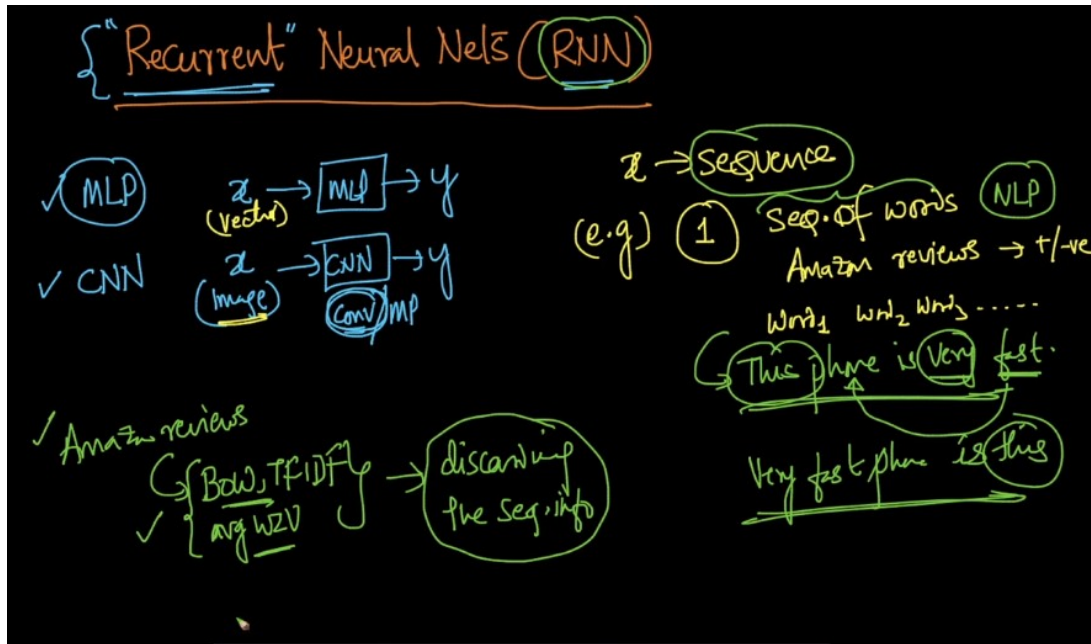Recurrent Neural networks

Sequence matters for making a meaningful sentence.
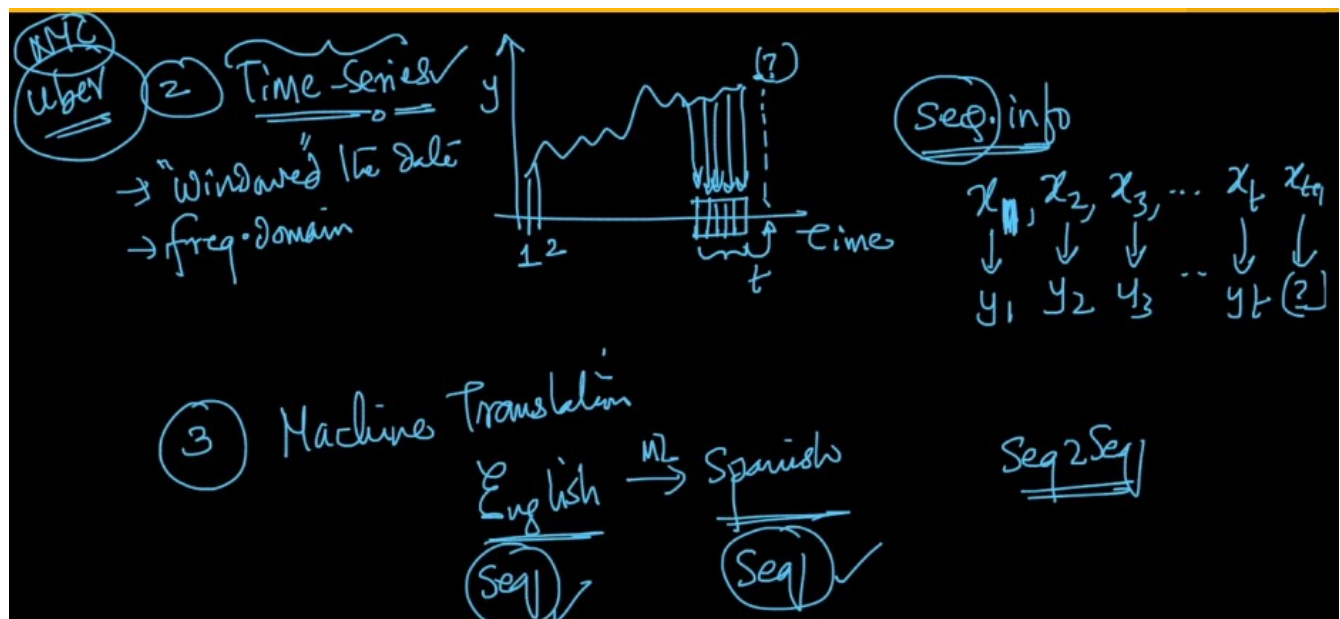Types of data various NN can handle:

In BOW,TFIDF,Aavg w2v ignores the sequence information. This is actually important to make most out of the sequence of words.



In case of time series data, we windowed the data and we try to predict the next point. Time series is also a sequence information.
We partially used the sequence information using the windowing technique.
Machine translation is the sequence information, sequence input and sequence output.

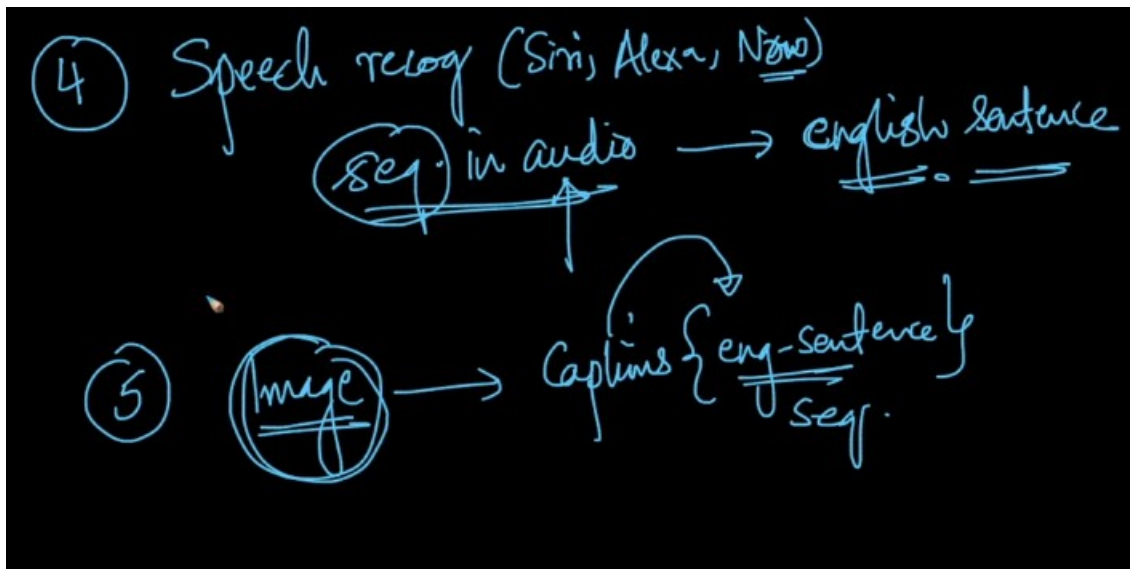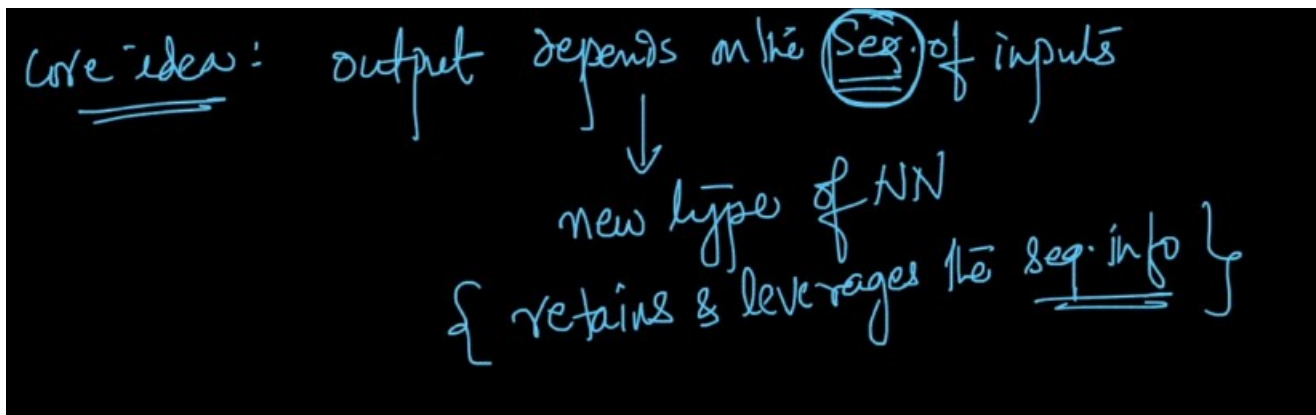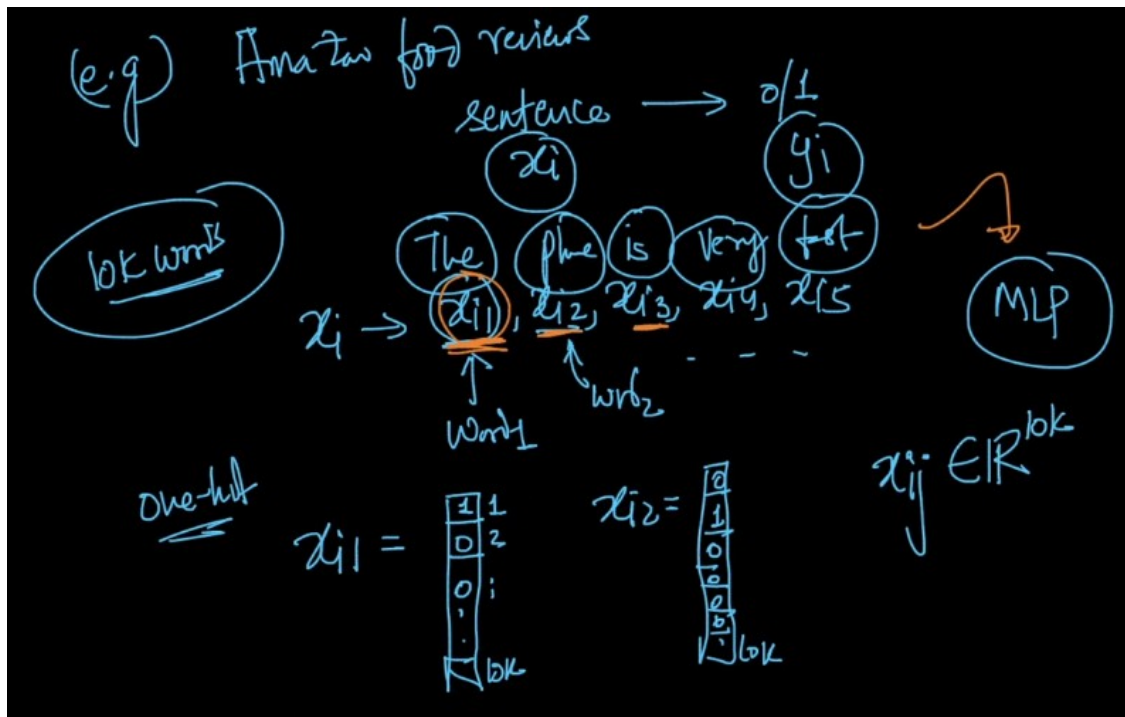Speech recognition and Image captions:



