

• One of the advantages of the skip-gram architecture is that it can handle rare words more effectively compared to CBOW.

Since the skip-gram predicts multiple context words for a single target word, it can better capture the diverse context in which a word appears.

This property allows skip-gram to generate more accurate word embeddings, particularly for infrequent words, or words with limited occurrences in the training data.

## # Part of Speech Tagging (POS)

POS tagging is a task of labelling each word in a sentence with its appropriate part of speech.

why	not	tell	someone	?
adverb	adverb	verb	noun	punctuation mark, sentence closer

POS is a preprocessing task or step

### Uses

- Text Analysis
- Machine Translation
- Named Entity Recognition (NER)
- Information Retrieval
- Speech Recognition
- Chatbot
- Word sense disambiguation

I left the room

left of the room

# Hidden Markov Models (HMMs): Part 13 - Part 1

Ex

- (S) <sup>(N)</sup>nitish <sup>(V)</sup>loves <sup>(N)</sup>campusx <E>
- (S) <sup>(M)</sup>can <sup>(N)</sup>nitish <sup>(V)</sup>google <sup>(N)</sup>campusx <E>
- (S) <sup>(M)</sup>will <sup>(N)</sup>ankita <sup>(V)</sup>google <sup>(N)</sup>campusx <E>
- (S) <sup>(N)</sup>ankita <sup>(V)</sup>loves <sup>(M)</sup>will <E>
- (S) <sup>(N)</sup>will <sup>(V)</sup>loves <sup>(N)</sup>google <E>

N - noun  
V - verb  
M - model

vocabulary	N	M	V
nitish	2/10	0	0
loves	0	0	3/5
campusx	3/10	0	0
can	0	1/2	0
google	1/10	0	2/5
will	0/10	1/2	0
ankita	2/10	0	0

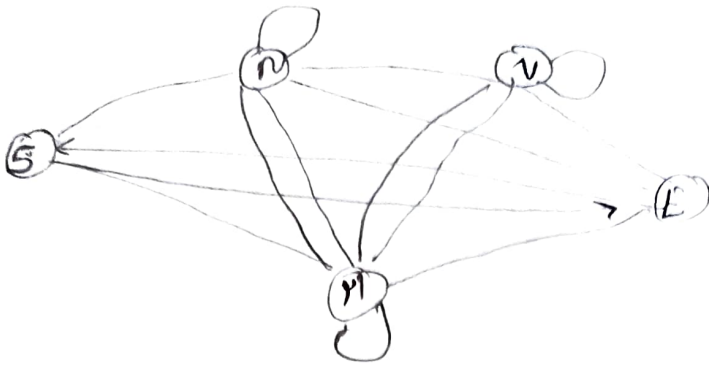
← Emission probability →

- ① nitish as a noun daya hai
- ② loves 3 baar daya hai, thirud baar as a verbs
- ③ google aik baar as a noun and 2 time as a verb daya hai

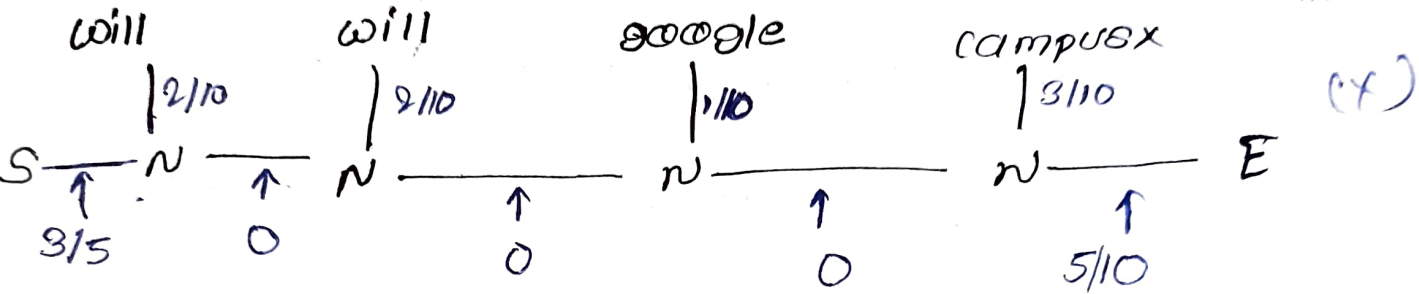
Total noun = 10  
model = 2  
verb = 5

## Transition Table

	N	M	V	E	
S	3/5	2/5	0	0	→ 5 sum
N	0	0	5/10	5/10	→ 10 sum
M	2/2	0	0	0	→ 2
V	5/5	0	0	0	→ 5

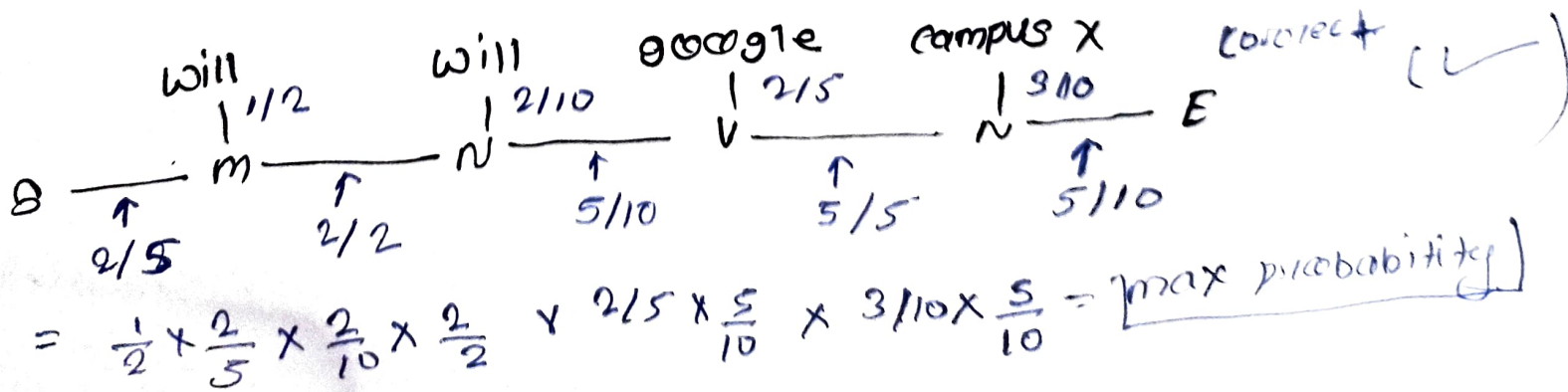
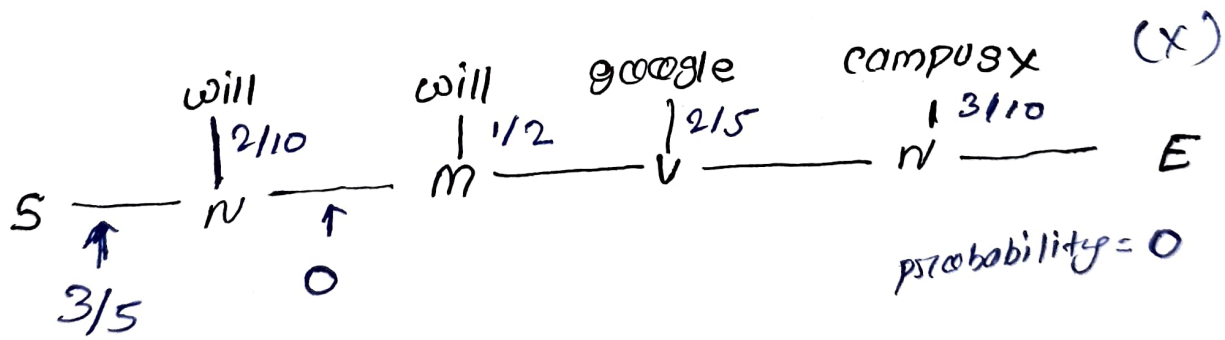


