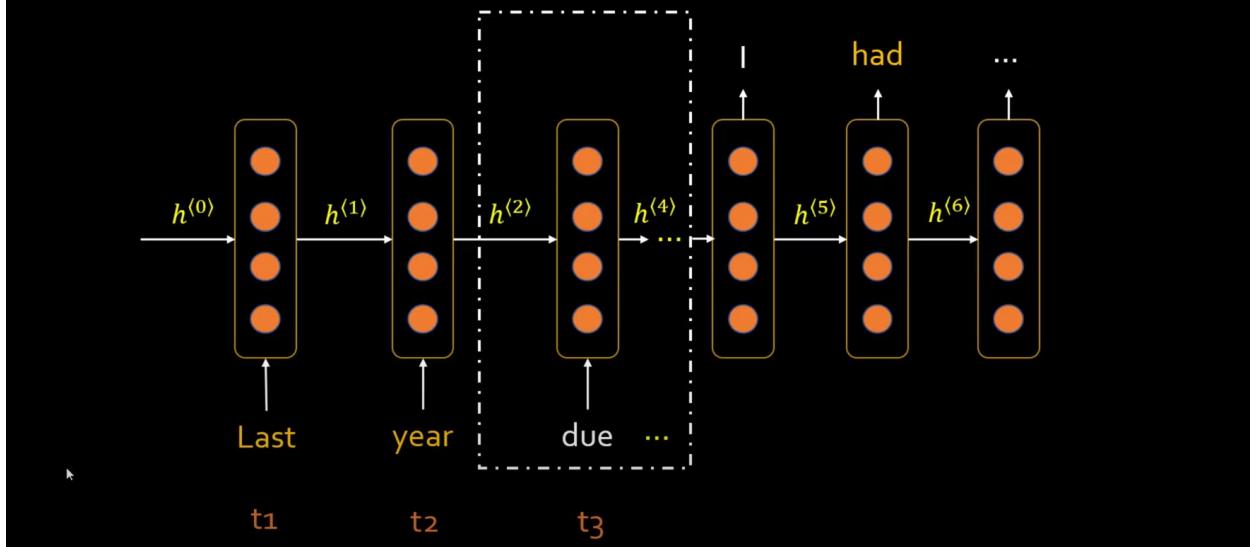
LSTM

LSTM or long short term memory is a special type of RNN that solves traditional RNN's short term memory problem. In this video I will give a very simple explanation of LSTM using some real life examples so that you can understand this difficult topic easily. Also refer to following blogs to explore math and understand few more details.

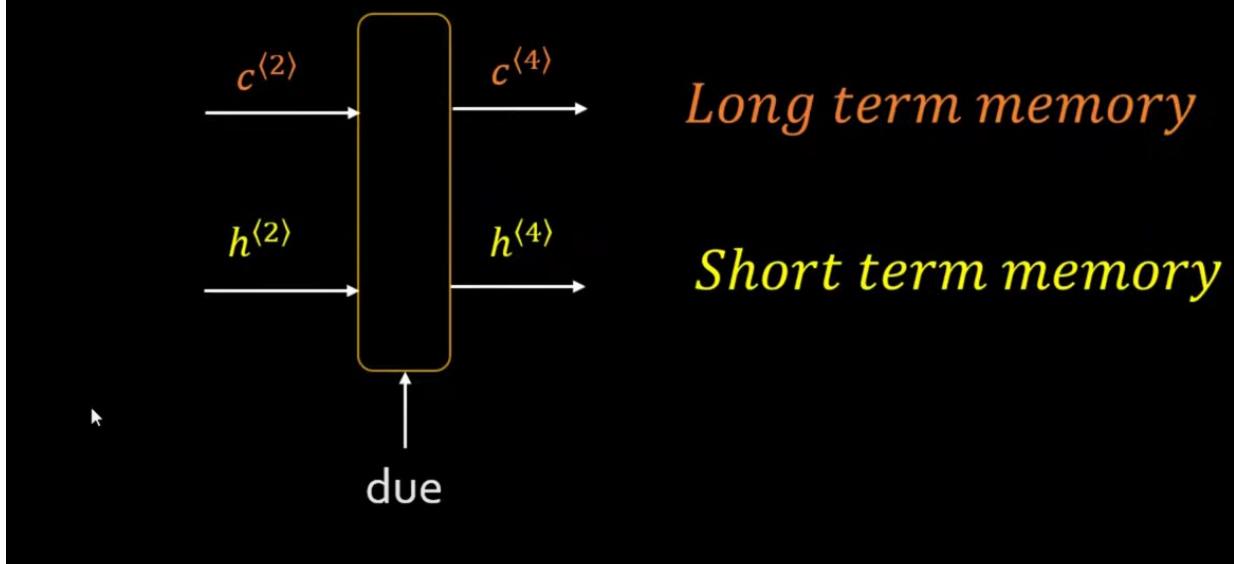
<http://colah.github.io/posts/2015-08-...>

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

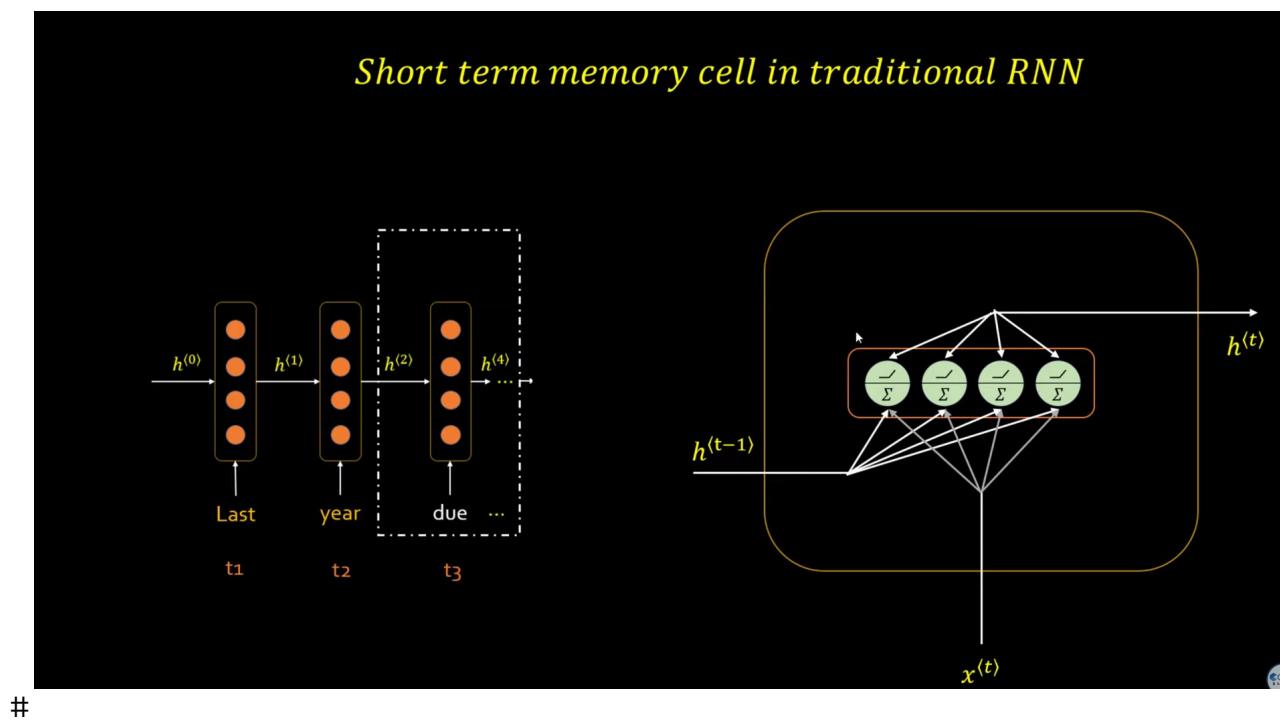
Last year, due to my current job situation and family conditions, I had to take a loan.