Image captioning can be used in google image search, google text search.
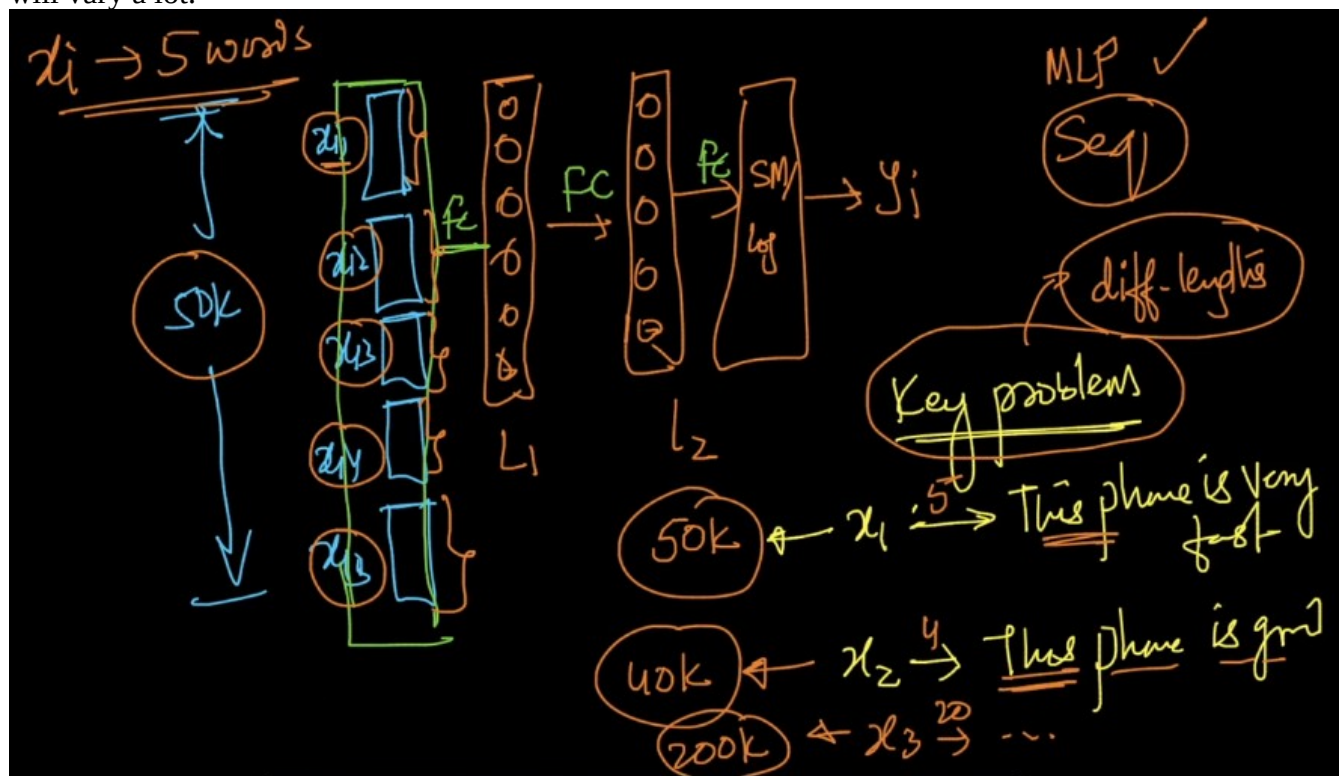
The importance of RNNs:

Representing the words in the vector.
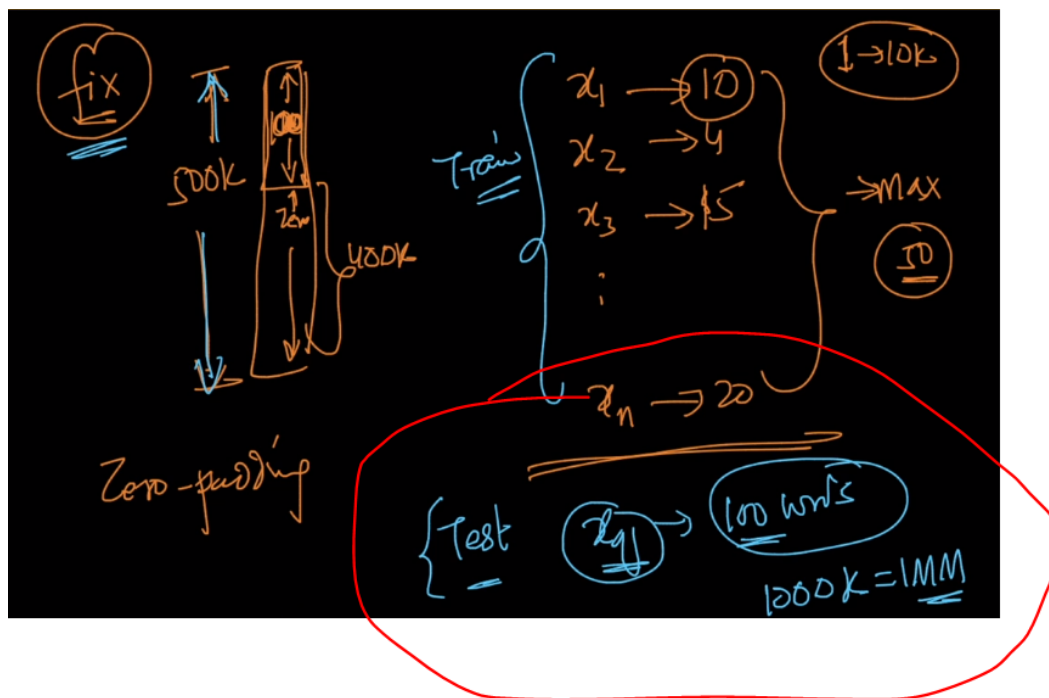


Why not we train MLP ?
Process of training MLP on the sequence data. Every word is one-hot encoded.
Xi has 5 words and there are 10k unique words in the data set. The number of words in the sentence will vary a lot.
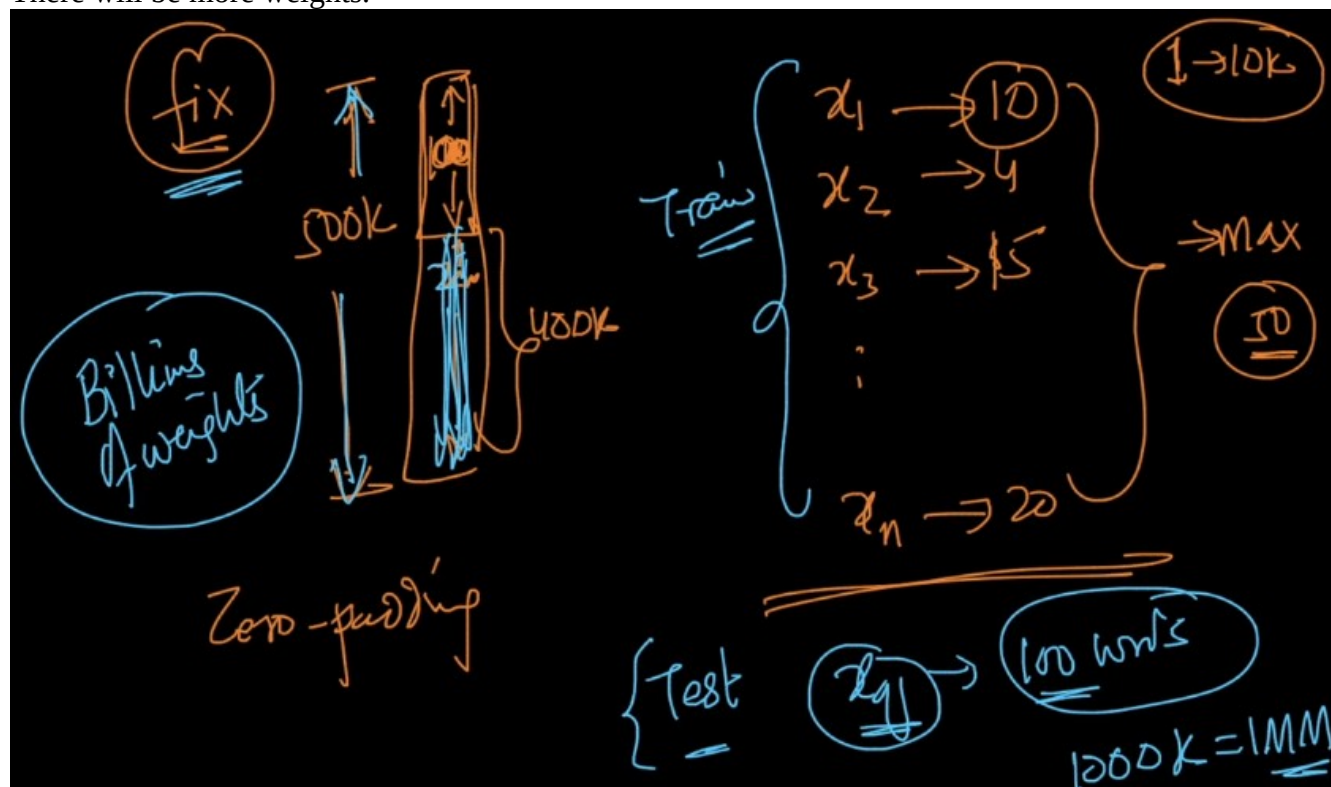
We can fix the variable length of the sentences, We will pad the sentence with zeros by taking the maximum length of the sentence.