Let assume all would cost money



now multiply all and look for highest (biggest number encountered all possibility)

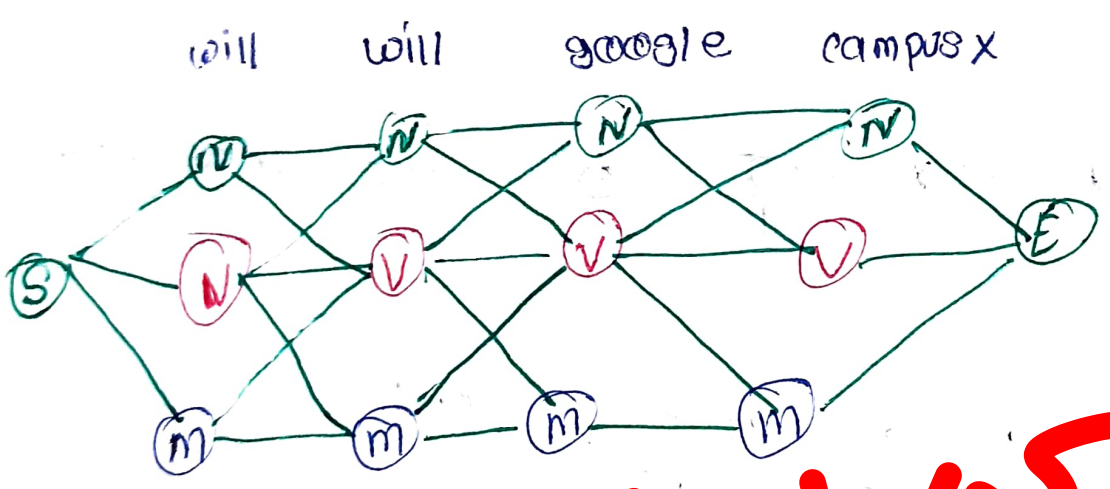
$$= \frac{2}{10} \times \frac{3}{5} \times \frac{2}{10} \times 0 \times \frac{1}{10} \times 0 \times \frac{3}{10} \times \frac{5}{10} = 0 \text{ (means not possible)}$$



will	will	google	campusx
(N)	(N)	(N)	(N)
(V)	(V)	(V)	(V)
(M)	(M)	(M)	(M)

Total combination =  $3^4$  (length of sentence) part of speech used on sentence

So, it is bruteforce method not good for large problem



Vishal  
Rishabh  
Kumar

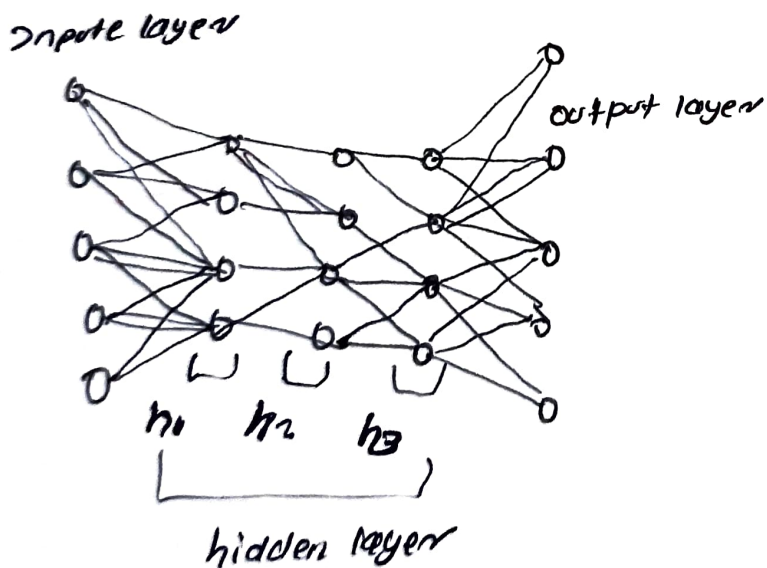


# [ deep learning vs neural network ]

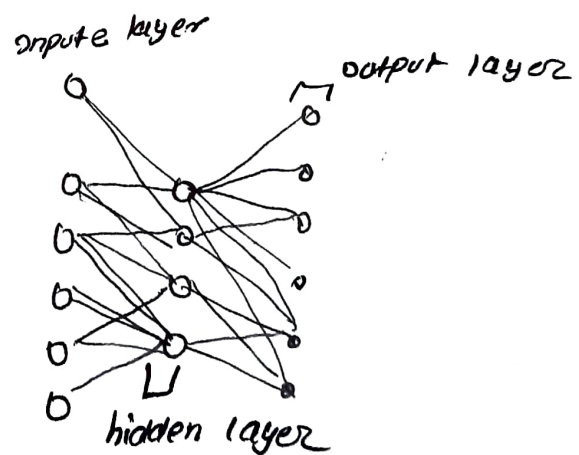
neural network: A neural network is a computational model inspired by the structure and functioning of biological neural networks, such as human brain.

Deep learning: Deep learning refers to a specific approach within machine learning that involves training deep neural networks. Deep learning focuses on neural networks with multiple hidden layers (hence term deep) allowing them to learn hierarchical representation of data.

## deep learning



## neural network



## # TRAX

Trax is an open-source deep learning library, developed by Google Brain's team. It is primarily designed for training and evaluating neural networks, particularly sequence models.

Trax aims to provide a simple and efficient framework for building, training and deploying machine learning models.

Trax is built on JAX and TensorFlow.

Trax combine the power of two deep learning libraries: JAX and TensorFlow. JAX provide the computational backend and high numerical operation, while TensorFlow is used for various functionality and utilities.