Memory cell

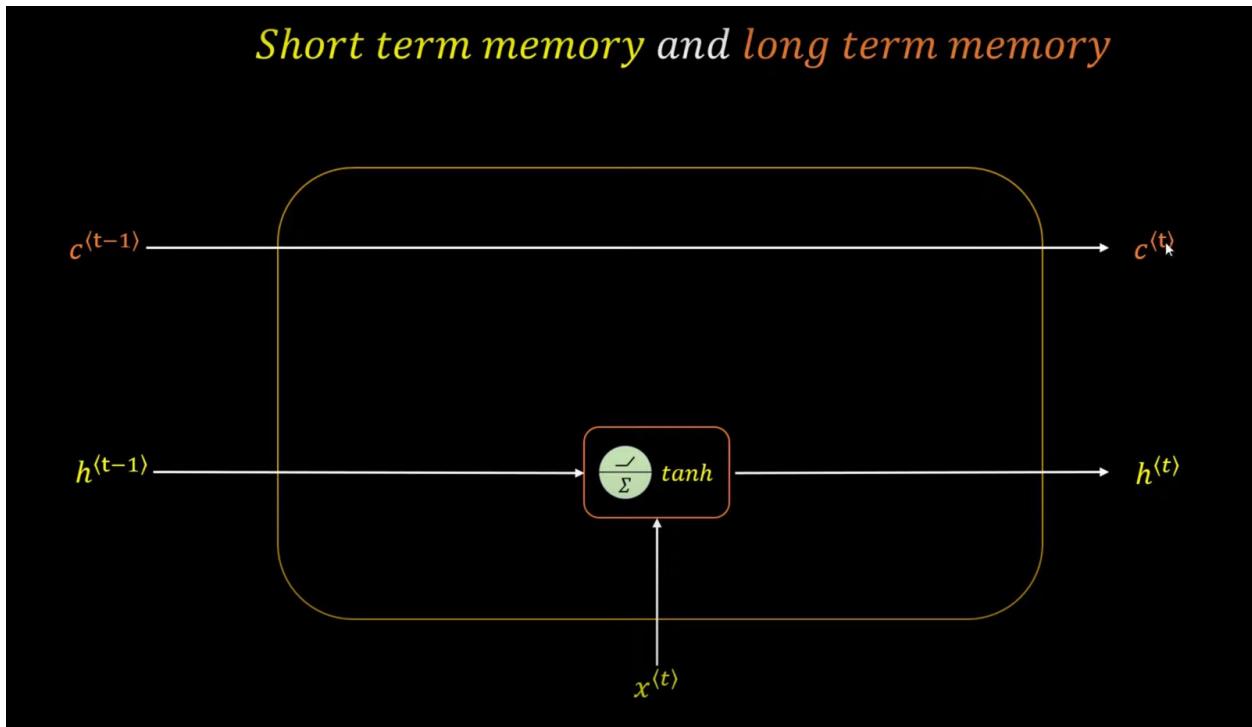


Short term memory cell in traditional RNN



#

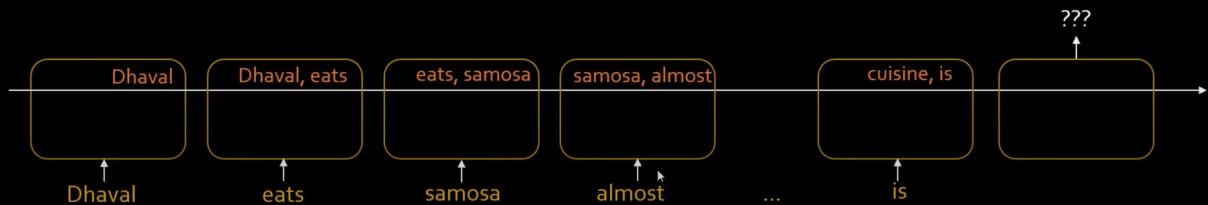
Short term memory and long term memory



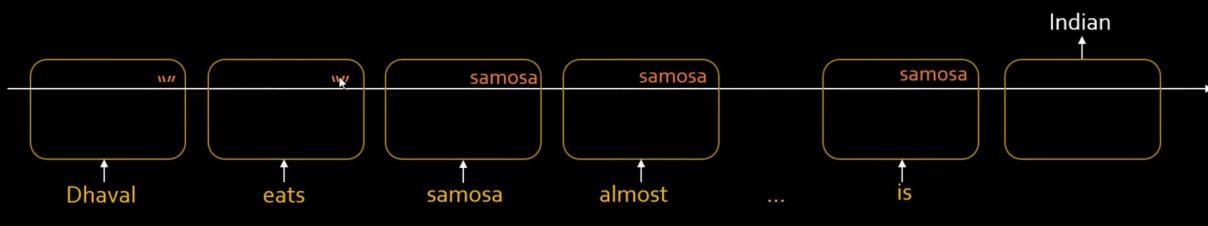
Dhaval eats samosa almost everyday, it shouldn't be hard to guess that his favorite cuisine is Indian