Even with the fixed length, there will bw problem in case of the text data if we have the large dimensions.
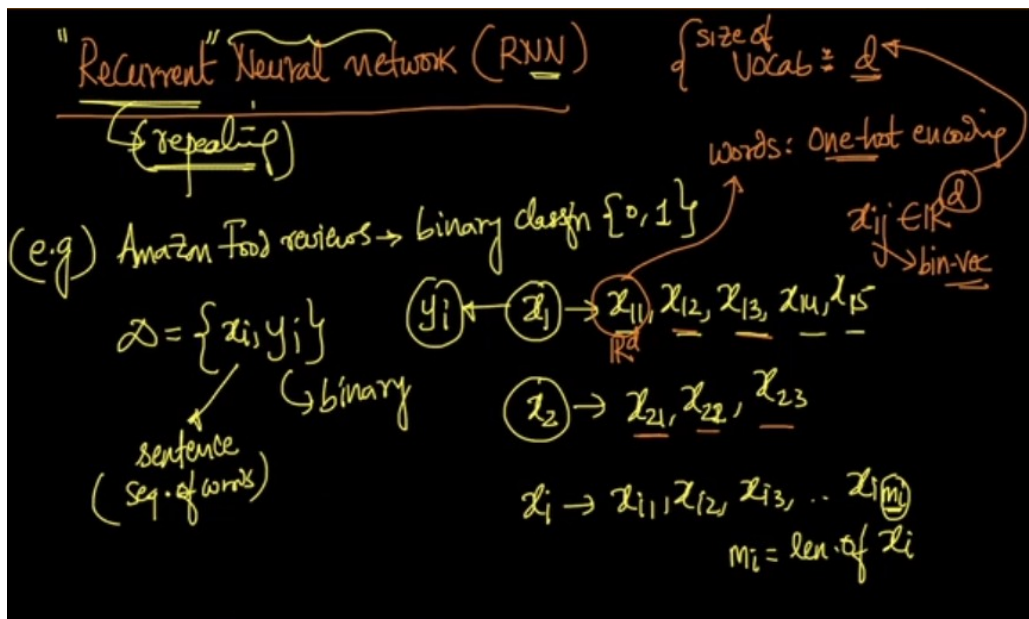


The input size will become massive. We are building un necessarily large vectors.
There will be more weights.

1We can use all the word vectors for RNN's

Amazon fine food reviews is the binary classification task (0, 1).

We are representing the words in the one hot encoding, let say the size of vocab has the d words and each word belongs to a sentence.

Task: To classify the seq of words into two classes.
In the first stage we take the first word and then second word and so on.

The same weight for the first word hidden layer is used for the next word and the out put of the first word.
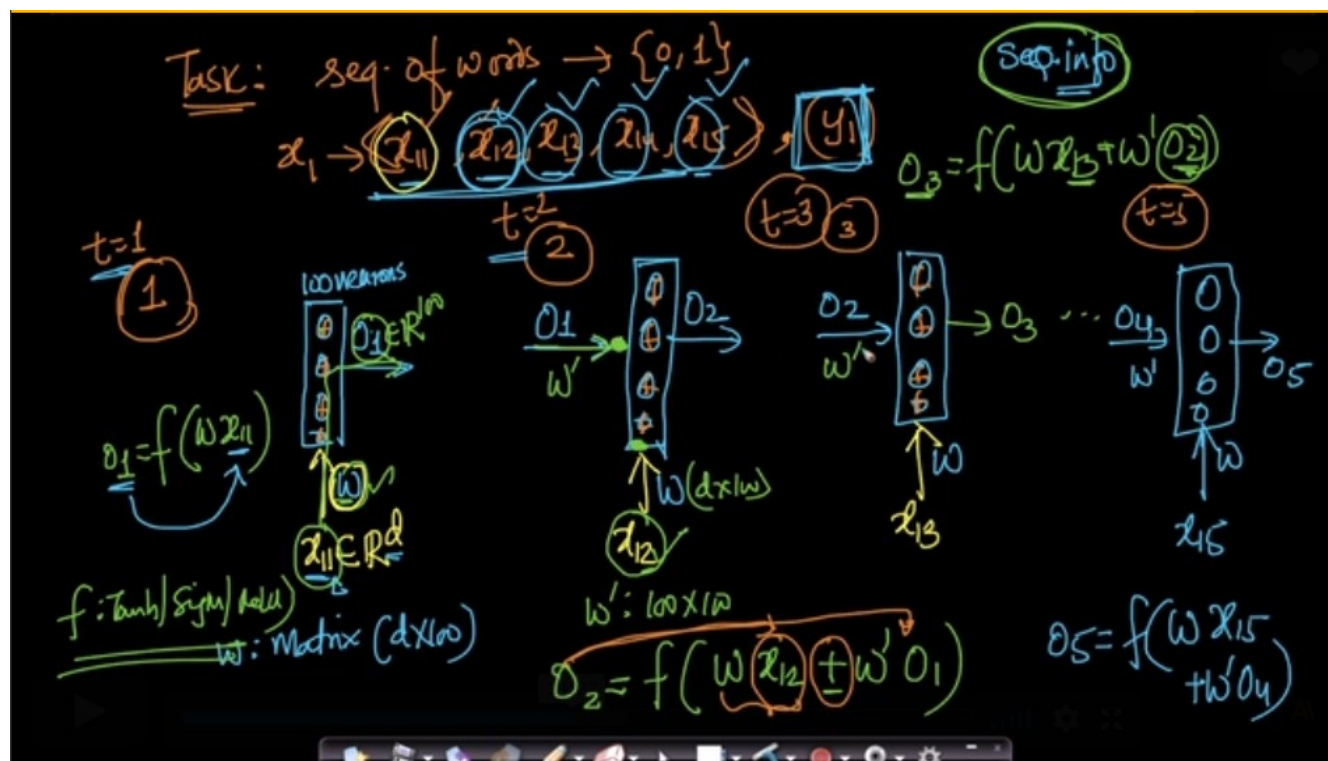
We want to keep the seq information. So, we use the output of the first layer and pass it to the second layer with also a next word in the sequence.

Recurrent means repeating.

Explanation of recurrent NN:
Each word in the sentence is represented as one hot encoding vectors. They are binary vectors.
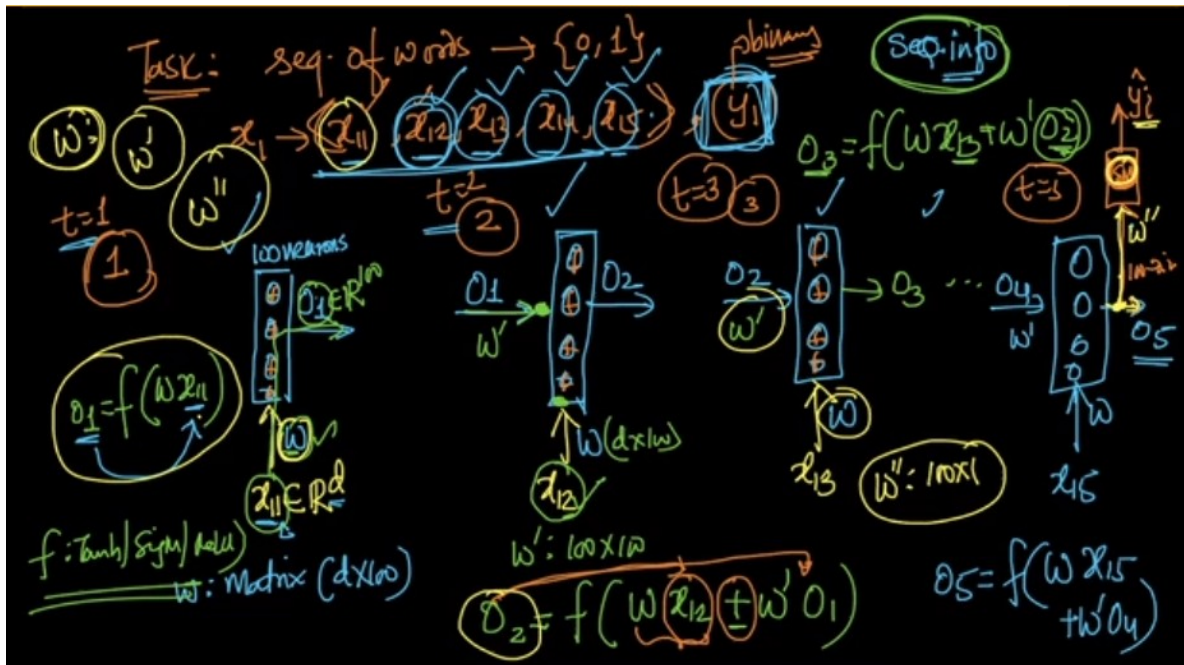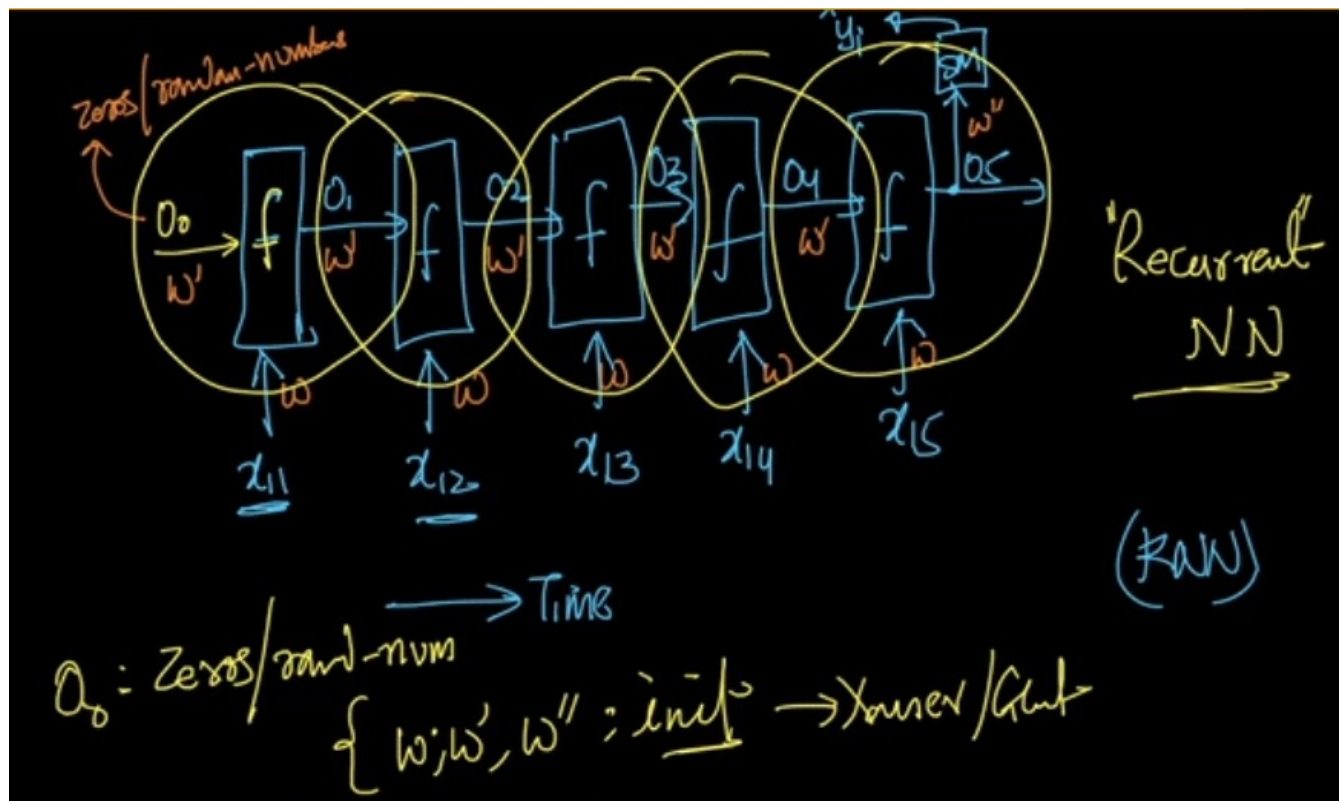
Training the sentence:

Determining the label yi of the sequence of words:
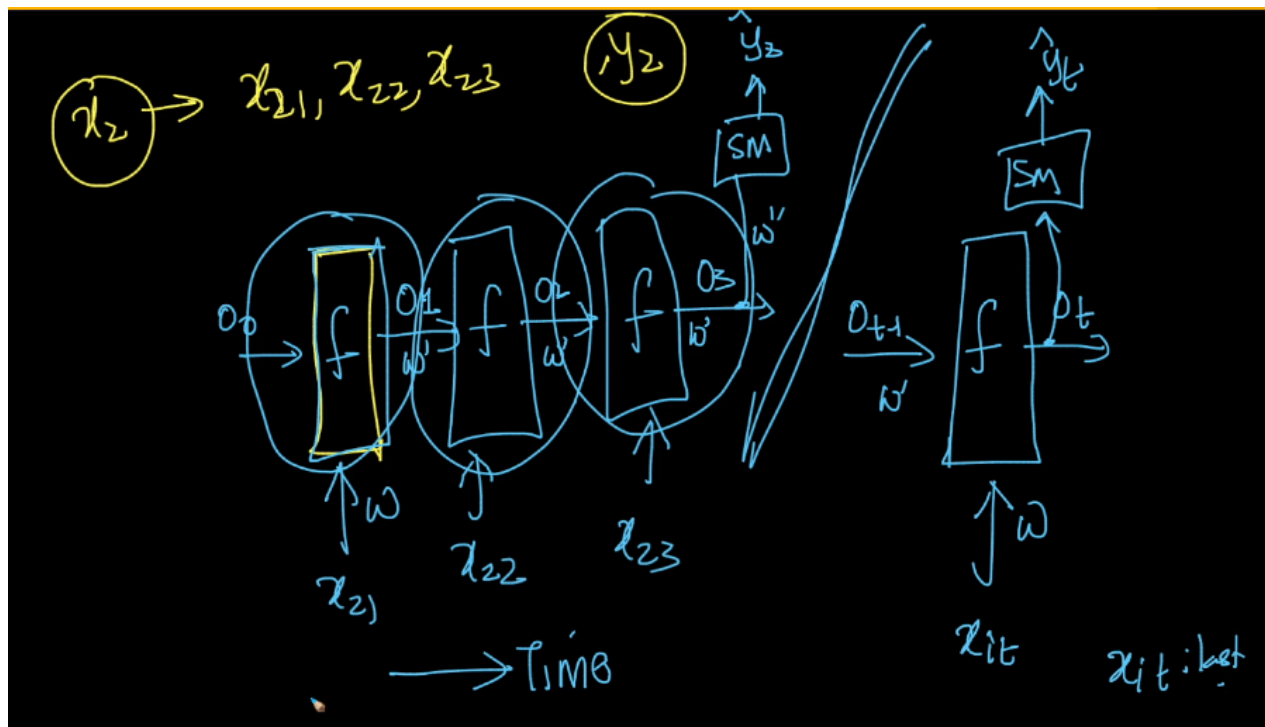We take the output of the last layer and we put a softmax layer to give the label yi.

We will have the weight vector at the end, this is used for determining the label of the sequence.
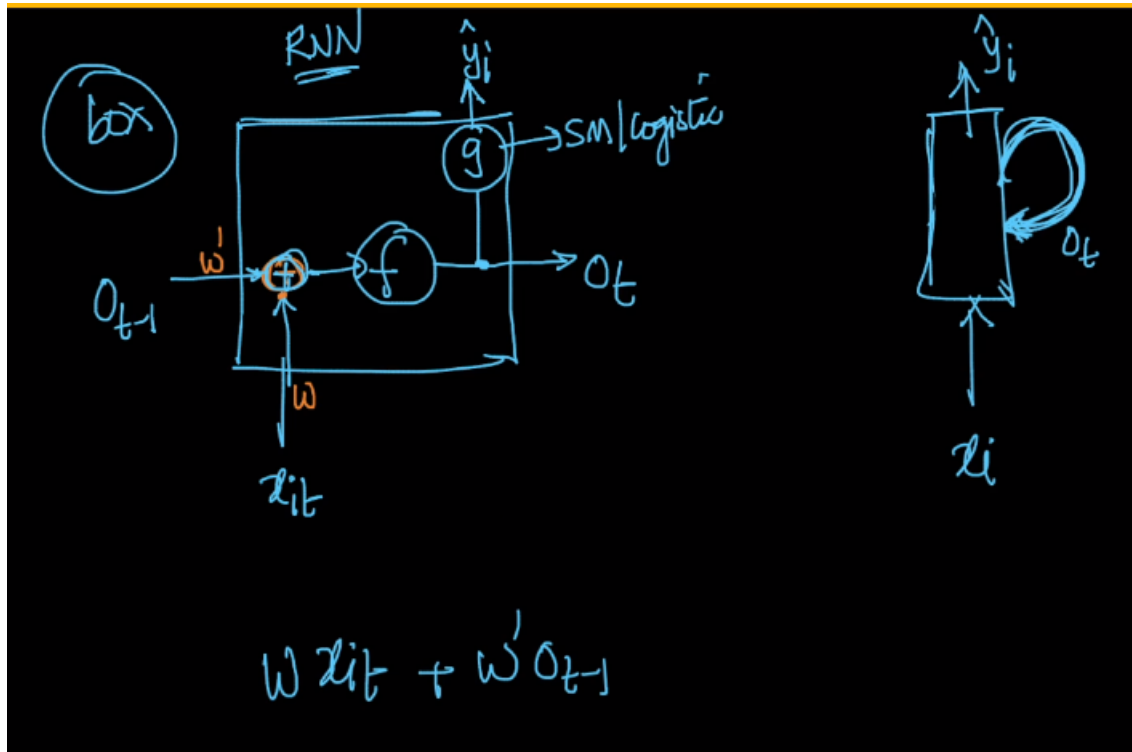


Simplified notation of RNN:

Training with the next sequence of words:
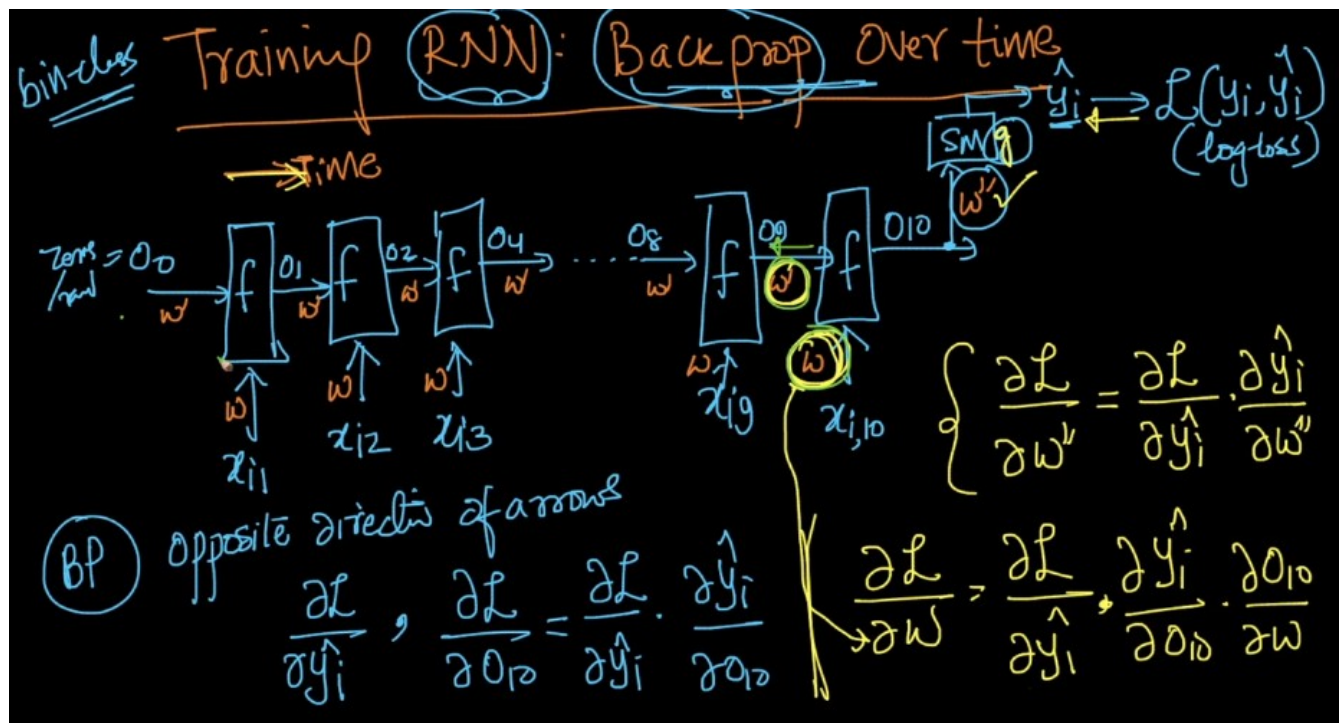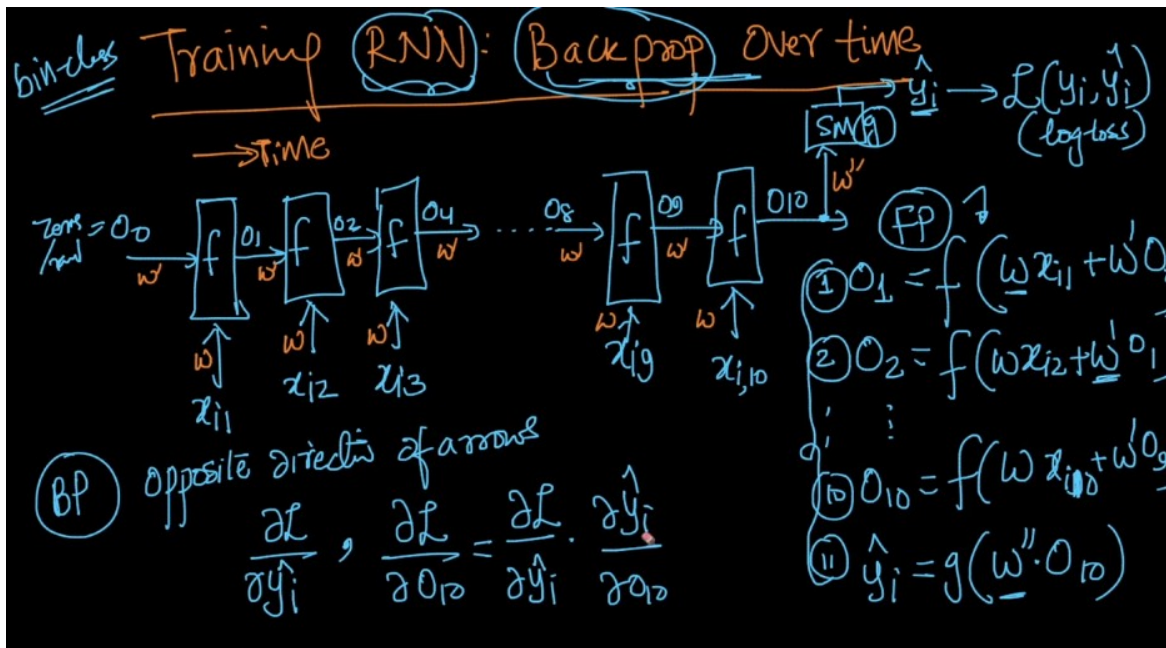