Trax take the advantage of GPU and TPU hardware accelerators.

Trax focus on sequence model, such as recurrent network (RNNs), transformer, and other variants.

[Advantage:]

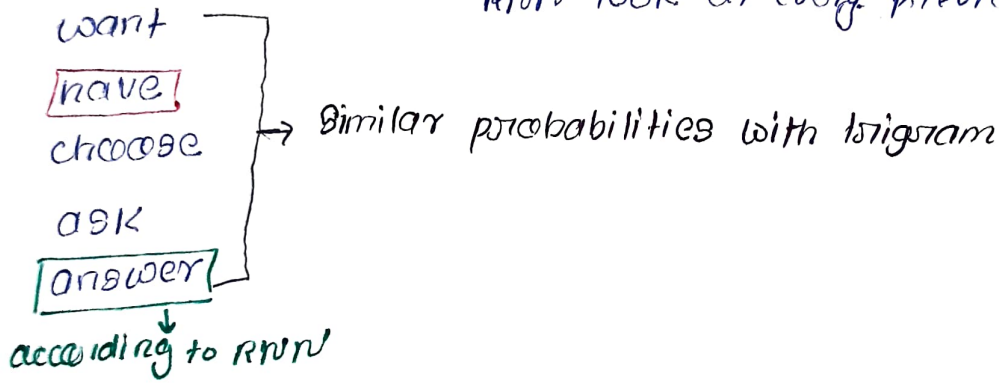
- Runs fast on CPU, GPU and TPUs
- parallel computing
- Record algebraic computation for gradient evaluation

# # Recurrent Neural Networks

- It is a type of neural network architecture specifically designed to process sequential model data.
- It is a type of neural network where the output from the previous step is fed as input to the current step.

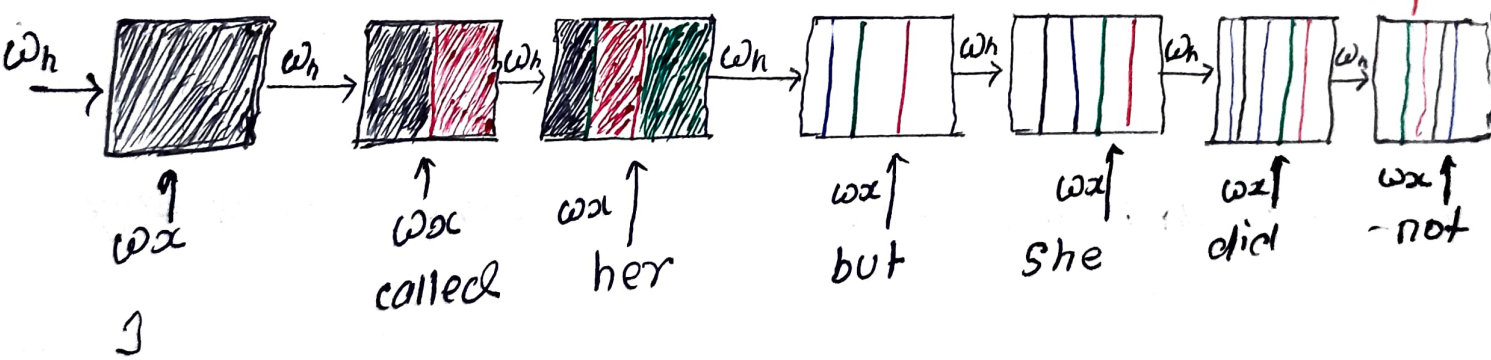
Rita was supposed to study with me. I called her but she did not \_\_\_\_\_

RNN look at every previous word



I called her but she did not \_\_\_\_\_

⇒ answer





## #way to implement an RNN model

one to one: Given score of chop, you can predict the winner

one to many: given an image, we can predict the caption it's going to be

~~many to many: given tweet, we can predict the caption~~

many to one: given tweet, we can predict sentiment

many to many: given an english sentence, we can translate in another

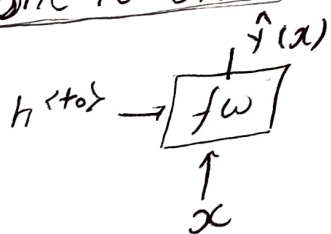
## advantage of RNNs

- Sequential model
- contextual understanding
- Flexible in input/output length
- parameter sharing

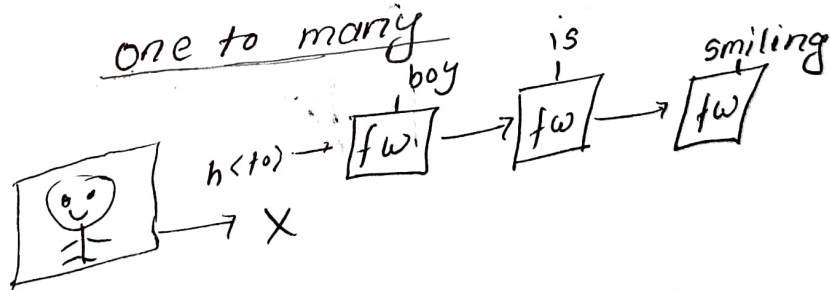
## Application of RNNs

- NLP
- Speech Recognition
- Time Series Analysis/forecasting
- Image Caption: RNN can combine with CNN
- Music Generation