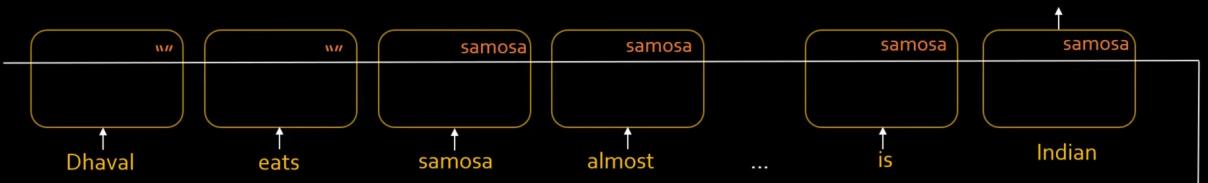
Traditional RNN

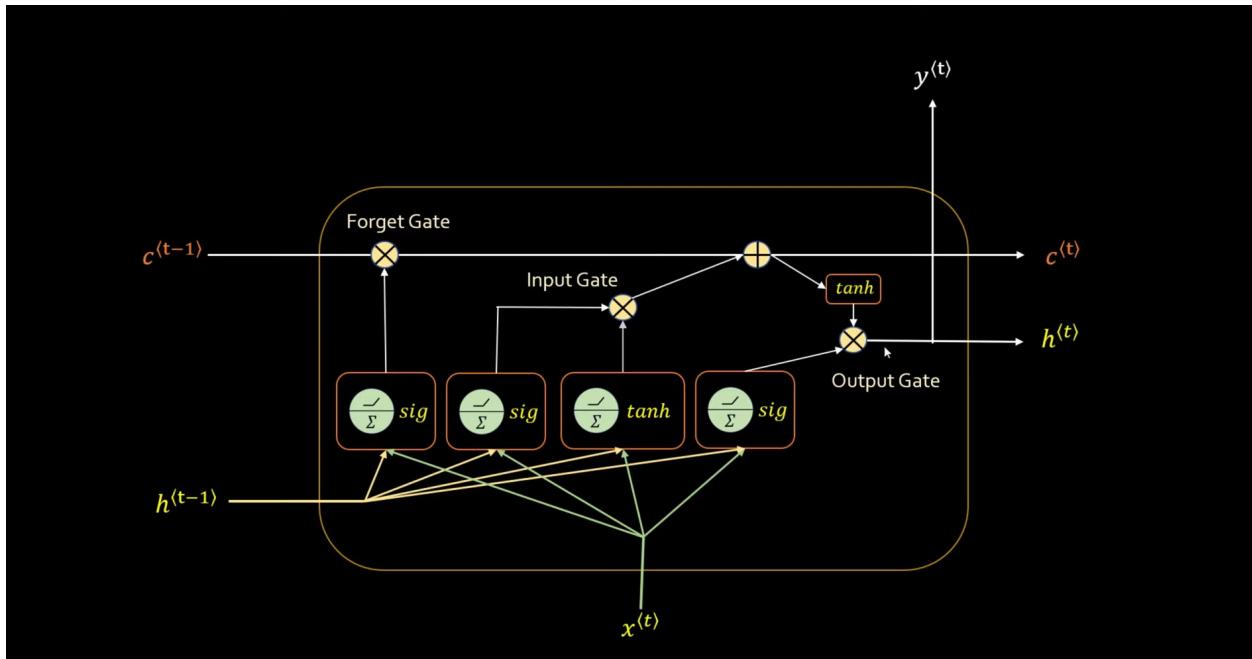


Dhaval eats samosa almost everyday, it shouldn't be hard to guess that his favorite cuisine is Indian



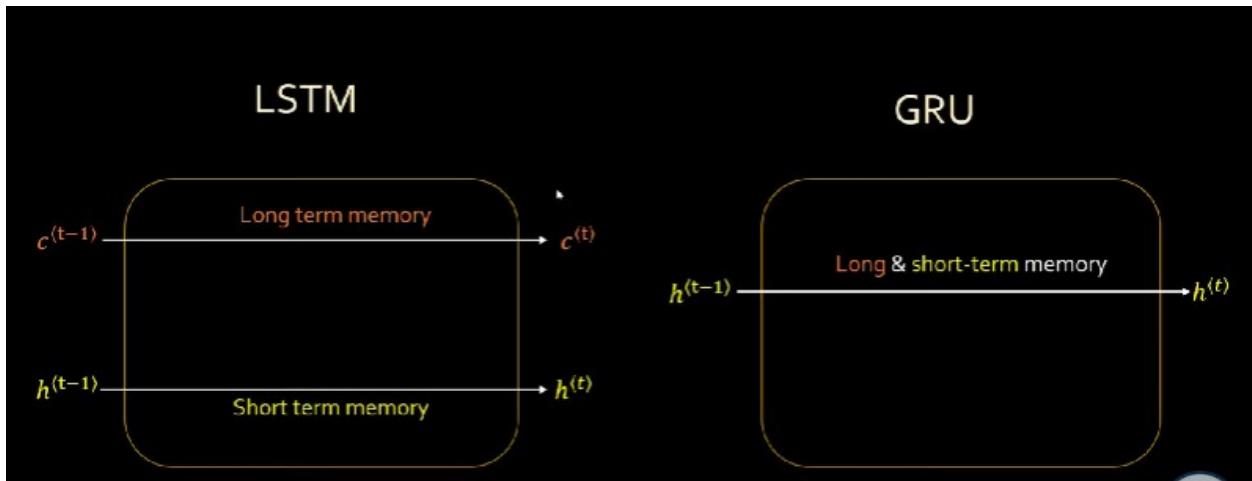
Dhaval eats samosa almost everyday, it shouldn't be hard to guess that his favorite cuisine is Indian. His brother Bhavin however is a lover of pasta and cheese that means Bhavin's favorite cuisine is Italian