The other way of representing the RNN's:
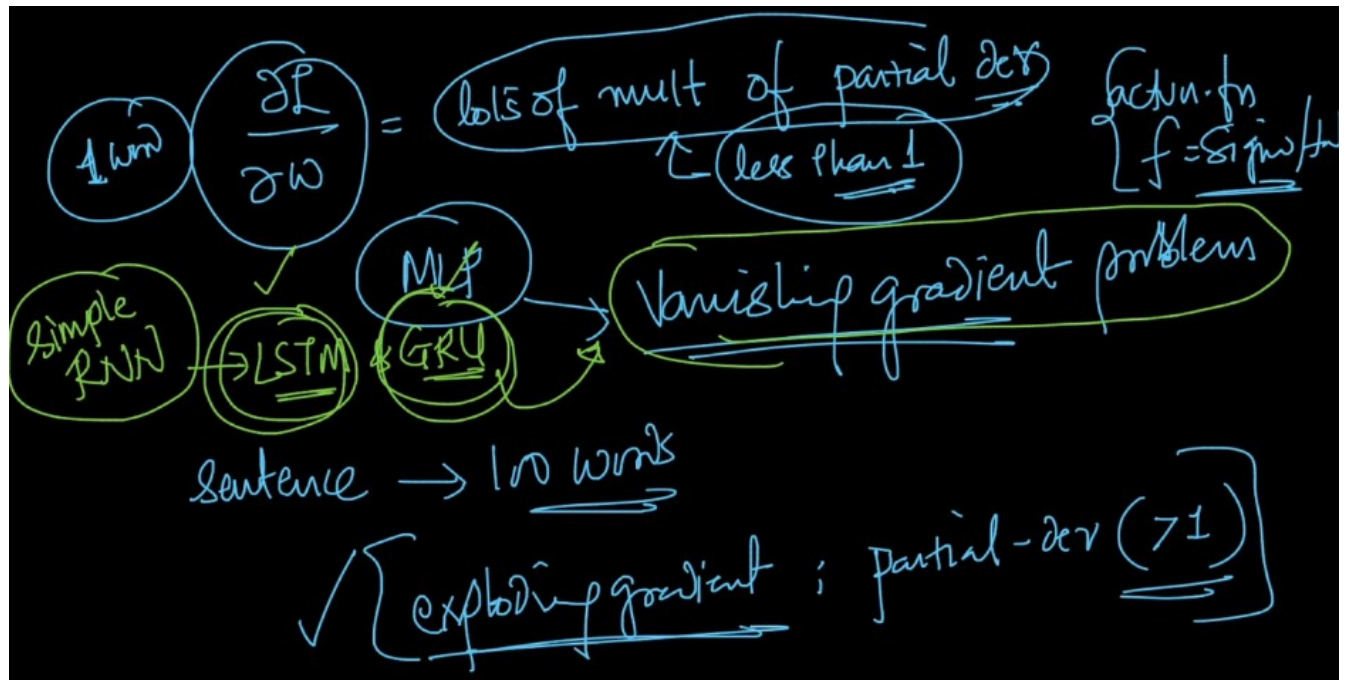
Training RNN's using back propagation:
This is a modification over the normal back propagation. Forward prop use the direction of arrows, to make training.
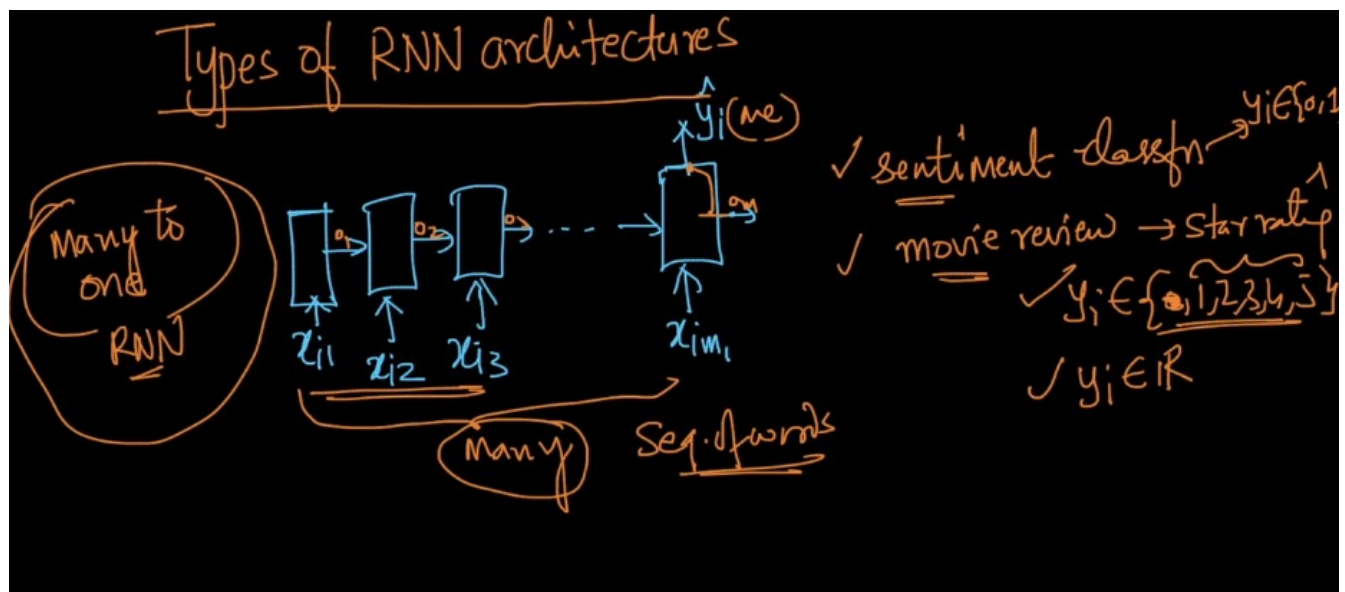Backward prop of the opposite of the arrows.

This is over time and we have the three derivatives w, w1, w11.

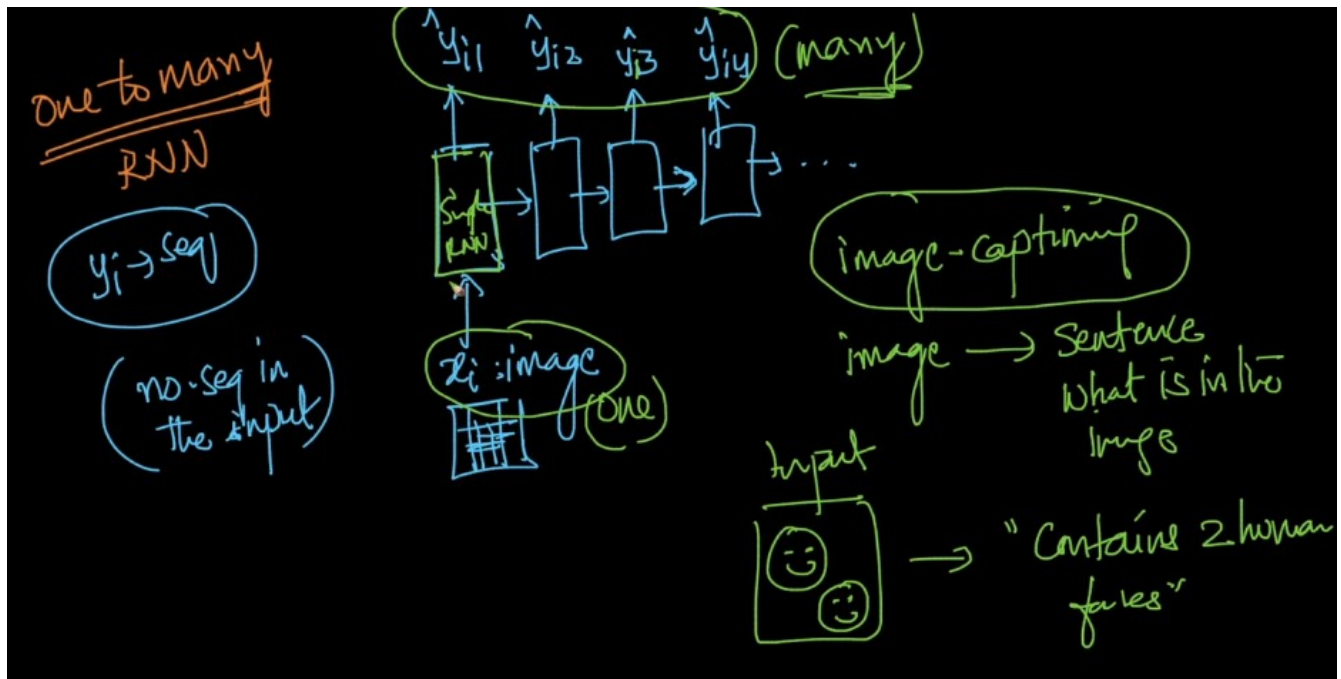Due to the Vanishing gradient problems the LSTMs and GRUs are used to avoid this.