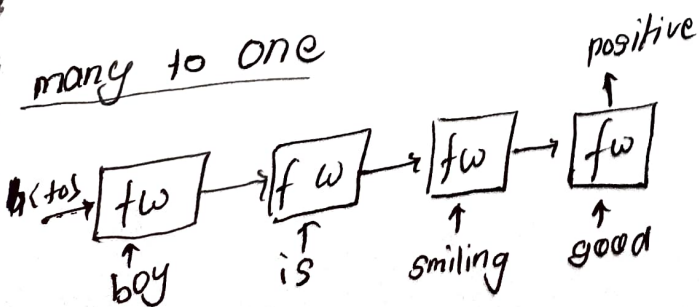
### one to one



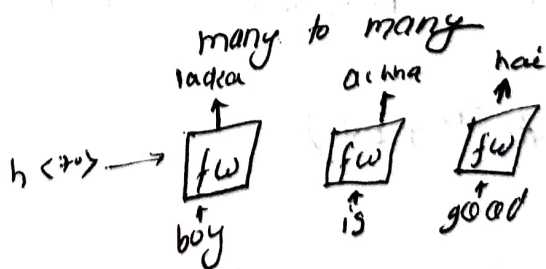
### one to many



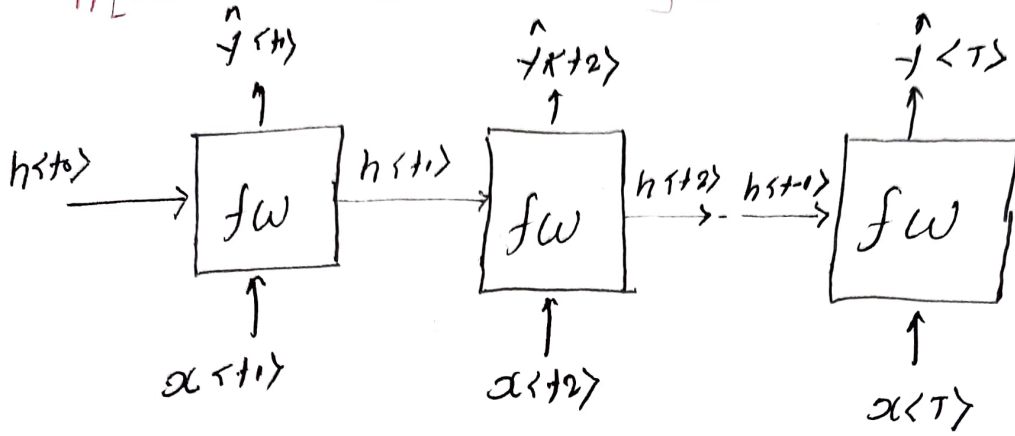
### many to one



### many to many



# # [Math in simple RNNs] [vanilla RNNs]



$$h \langle t \rangle = g(W_{hh} [h \langle t-1 \rangle, x \langle t \rangle] + b_h)$$

$$h \langle t \rangle = g(W_{hh} h \langle t-1 \rangle \oplus W_{hx} x \langle t \rangle + b_h)$$

$W_{hh}$  → weight of recurrent neuron

$W_{hx}$  → weight of input neuron

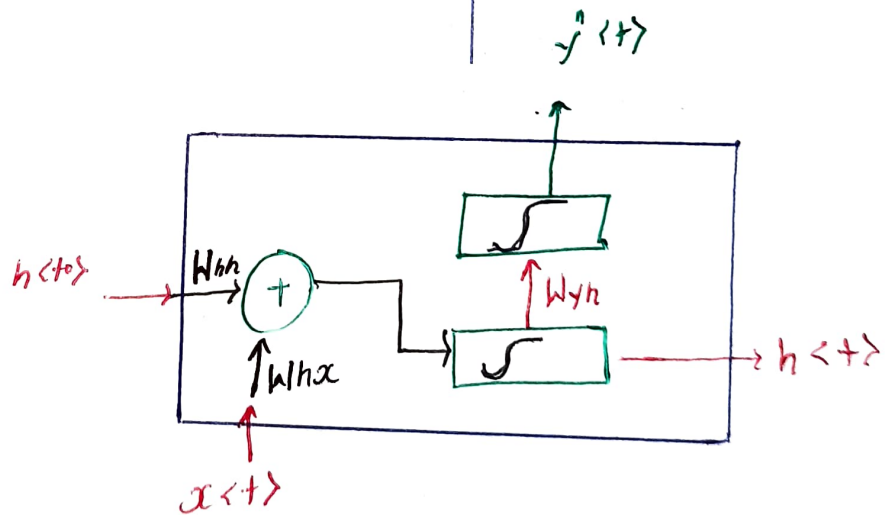
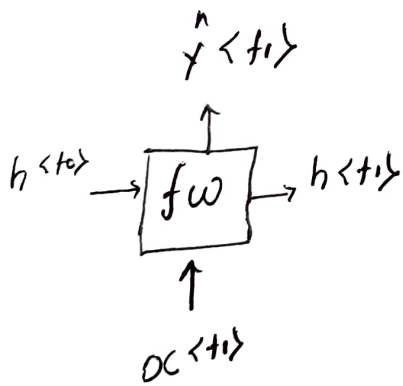
Formula for calculating current state

$$h_t = f(h_{t-1}, x_t)$$

$h_t$  → current state

$h_{t-1}$  → previous state

$x_t$  → input state



$$h \langle t \rangle = g(W_{hh} h \langle t-1 \rangle \oplus W_{hx} x \langle t \rangle + b_h)$$

$$\hat{y} \langle t \rangle = g(W_{yh} h \langle t \rangle + b_y)$$

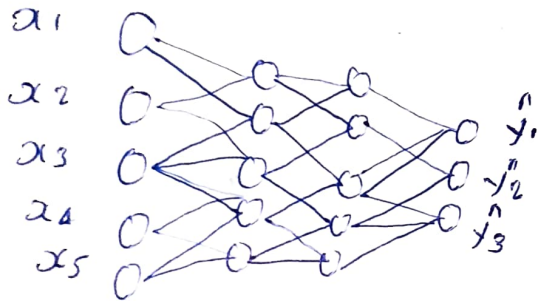
should be the size of matrix wh. if  $h^{(t)}$  had size  $4 \times 1$   
 and  $x^{(t)}$   $1 \times 1$

a)  $4 \times 14$

b)  $14 \times 14$

c)  $4 \times 10$

## COST FUNCTION FOR RNN'S



The cost function used in RNN is the cross entropy loss

$K$ -classes or possibilities

$$\text{Loss}(y) = - \sum_{j=1}^K y_j \log \hat{y}_j$$

Either 1 or 0

Looking at a single example ( $x, y$ )

Computing loss over several steps

$$J = - \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^K y_j^{(t)} \log \hat{y}_j^{(t)}$$

**vanilla RNN:** Also known as simple RNN, refers to the basic and original form of the recurrent layer in deep learning.

- It suffers from vanishing gradient problem

To address this vanishing gradient problem we use LSTM and GRU.

## # GRU (Gate recurrent unit)

The Gate Recurrent Unit (GRU) is a type of recurrent neural network (RNN) architecture that addresses the limitations of the traditional RNN, such as the vanishing gradient problem and difficulties in capturing long-term dependencies.

Three main components of GRU

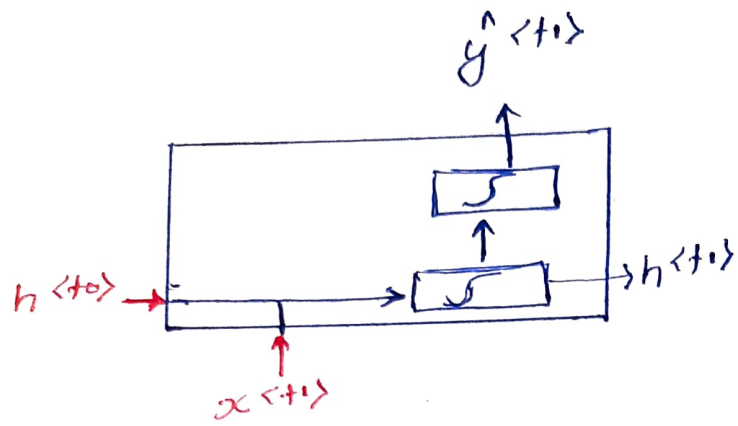
- 1) update gate ( $U$ ): It decides how much of the previous hidden state should be kept and how much of the new information should be added to the current hidden state.
- 2) Reset Gate ( $r$ ): It controls how much of the previous hidden state should be forgotten.
- 3) candidate Hidden state ( $h_c$ ): It represents the new information that will be added to the hidden state.
- 4) current Hidden state ( $h$ ): It is the output of the GRU layer and serves as the hidden state for the next time step.