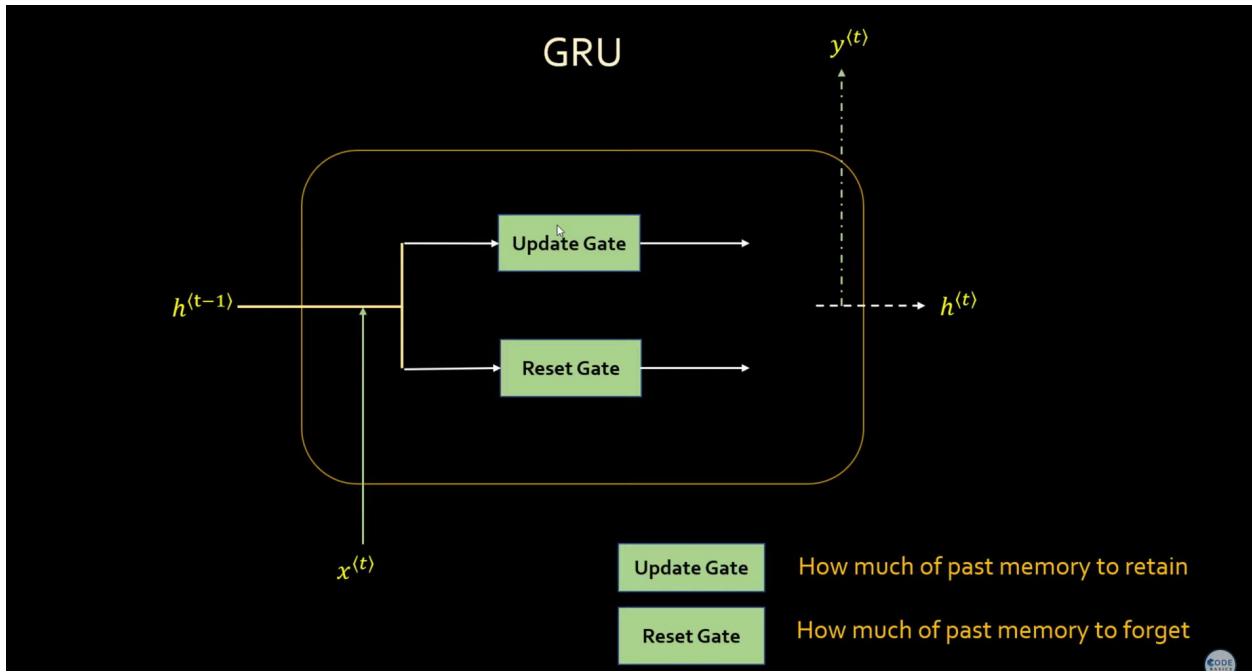


GRU

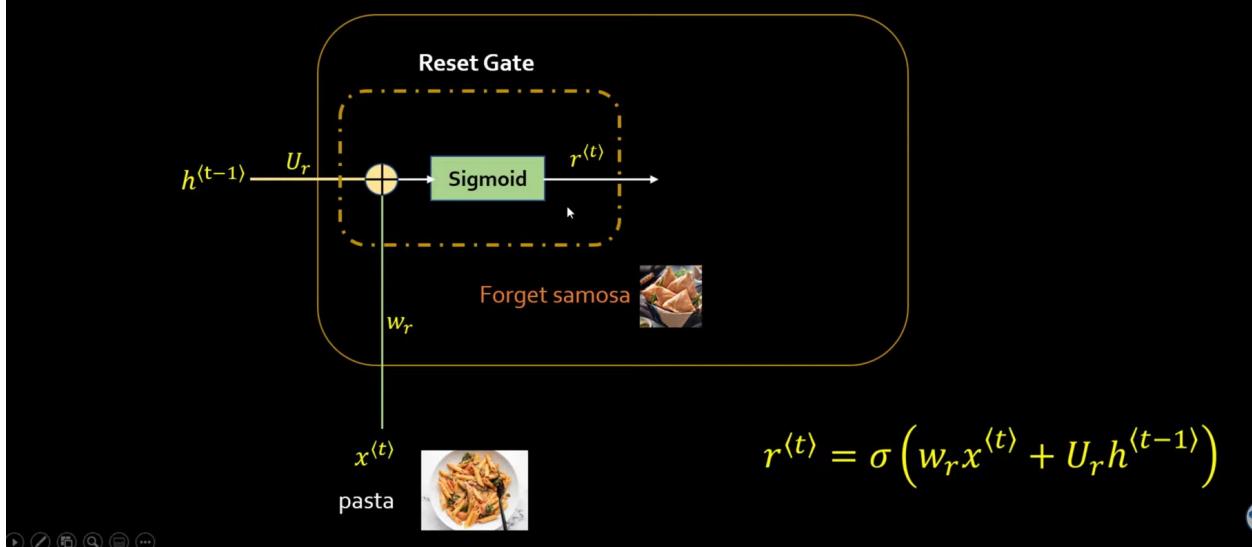
Simple Explanation of GRU (Gated Recurrent Units): Similar to LSTM, Gated recurrent unit addresses short term memory problem of traditional RNN. It was invented in 2014 and getting more popular compared to LSTM.

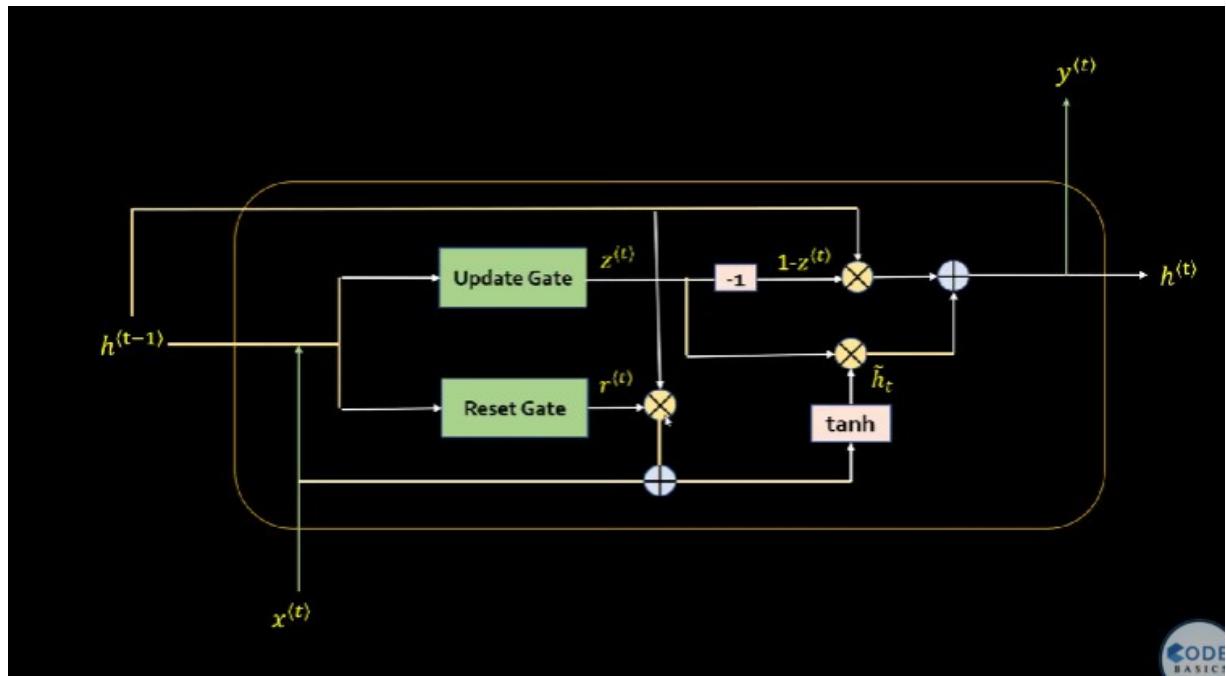
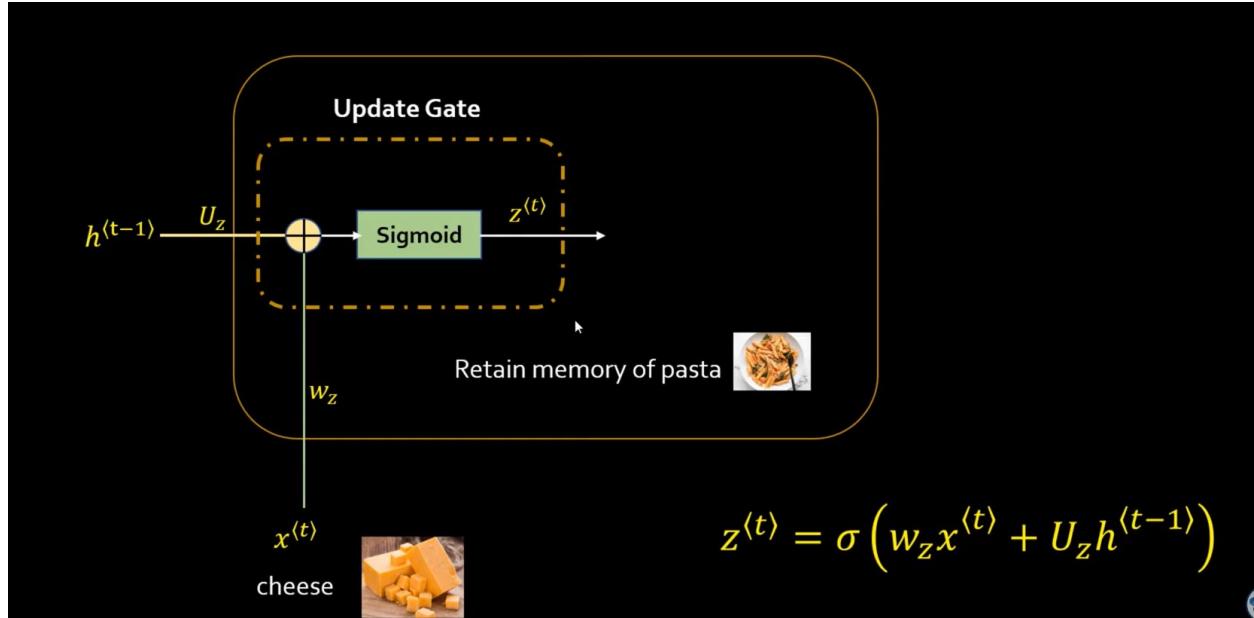


LSTM has two cell state (one is hidden state) whereas GRU has only hidden state as cell state which can combine both long and short term memory. LSTM has three gates:input,output and forget .GRU has two gates:Update gate and reset gate.