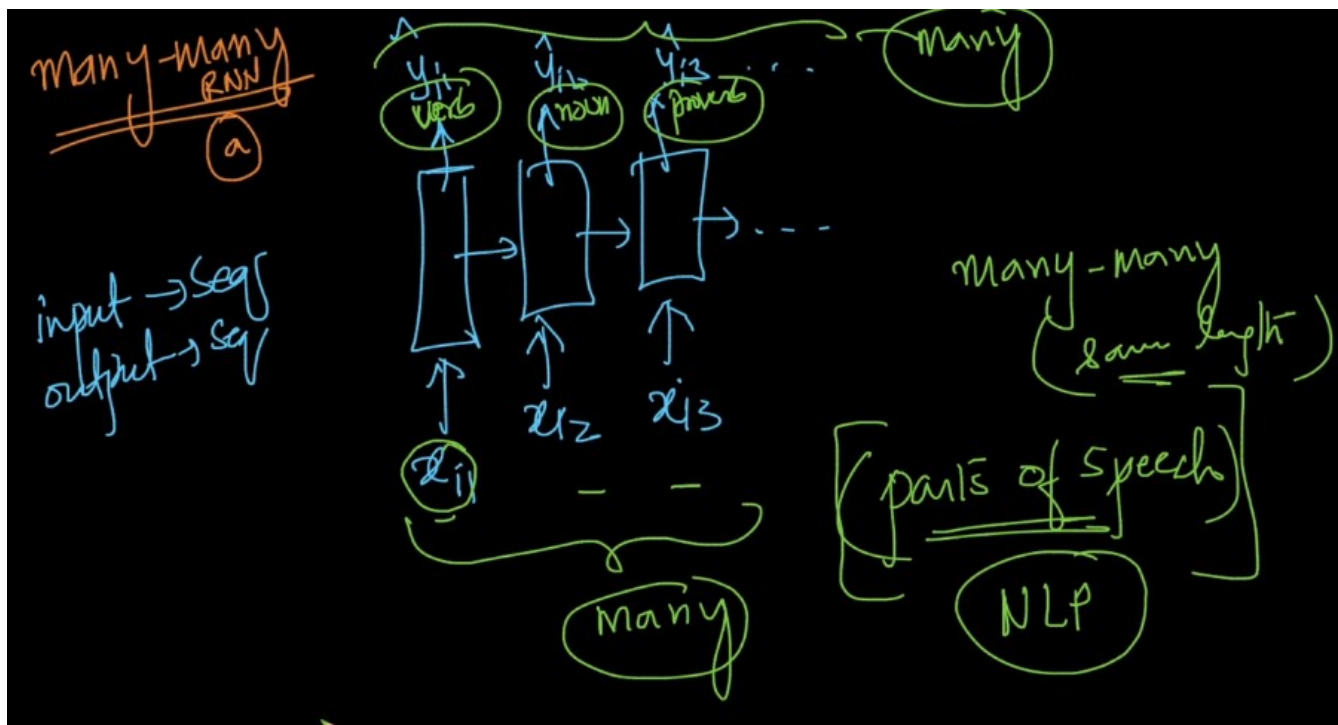
Types of RNN's:

Many – one Rnns:
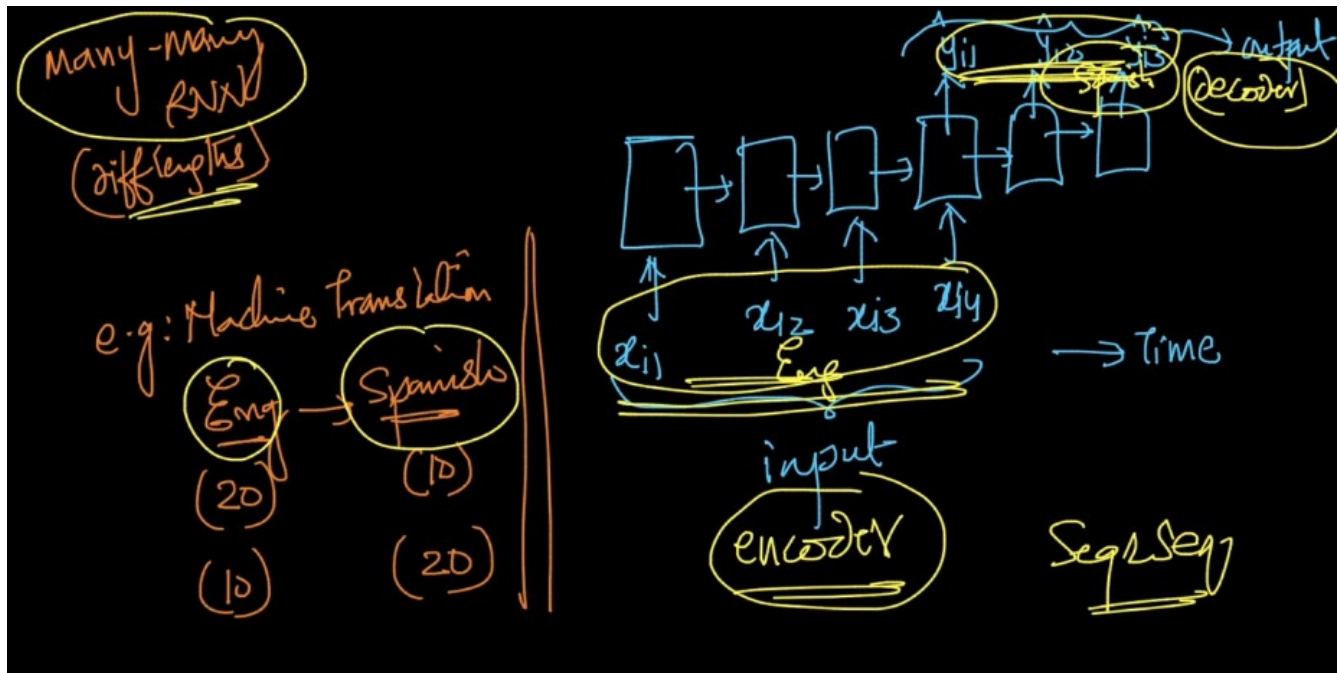
One to many:



Many to Many:
Same length.
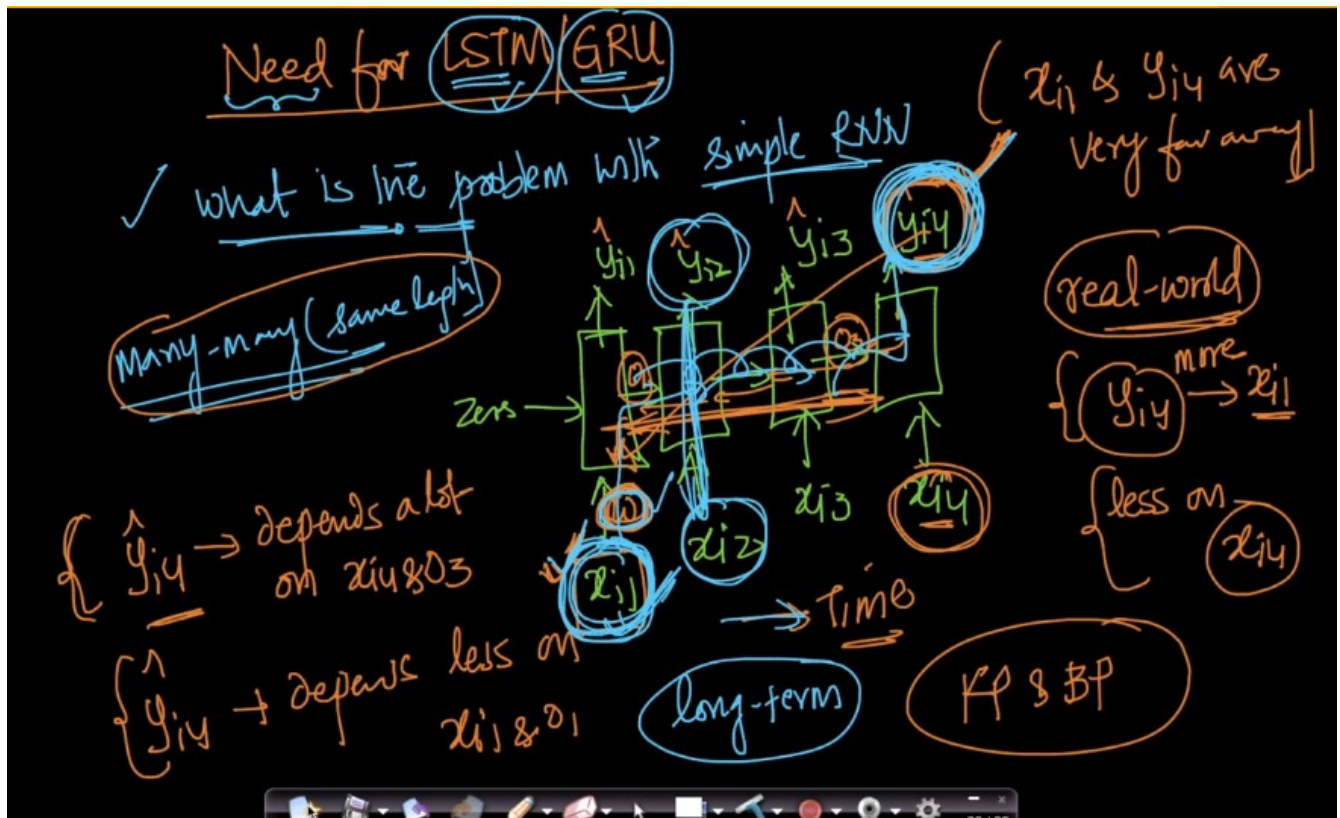For determining the parts of speech of a given sequence.
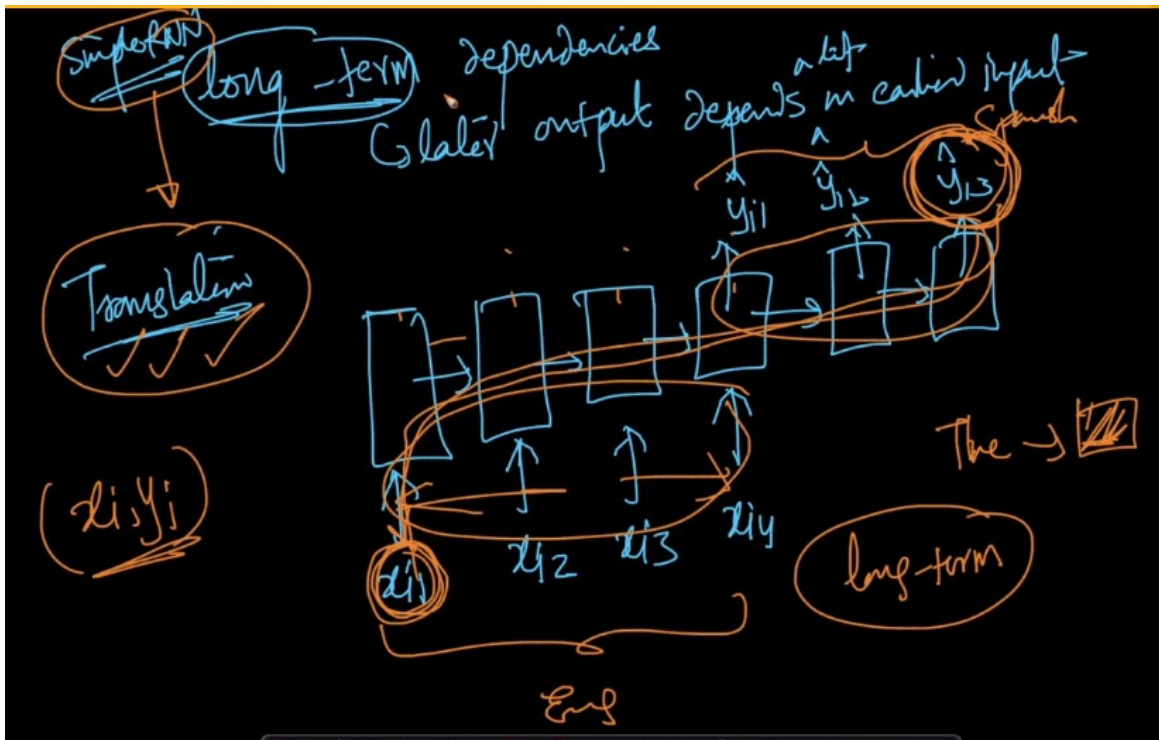
Different length.