Gates to keep/update relevant information in the hidden state

$$\Gamma_r = \sigma(w_r [h^{(t)}, x^{(t)}] + b_r)$$

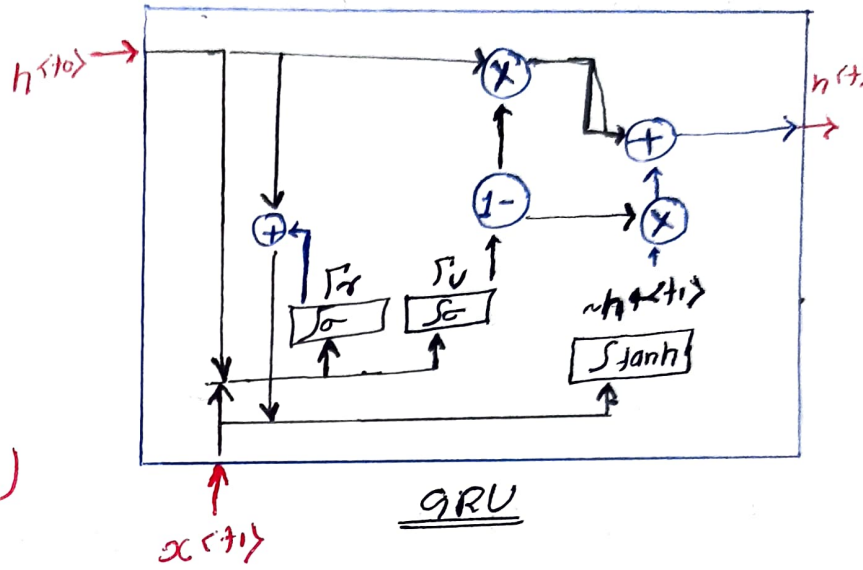
$$\Gamma_u = \sigma(w_u [h^{(t)}, x^{(t)}] + b_u)$$

# # Vanilla RNN VS GRU's



$$h^{(t+1)} = g(W_h [h^{(t+1)}, x^{(t+1)}] + b_h)$$

$$y^{(t+1)} = g(W_y h^{(t+1)} + b_y)$$



Gates to keep/update relevant information in the hidden state

$$\begin{aligned} \text{Forget gate} &\rightarrow \Gamma_r = \sigma(W_r [h^{(t)}, x^{(t+1)}] + b_r) \\ \text{Update gate} &\rightarrow \Gamma_u = \sigma(W_u [h^{(t)}, x^{(t+1)}] + b_u) \end{aligned}$$

$$\text{hidden state candidate} \rightarrow \tilde{h}^{(t+1)} = \tanh(W_h [\Gamma_r * h^{(t)}, x^{(t+1)}] + b_h)$$

$$h^{(t+1)} = (1 - \Gamma_u) * h^{(t)} + \Gamma_u * \tilde{h}^{(t+1)}$$

$$y^{(t+1)} = g(W_y h^{(t+1)} + b_y)$$

update gate  
lower + hidden  
state

# # LSTM (Long Short-Term Memory)

Long Short-Term memory is a type of recurrent neural network architecture that address the limitation of traditional RNNs in capturing and preserving long-term dependencies.

LSTM are designed to overcome vanishing and exploding problem.

Memory cell: The memory cell is responsible for storing and carrying information across time steps.

## application

- ⇒ input gate
  - ⇒ forget gate
  - ⇒ output gate
- } Similar to GRU

- next-character prediction
- chatbots
- image captioning
- speech recognition
- music composition

GRU + Memory cell ⇒ LSTM (we can say that)

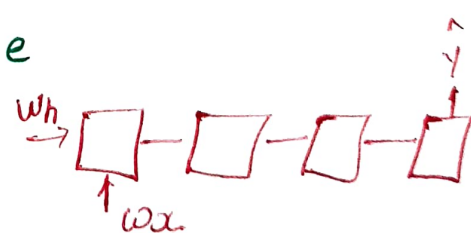
## ##

How vanishing and exploding gradient happen

→ Back propagation through time

$$\frac{\partial L}{\partial w_h} \propto \sum_{k \leq t} \left( \prod_{i=k}^t \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial w_h}$$

→ contribution of hidden state k



Length of the product proportional to how far k is from t

Partial derivative < 1

contribution goes to 0 → vanishing gradient

partial derivative > 1

contribution goes to infinity → exploding gradient

## # Solution to vanishing gradient problems

1) Identity RNN with ReLU activation

$$\begin{bmatrix} 1 & -1.0 & -0.01 \\ -0.1 & 1 & -0.1 \\ 0 & 0 & -0.2 \\ 1.1 & 1 & 1 \end{bmatrix} \xrightarrow{\text{ReLU activation}} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} \quad \cdot 1 \rightarrow 0$$

$$\text{ReLU}(x) = \max(0, x)$$

2) Gradient clipping: Gradient clipping is a technique where the gradient are clipped or capped at a certain threshold during back propagation. This prevent the gradient from exploding and help to stabilize training.

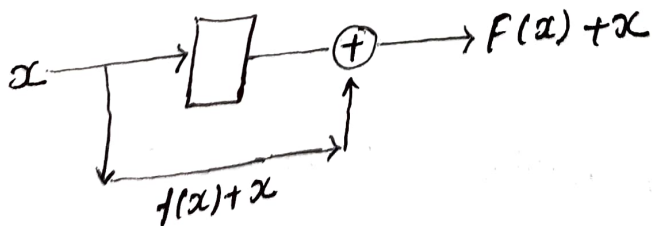
ex. threshold set = 25

$$\text{then } 32 \rightarrow 25$$

$$46 \rightarrow 25$$

$$6 \rightarrow 6$$

3) skip connections (skip activation function)