Dhaval eats samosa almost everyday, it shouldn't be hard to guess that his favorite cuisine is Indian. His broth Bhavin however is a lover of pasta and cheese that means Bhavin's favorite cuisine is Italian





$$\begin{aligned}
 z_t &= \sigma_g(W_z x_t + U_z h_{t-1} + b_z) \\
 r_t &= \sigma_g(W_r x_t + U_r h_{t-1} + b_r) \\
 \hat{h}_t &= \phi_h(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h) \\
 h_t &= (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t
 \end{aligned}$$

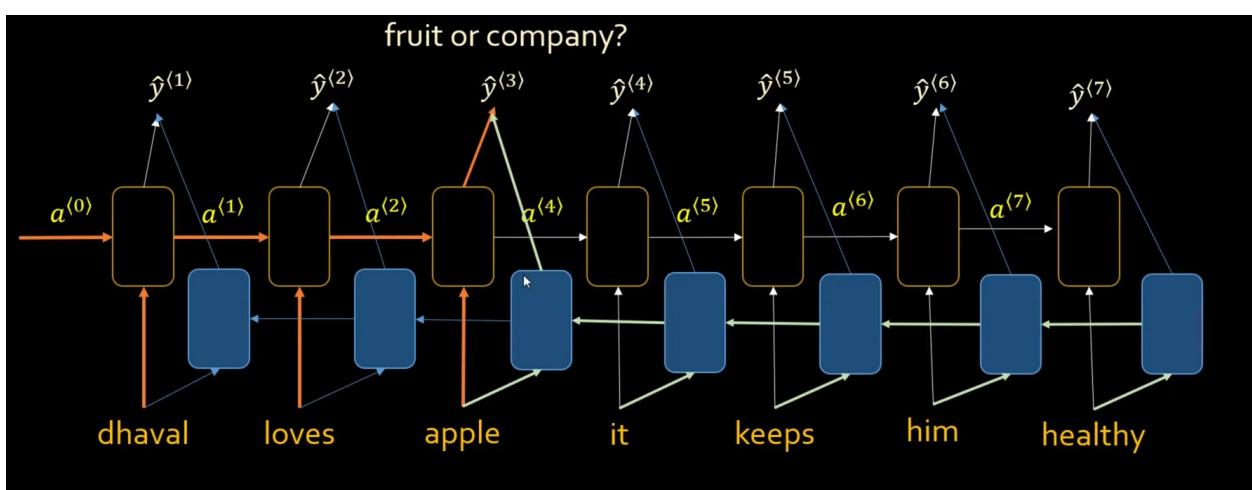
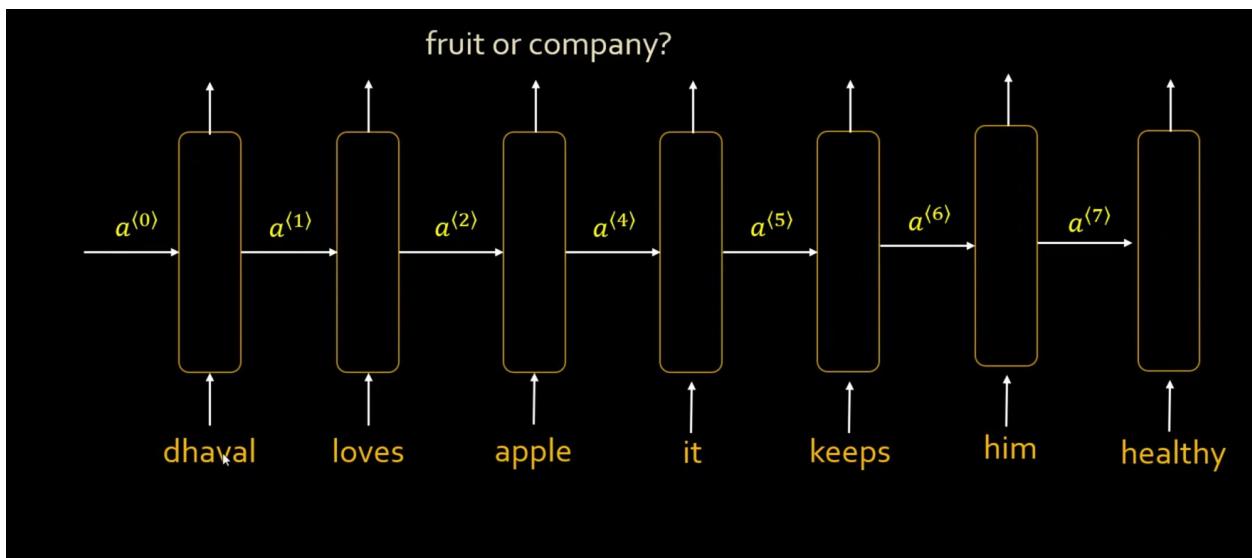
Variables

- x_t : input vector
- h_t : output vector
- \hat{h}_t : candidate activation vector
- z_t : update gate vector
- r_t : reset gate vector
- W, U and b : parameter matrices and vector

LSTM	GRU
3 Gates: Input, output, forget	2 Gates: reset, update
More accurate on longer sequence, less efficient	More efficient computation wise. Getting more popular
Invented: 1995 - 1997	Invented: 2014