Machine Translation is example of different lengths. This is also called the Encoder – Decoder Rnn's.
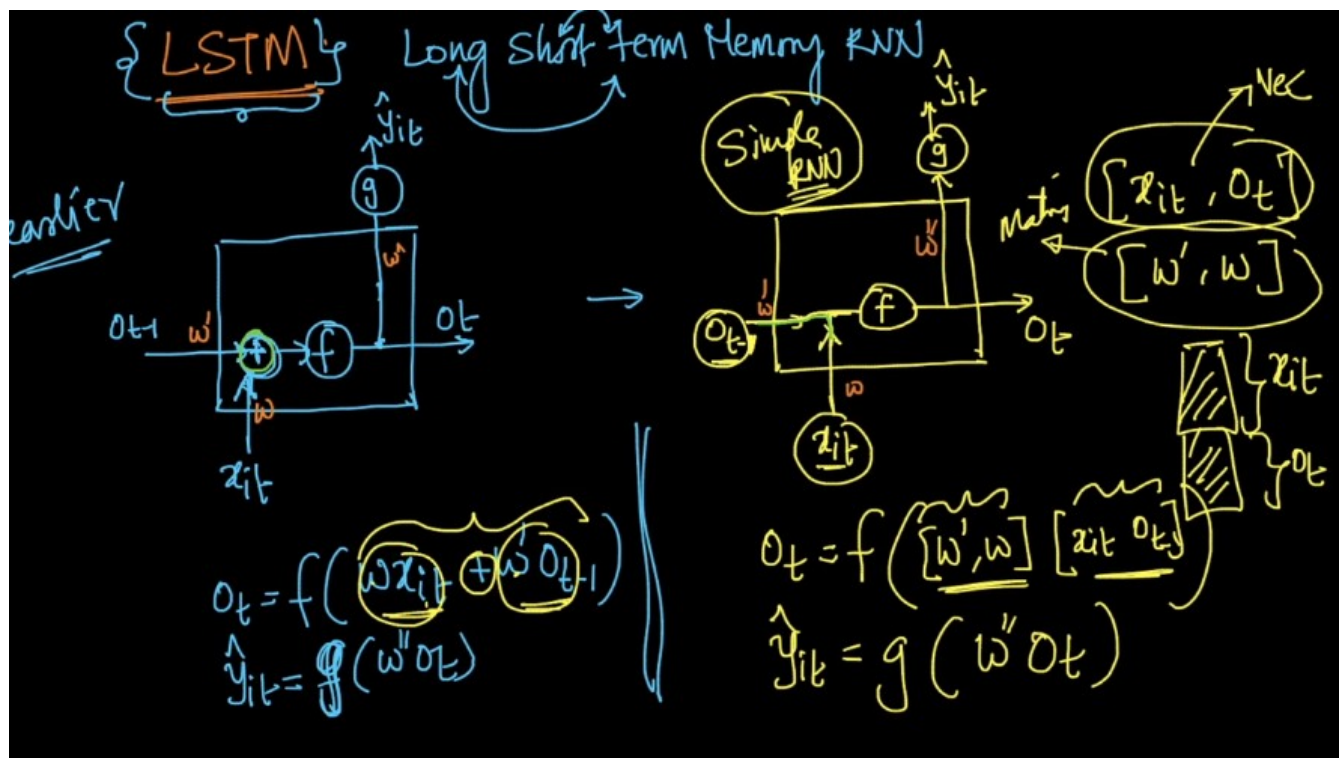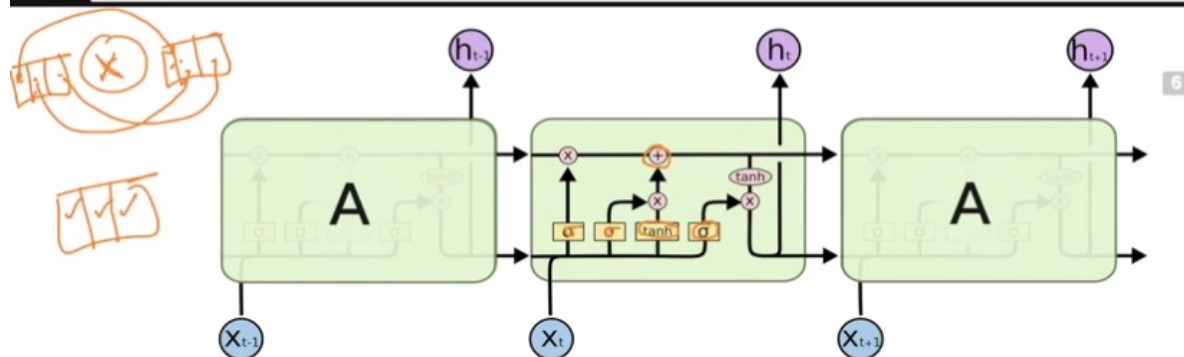


One – one is also present.
Need for LSTM/GRU:

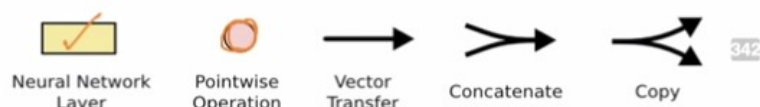LSTM's: In lstms we concatenate instead of summation.

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



The repeating module in an LSTM contains four interacting layers.

Don't worry about the details of what's going on. We'll walk through the LSTM diagram step by step later. For now, let's just try to get comfortable with the notation we'll be using.

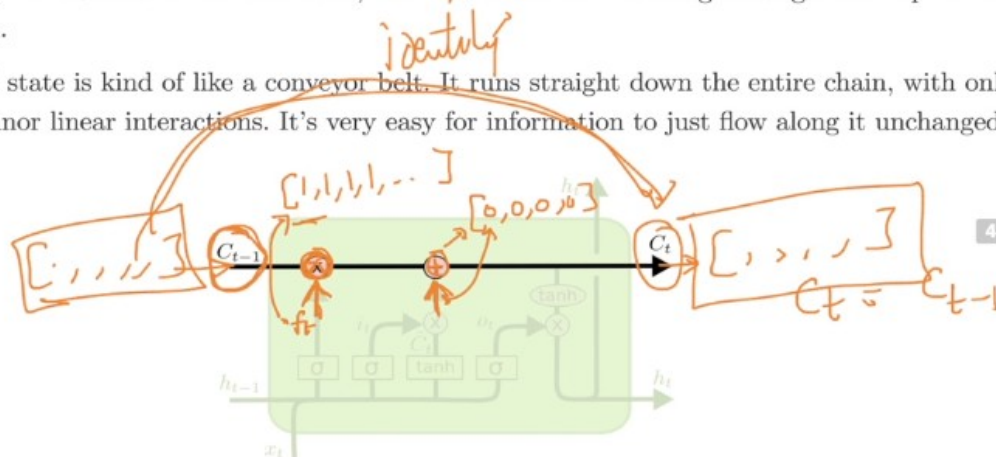| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

First path:



## The Core Idea Behind LSTMs

The key to LSTMs is the cell state, the horizontal line running through the top of the diagram.

The cell state is kind of like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions. It's very easy for information to just flow along it unchanged.
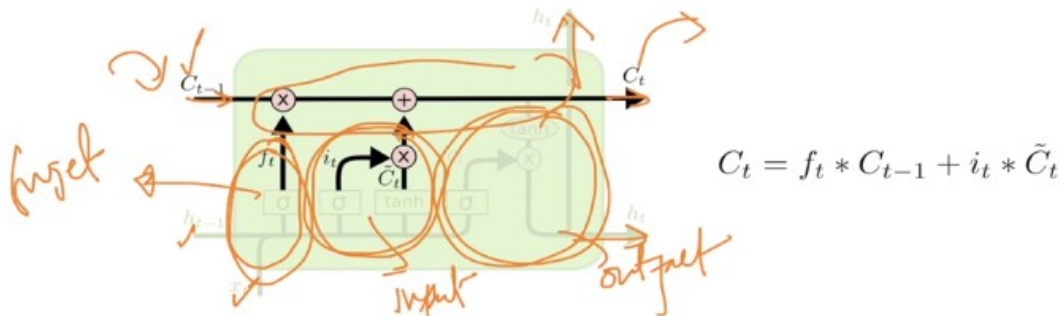
The LSTM does have the ability to remove or add information to the cell state, carefully regulated by structures called gates.