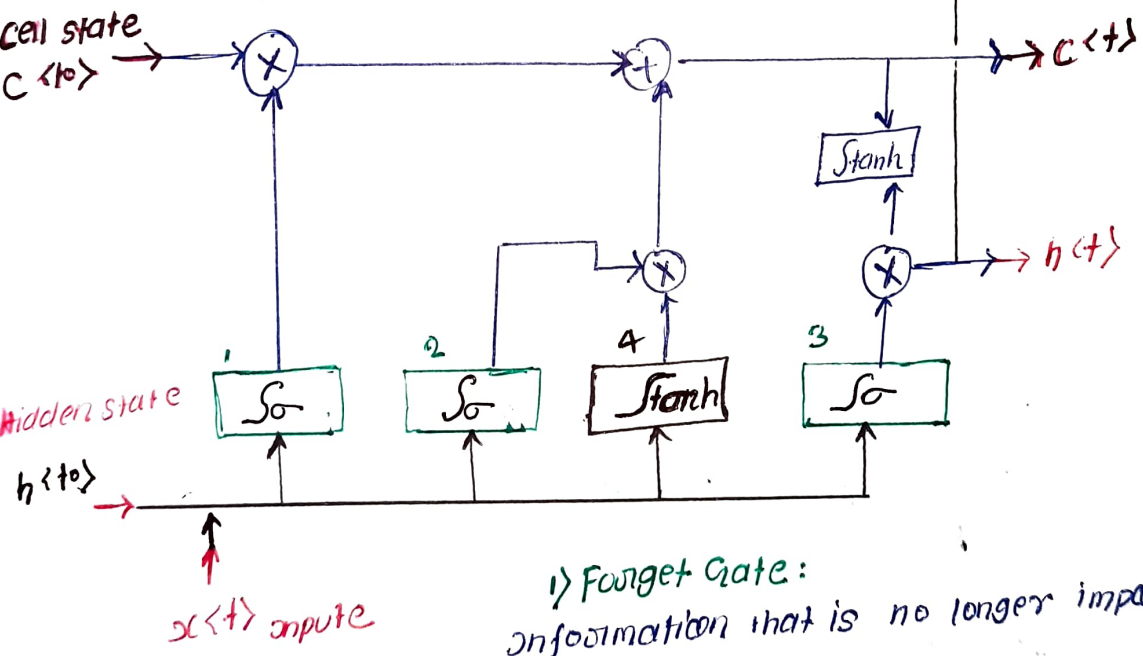
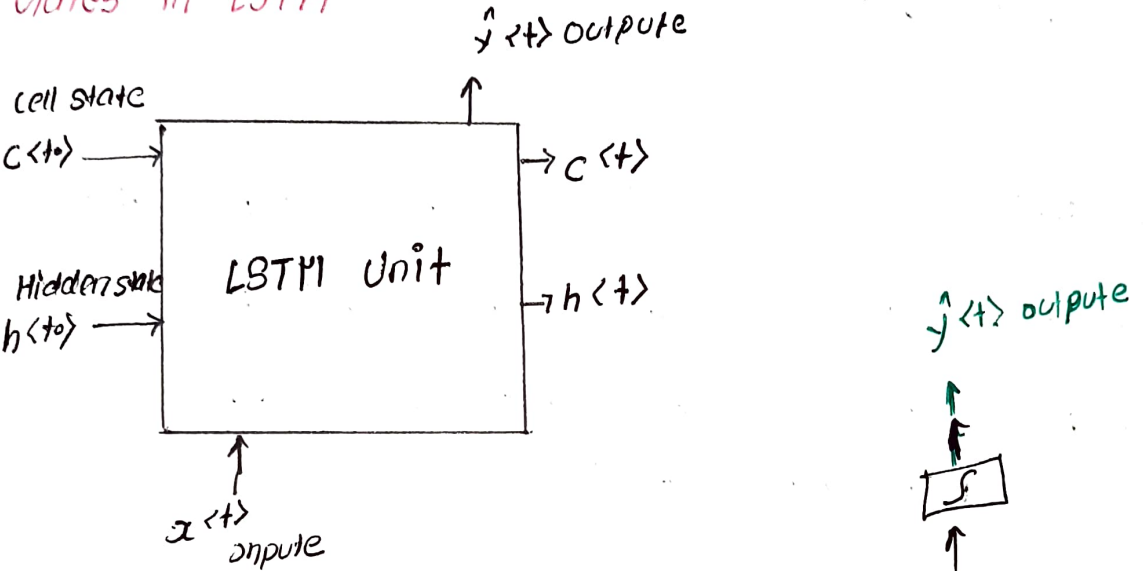
(a) order of gate that information flows through in an LSTM unit.

a) Forget  $\rightarrow$  output  $\rightarrow$  input gate

b) input  $\rightarrow$  Forget  $\rightarrow$  output gate

c) Forget  $\rightarrow$  input  $\rightarrow$  output gate

### Gates in LSTM



- 1) Forget Gate: information that is no longer important
- 2) input gate information to be stored
- 3) output gate information to use at current step



## Sigmoid output between 0 and 1

0 → close gate (information does not get through)

1 → information get through freely

4

$\tanh$  → candidate cell state

information from the previous hidden state and current input

$\tanh$  shrinks argument to be between -1 and 1

## new cell state

add information from the candidate cell state using forget and input gate.

## Summary

LSTM use a series of gates to decide which information to keep:

- forget gate decide what to keep and what to forget
- input gate decides what to add
- output gate decides what the next hidden state will be

## # Name Entity Recognition (NER)

is a NLP task that involves identifying and classify name entities in text into predefined categories such as person name, location, date, time etc.

Ex: pooran is going to see Taj Mahal tomorrow sunday at 8 pm  
person location date time

## Application

- search engine efficiency
- recommendation engine

- customer service
- automatic trading

## Training NERFs: data processing

- convert word and Entity classes into arrays:
- pad with tokens
- create a data generator

## Training on NER

- create a tensor for each input
- put them into batch (32, 64, 128...)
- feed into LSTM unit
- run the output through a dense layer
- predict using log softmax over k classes

—END of week—

## # Siamese network

A siamese network is a type of neural network architecture that are designed to compare and measure similarity or dissimilarity between two input sample.