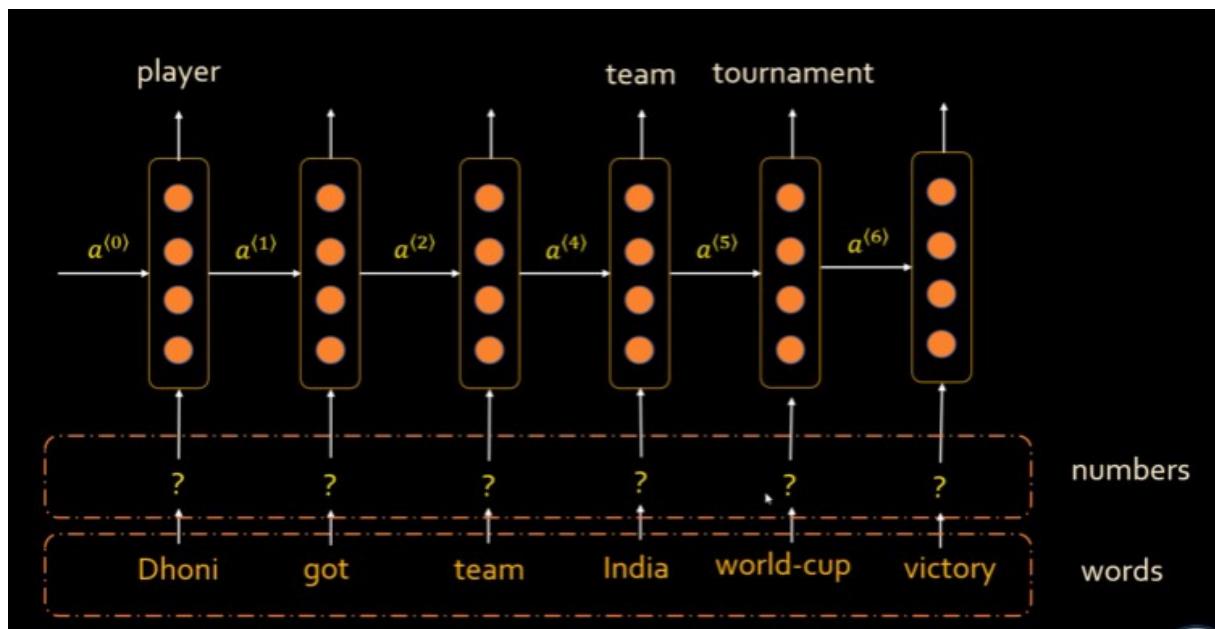
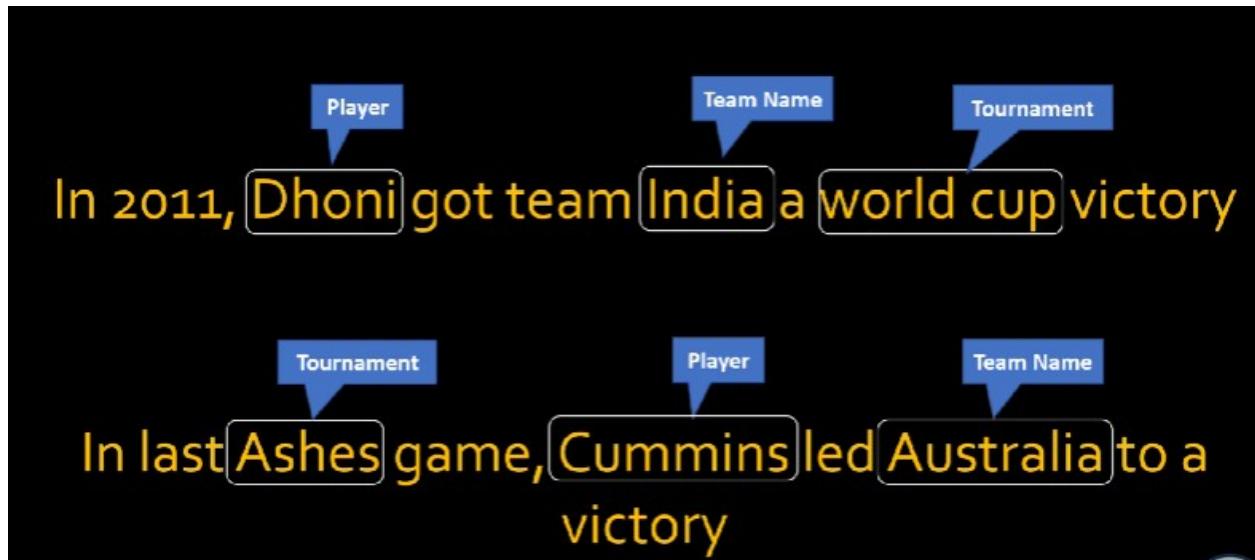
Bi directional RNNs

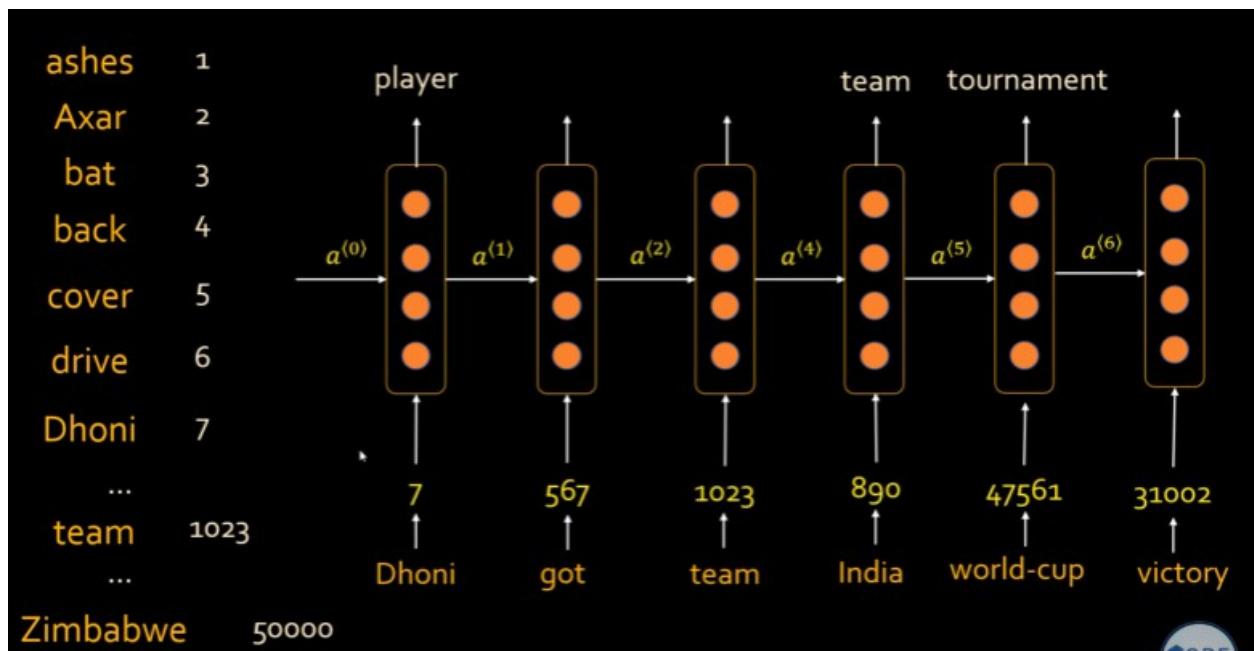
Bi directional RNNs are used in NLP problems where looking at what comes in the sentence after a given word influences final outcome. In this short video we will look at bi directional RNN architecture using a very simple example of named entity recognition.



Machine learning models don't understand words. They should be converted to numbers before they are fed to RNN or any other machine learning model. In this tutorial, we will look into various techniques for converting words to numbers. These techniques are, 1) Using unique numbers 2) One hot encoding 3) Word embeddings

Word embeddings





Issue with option 1: unique numbers

Numbers are random. They don't capture relationship between words

Issue with option 2: one hot encoding

1. Doesn't capture relationship between words
2. Computationally in-efficient





How can we capture similarities between two words?