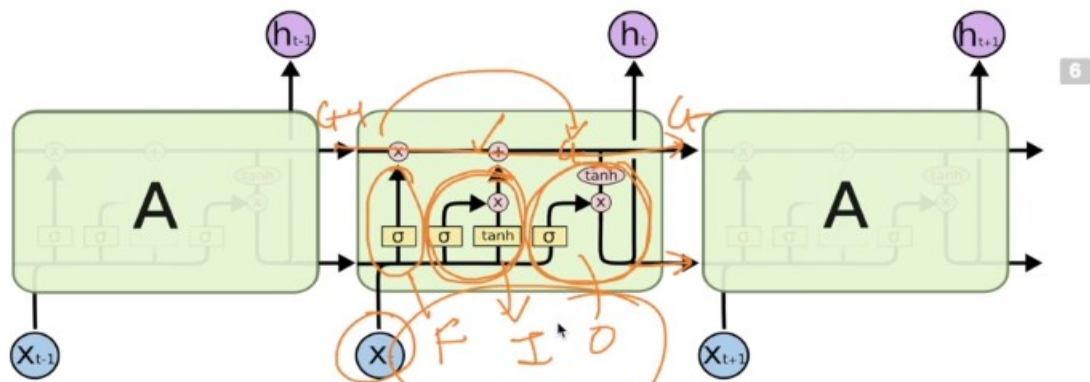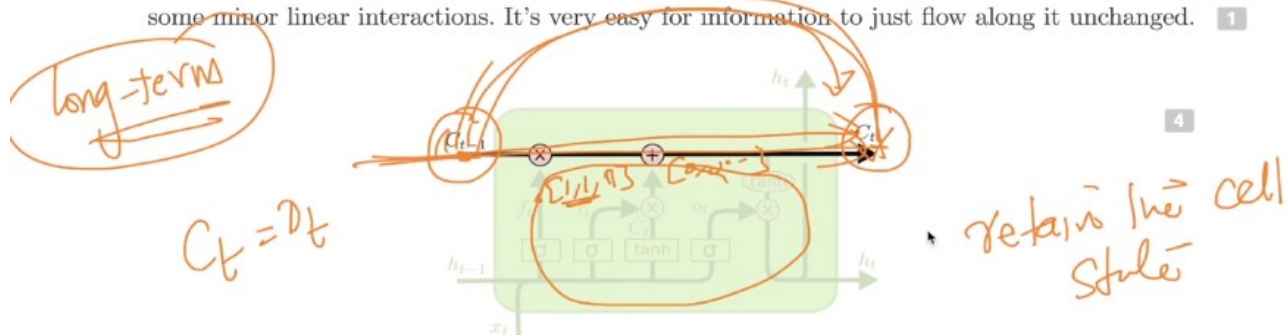Gates are a way to optionally let information through. They are composed out of a sigmoid neural net

We multiply the old state by $f_t$, forgetting the things we decided to forget earlier. Then we add $i_t * \tilde{C}_t$. This is the new candidate values, scaled by how much we decided to update each state value.

In the case of the language model, this is where we'd actually drop the information about the old subject's gender and add the new information, as we decided in the previous steps.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Finally, we need to decide what we're going to output. This output will be based on our cell state, but will be a filtered version. First, we run a sigmoid layer which decides what parts of the cell state we're going to output. Then, we put the cell state through tanh (to push the

special way.



**The repeating module in an LSTM contains four interacting layers.**

Don't worry about the details of what's going on. We'll walk through the LSTM diagram step by step later. For now, let's just try to get comfortable with the notation we'll be using.

We have the bypassing the input, by this the long term dependency is achieved.
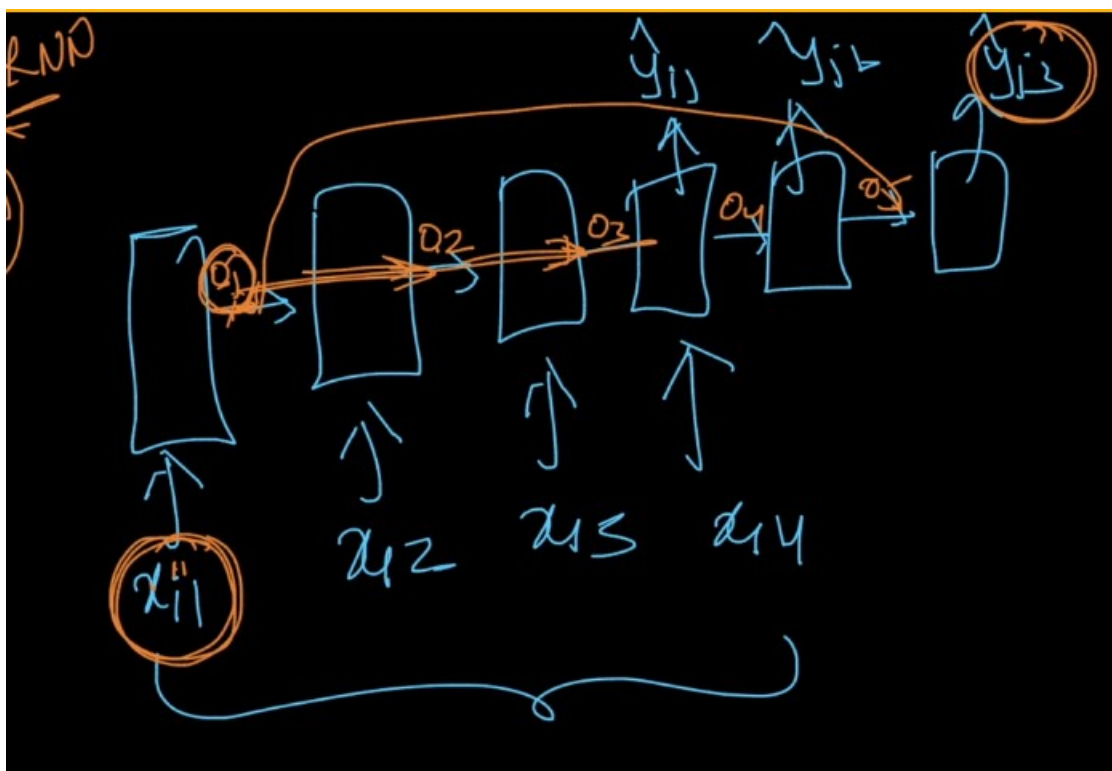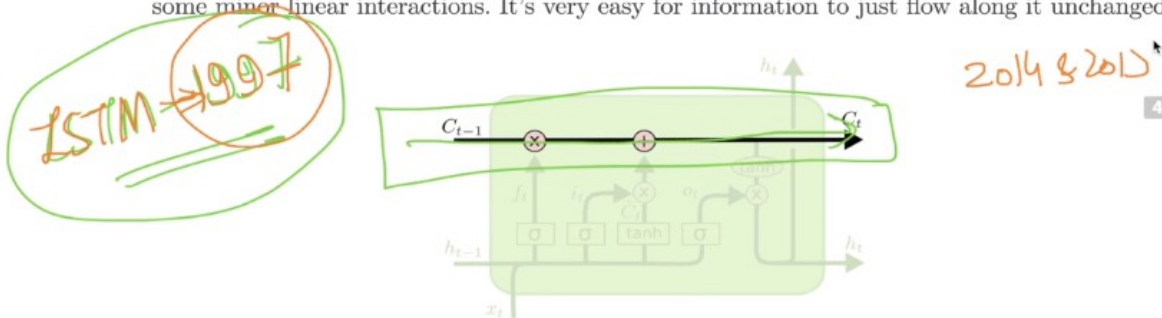
If the last layer needs the first layers output information, we will run into vanishing gradient problem. That is why we use the LSTM network to bypass some of the information. By this long term dependency worked out.

The key to LSTMs is the cell state, the horizontal line running through the top of the diagram.

The cell state is kind of like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions. It's very easy for information to just flow along it unchanged. [1]

The LSTM does have the ability to remove or add information to the cell state, carefully regulated by structures called gates. [1]

**By short circuiting the RNN network, we can impact something which is far away.**



The key to LSTMs is the cell state, the horizontal line running through the top of the diagram.

The cell state is kind of like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions. It's very easy for information to just flow along it unchanged.

The LSTM does have the ability to remove or add information to the cell state, carefully regulated by structures called gates.

Gates are a way to optionally let information through. They are composed out of a sigmoid

GRU's:
These are more modern recurrent NN.

Most often we use the set of LSTMs or GRU cells than simple RNN's.

Bidirectional RNN's:



These are more common in NLP tasks.

When the sequences are of different lengths, it is hard to learn the weights because, we have the arbitary length of sequences.
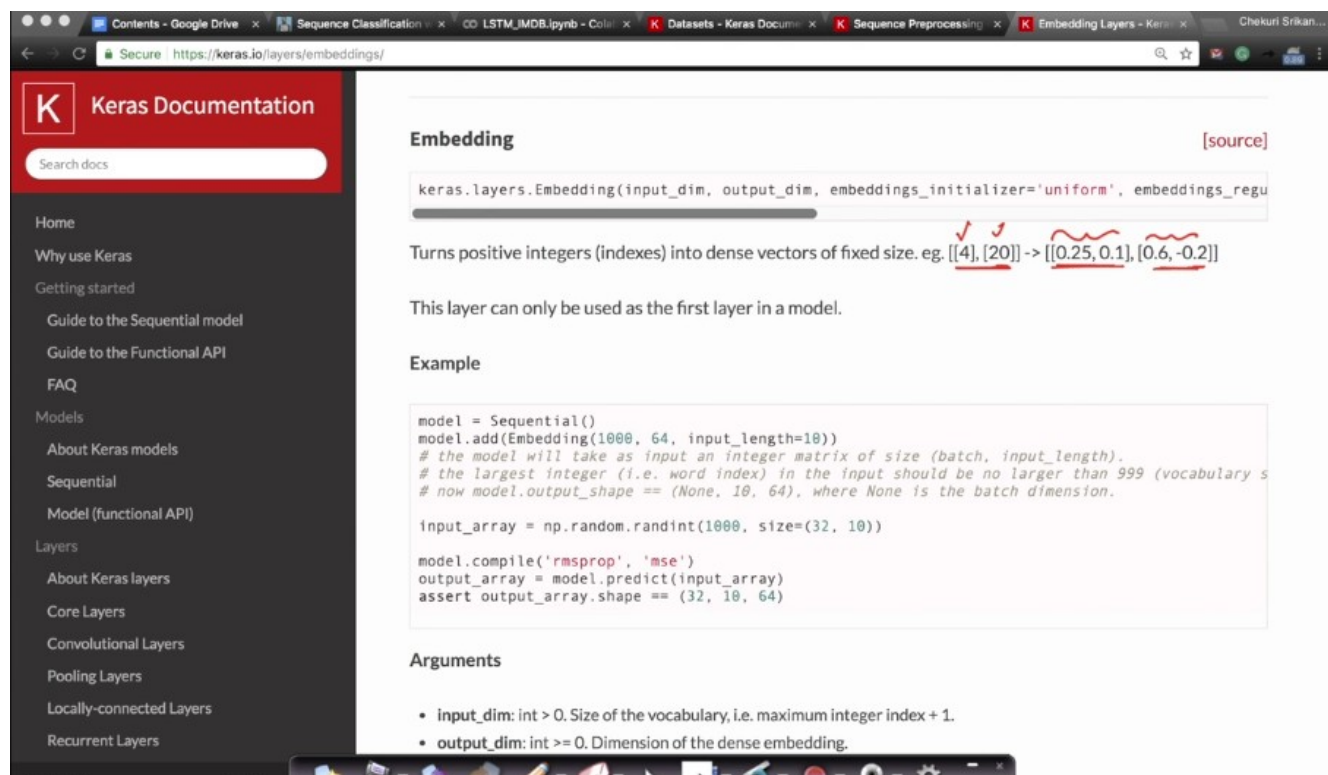
So, batch would be better option to train the RNN's, If we pad all the sequences. This will speed up the LSTM because the sequences are of same length.

The first layer we apply Is the embedding layer. These are used a lot in NLP tasks.



Each of my words can be represented as a dense vector.
LSTM 100 means, each of the LSTM takes the complete sequence of input one word after the other.