Ex. How old are you = what is your age

where are you from  $\neq$  where are you going

## Application

- Handwritten checks

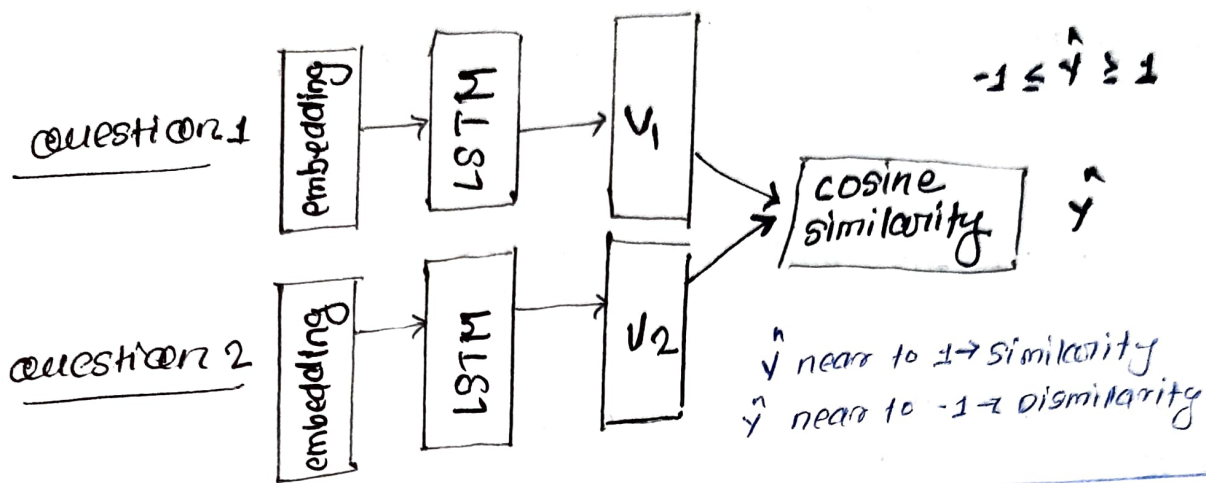
B      Mr

- question duplicates

How old ~~are~~ you?  
what is your age?

- queries

# Siamese network Architecture



- 1) Embedding
- 2) LSTM
- 3) vectors
- 4) cosine similarity

The architecture of a siamese network consists of two or more identical subnetworks, referred as 'twin' or 'siamese twin' share the same set of weight and parameters.

## Cost Function

How old are you? Anchor

What is your age? positive

Where are you from? negative

$$\cos(v_1, v_2) = \frac{v_1 \cdot v_2}{|v_1| \times |v_2|}$$

cosine similarity

$\cos(A, P) \approx 1$  for good model  
 $\cos(A, N) \approx -1$  negative

$$\text{Loss} = \cos(A, N) - \cos(A, P) \leq 0$$

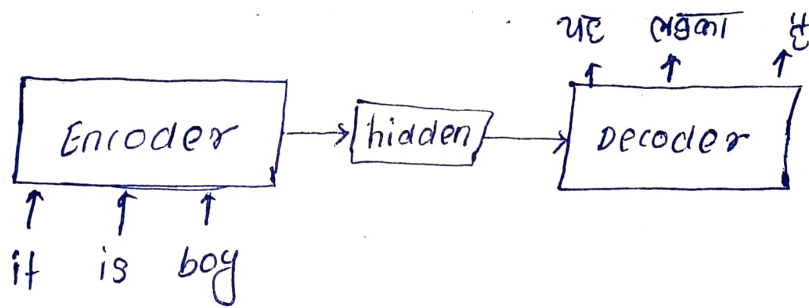
## One shot learning

The key reason for one shot learning is to be able to classify new classes without retraining any models.

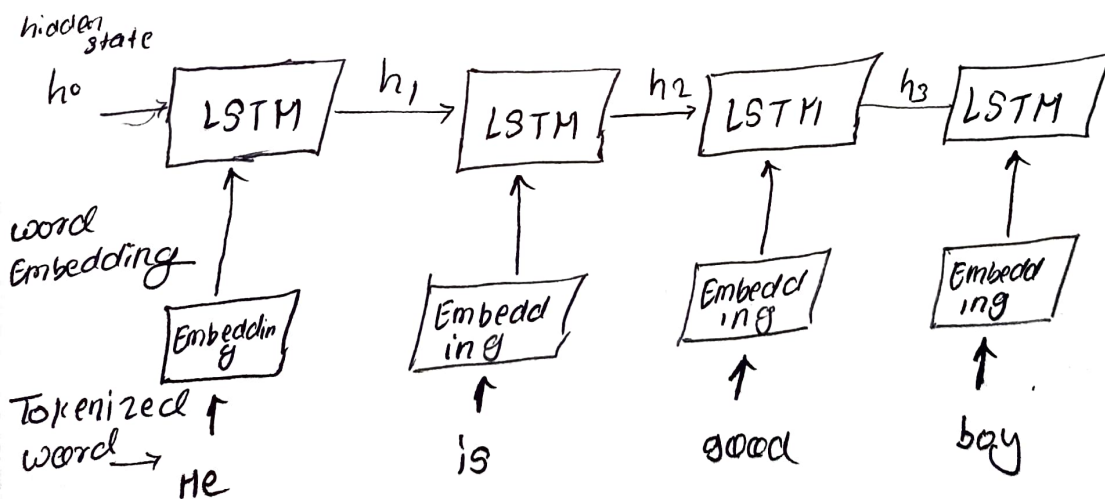
# # Seq2Seq Model

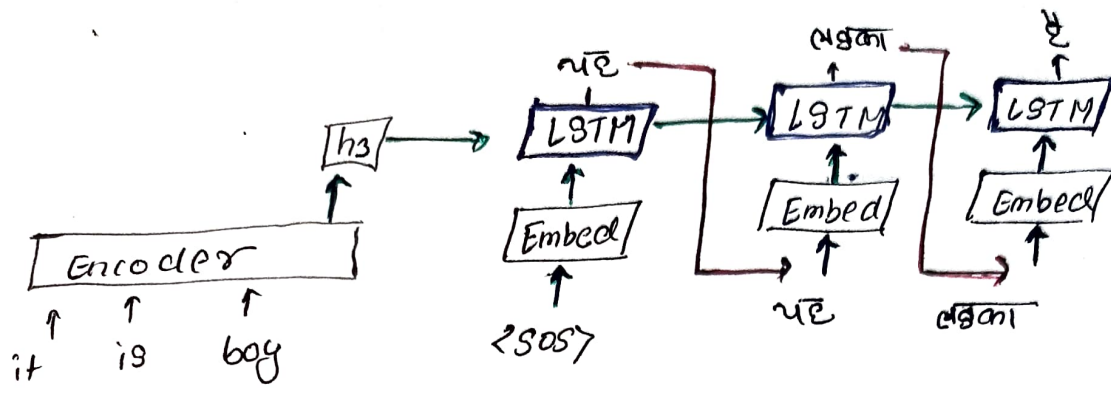
Sequence to sequence model commonly used in task like machine translation.