Dhoni



Cummins



Australia

Person: 1
Healthy/fit: 0.9
Location: 0
Has two eyes: 1
Has government: 0

Person: 1
Healthy/fit: 0.87
Location: 0
Has two eyes: 1
Has government: 0

Person: 0
Healthy/fit: 0.7
Location: 1
Has two eyes: 0
Has government: 1



Dhoni



Cummins



Australia

$$\begin{bmatrix} 1 \\ 0.9 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0.87 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 0.7 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

	ashes	Australia	Bat	Cummins	cover	Dhoni	World cup	..	Zimbabwe
Person	0	0.02	0.1	0.95	0.03	0.96	0.67	...	0.04
Country	0	0.97	0	0	0	0	0	...	1
Healthy & Fit	0	0	0.3	0.87	0	0.9	0	...	0
Event	1	0.1	0	0	0.4	0	1	...	0
gear	0	0	1	0	0	0	0	...	0

	ashes	Australia	Bat	Cummins	cover	Dhoni	World cup	..	Zimbabwe
Person	0	0.02	0.1	0.95	0.03	0.96	0.67	...	0.04
Country	0	0.97	0	0	0	0	0	...	1
Healthy & Fit	0	0	0.3	0.87	0	0.9	0	...	0
Event	1	0.1	0	0	0.4	0	1	...	0
gear	0	0	1	0	0	0	0	...	0



Converting words to numbers

1. Unique numbers
2. One hot encoding
3. Word embeddings

1. TF-IDF

2. Word2Vec

Embeddings are not hand
crafted. Instead, they are
learnt during neural network
training

Techniques to compute word embeddings

1. Using **supervised** learning
2. Using **self-supervised** learning
 1. Word2vec
 2. Glove

1. Using **supervised** learning

Take an NLP problem and try to solve it.
In that pursuit as a side effect, you get
word embeddings



nice food. The pasta dish was too good!



poor quality food. I would never go there again!



After food nice ... poor ... zonal

1

2

3

5000

$$\begin{bmatrix} 0.5 \\ 1.2 \\ 0.7 \\ 3.1 \end{bmatrix}$$

$$\begin{bmatrix} 4.3 \\ 0.1 \\ 0.9 \\ 5.5 \end{bmatrix}$$

$$\begin{bmatrix} 0.4 \\ 8.1 \\ 8.8 \\ 4.2 \end{bmatrix}$$

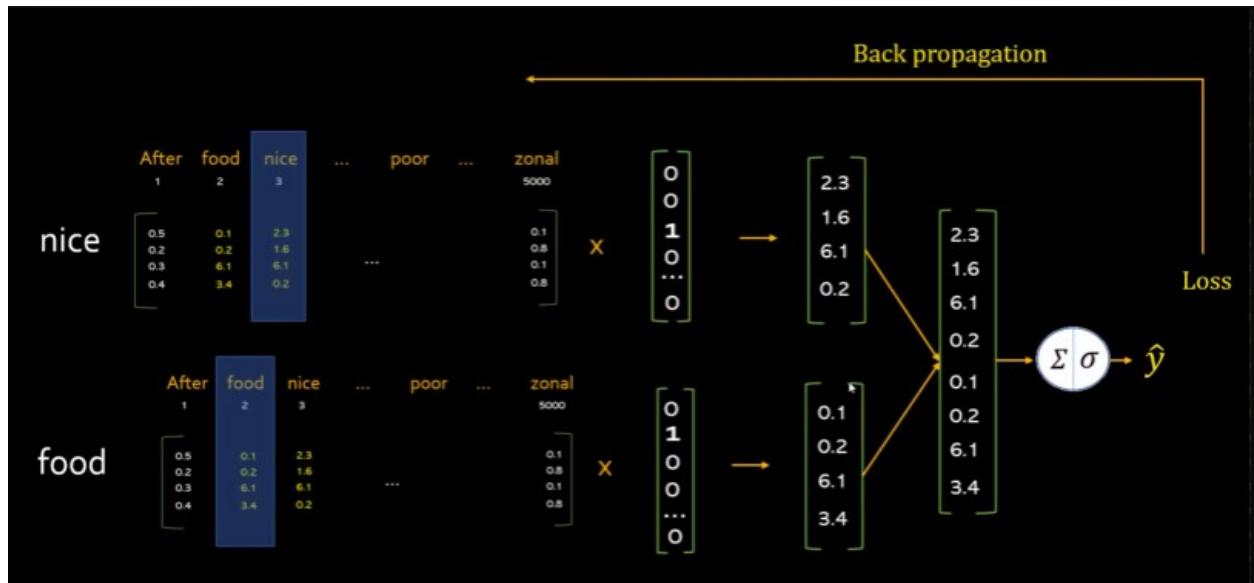
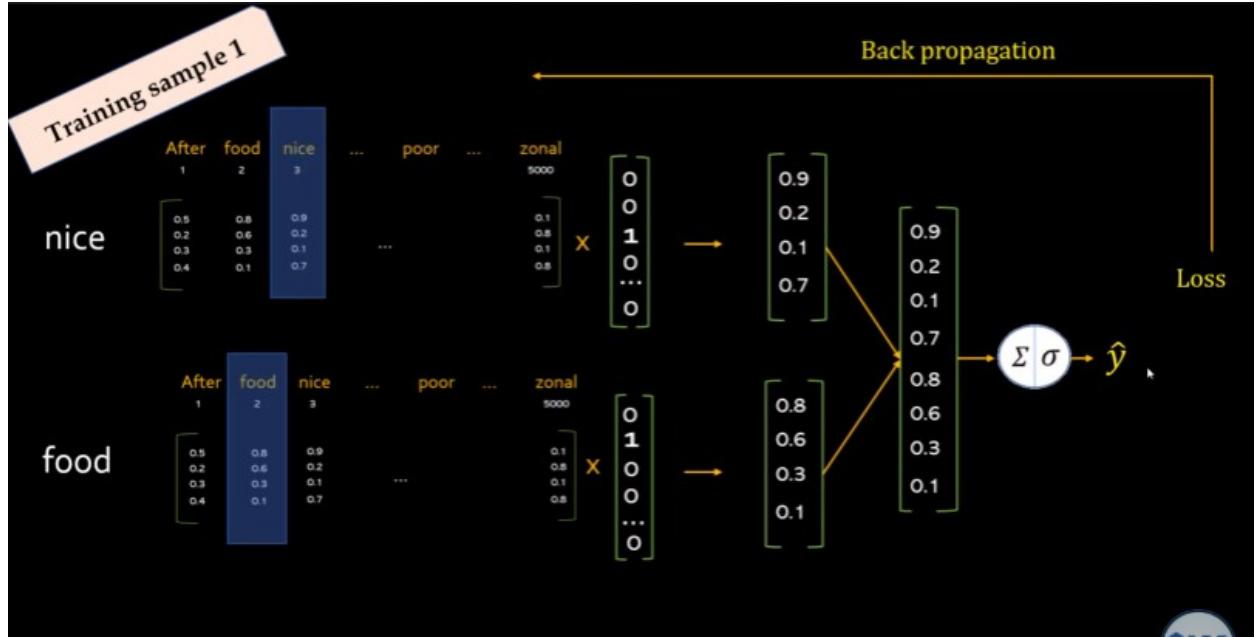
$$\begin{bmatrix} 9.8 \\ 0.6 \\ 2.2 \\ 1.3 \end{bmatrix}$$