- The model consist of two main component: an encoder and decoder.
- The **encoder** processes the input sequence and convert it into fixed length vector called **context vector**
- The **decoder** takes the context vector as input and generate the output sequence step by step



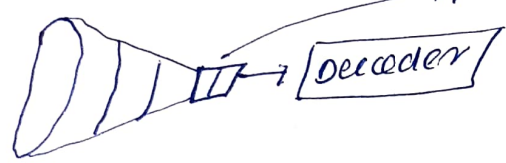
## Seq2Seq encoder





Problem

Information is bottlenecked → Fixed hidden state size



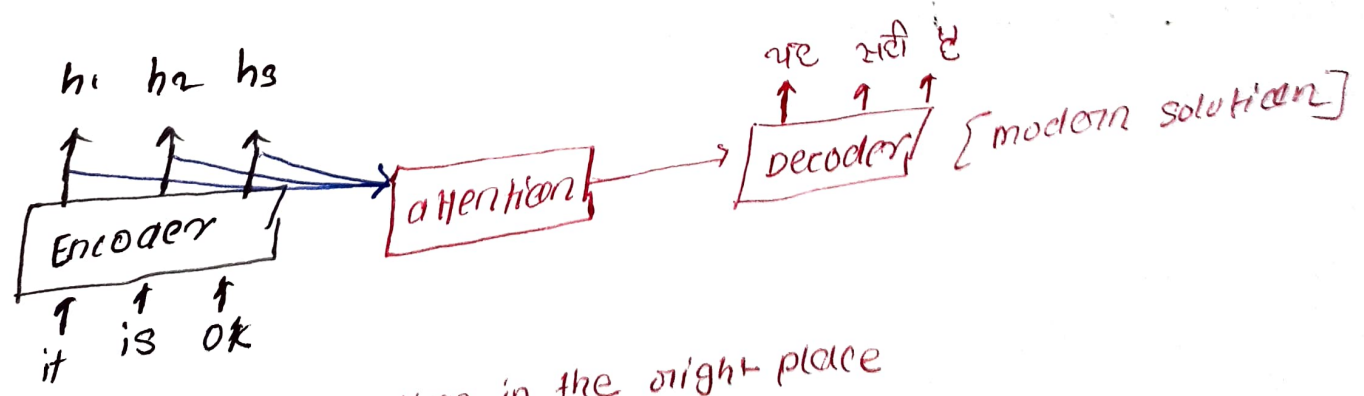
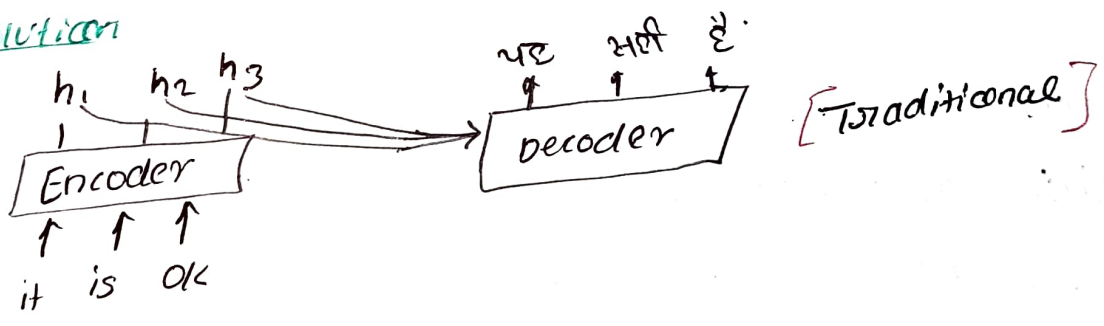
• A fixed amount of information goes to the decoder

(Q) Why are longer sequences problematic for traditional seq2seq model

Ans Because seq2seq relies on fixed-length memory

• As sequence size increases, model performance decreases

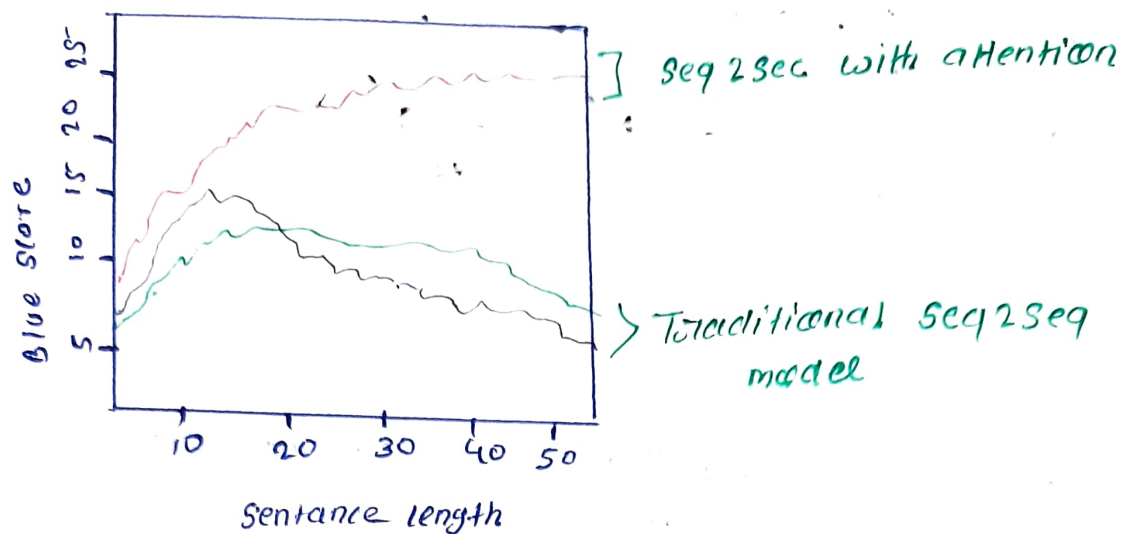
Solution



or Solution: focus attention in the right place

## # Performance

0505



on traditional model performance is decrease as word increasing.

- Attention is a layer that lets a model focus on what important
- Queries, values, and key are used for information retrieval inside the attention layer.

## Blue score (precision)

Blue score (Bilingual Evaluation Understanding) is a matrix commonly used to evaluate of the quality of machine translation output by comparing one or more references translation.

The BLEU score measures the similarity between the machine-generated translation and the reference translation, assigning a score between 0 and 1.

A higher BLEU score indicates a better match between the machine translation and the reference translation.

candidate = [ ] [ ] [ am ] [ ] - 4

reference 1 = younes said [ ] am hungry

reference 2 = He said [ ] am hungry

$$\text{BLEU score} = \frac{1 + 1 + 1 + 1}{4} = 1$$

First = ] present in both sent (1)  
 second = ] " " " (1)  
 third = am " " " (1)  
 fourth = ] " " " (1)

$$\text{Count-Total} = 4$$

### BLEU score (modified)

candidate [ ] [ ] [ am ] ]  
 reference 1 younes said [ ] [ am ] hungry  
 reference 2 He said [ ] [ am ] hungry

$$\text{BLEU score} = \frac{1 + 1}{4} = 0.5$$

First: ] present in both count (1), then delete ] from both  
 second: ] not present in both (0)  
 third: am present in both count (1) then delete it  
 fourth: ] not present =  $\frac{2}{4}$