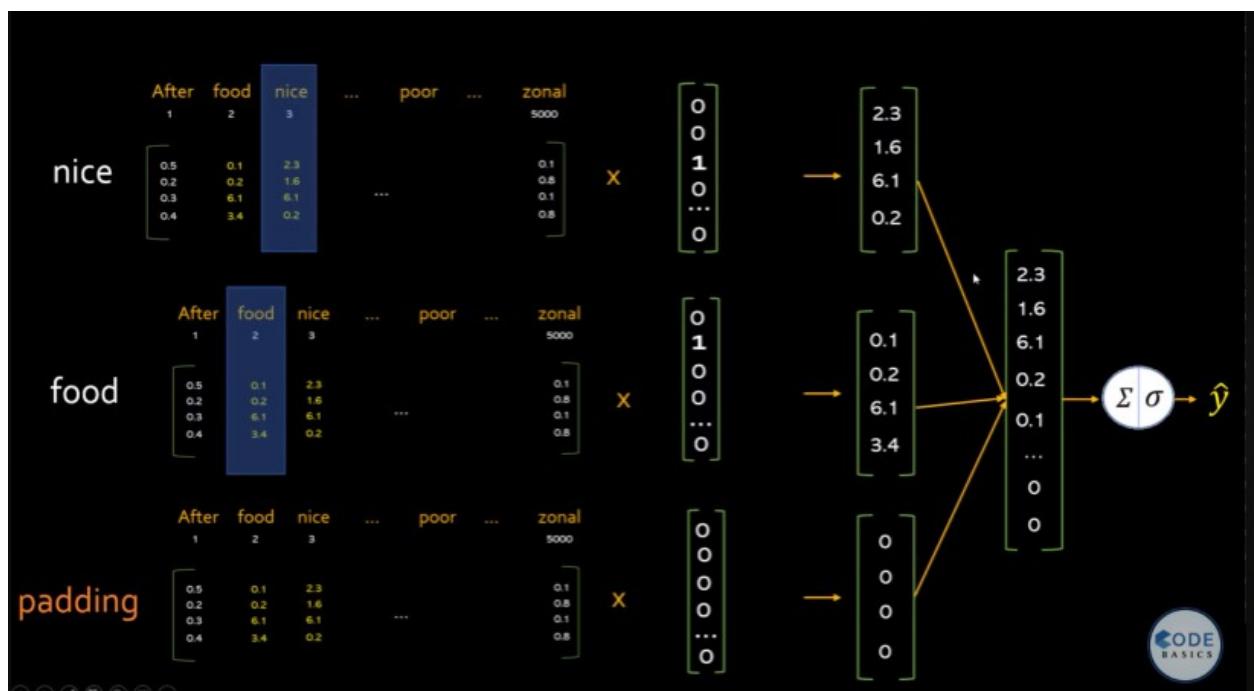
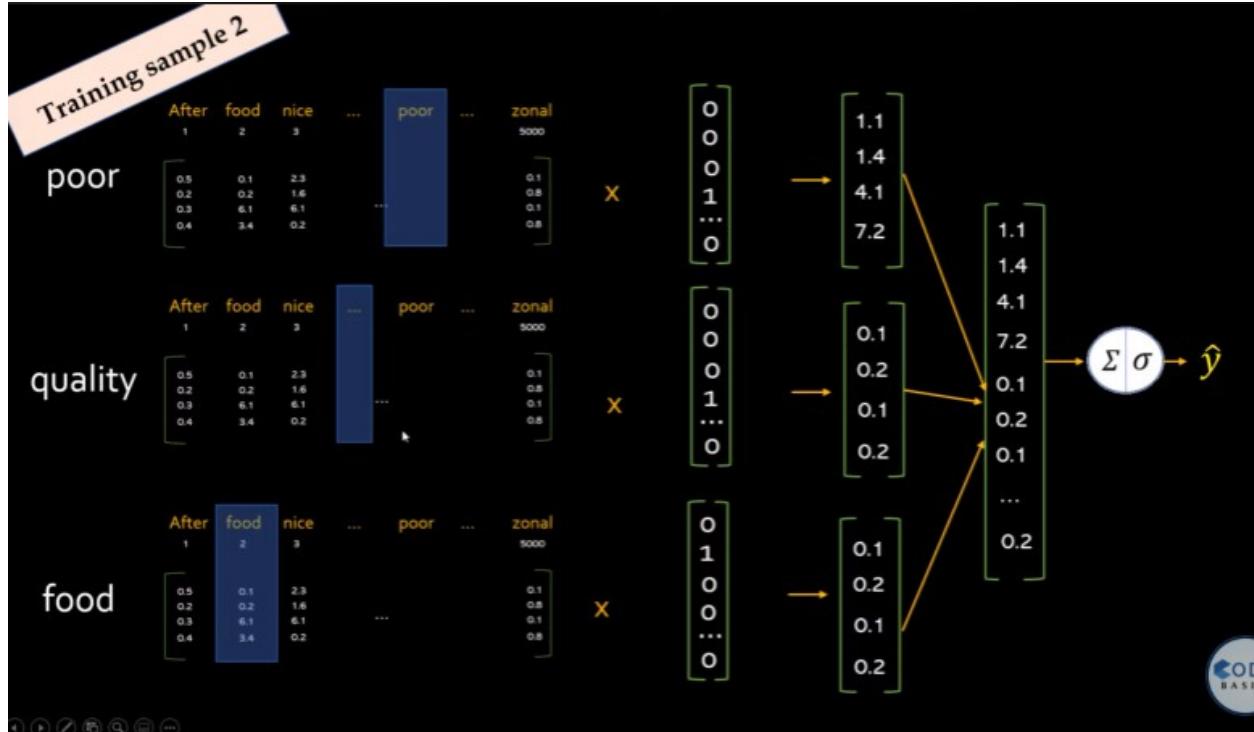
4-dimensional word embedding vector

After	food	nice	...	poor	...	zonal
1	2	3				5000
0.5	4.3	0.4				9.8
1.2	0.1	8.1				0.6
0.7	0.9	8.8	...			2.2
3.1	5.5	4.2				1.3

4 by 5000 Embedding Matrix: E

	After	food	nice	...	poor	...	zonal	5000
After	1	0	0		0		0	
food	0	1	0		0		0	
nice	0	0	1		0		0	
...								
poor	0	0	0		1		0	
...								
zonal	0	0	0		0		1	





After food nice good poor weak ... zonal
5000

0.5	4.3	0.4	0.38	2.3	2.2	9.8
1.2	0.1	8.1	8.2	1.1	1.0	0.6
0.7	0.9	8.8	8.9	6.7	6.5	2.2
3.1	5.5	4.2	4.1	2.0	1.9	1.3

4 by 5000 Embedding Matrix: E

Similar words have almost similar vectors.