

GRU, Seq to seq mapping

First see this practical ipynb file:
Deep_and_Bi_directional_RNN_LSTM_GRU.ipynb

Following are the 3 types of Sequence to Sequence mappings:

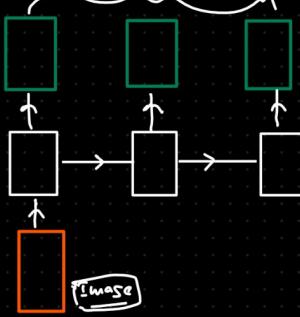
- ① One to many
- ② many to one
- ③ many to many

Encoder and decoder (2016)
Seq to Seq

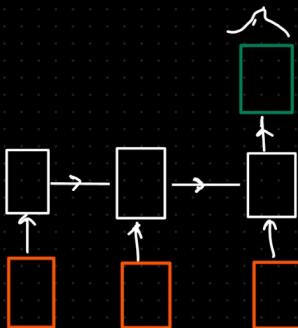
RNN \Rightarrow 1980-1990

Seq to Seq \Rightarrow RNN, LSTM, GRU

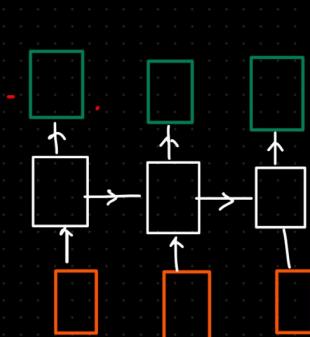
One to many }



many to one }



many to many ?



Application \Rightarrow Image captioning

Application \Rightarrow Sentiment analysis

Application \Rightarrow machine (Google translation)

Seq to Seq \Rightarrow Classical \Rightarrow RNN, LSTM, GRU

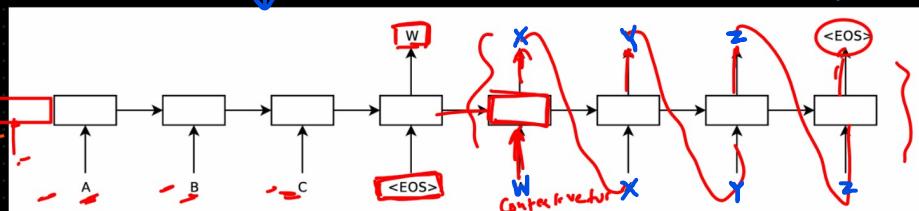
{ Fix length outputs }

Seq to Seq \Rightarrow variable length inputs

Encoder-Decoder architecture

I/P: A B C
O/P: W X Y Z
Variable Dimensionality

Var



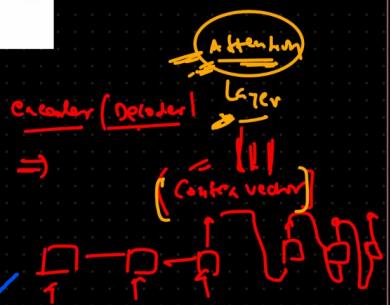
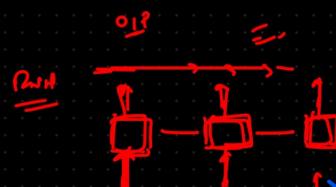
EOS represents the End of Sentence. Cells upto first EOS represents Encoder side whereas, beyond that cells are of Decoder side. Output of Encoder is denoted by W which is context vector which is then passed to the decoder side in cascading fashion (Observe passing of W, X, Y, Z)

RNN, LSTM, GRU are considered as the classical seq to seq mapping models. Let's understand the drawback with this classical model.

So, these models work fine in case of one to many (Image captioning where images can be flatten and then passed) and many to one use case. But in case of many to many these are going to work where we have fixed length input and output. But in real word scenario Input and output length may vary. For e.g: English to Hindi Google translator -

"I wanna be" \rightarrow "Mai banna cahta hu" we can observe that when we passed input of len 3 in get translated into hindi with len 4. Classical models will fail in such case since they work only in case of fixed length.

The problem that has been stated above gave birth to modern day used architecture of seq to seq mapping called Encoder Decoder.



RNN/LSTM/GRU were the timestamp based approach where in each timestamp we get an output. There is no single vector obtained containing the context of whole input passed. Whereas, In Encoder Decoder which is not a timestamp based approach, we are wrapping up the entire input into a single vector called context vector and then using decoder to basically decode it.

55
65
70
100
Covars

30
30-40
O1P

If possible go on a high level through Sequence to Sequence Neural Network 2014 research paper which is for Encoder Decoder.

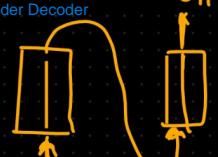
Encoder-Decoder \Rightarrow LSTM, GRU, RNN

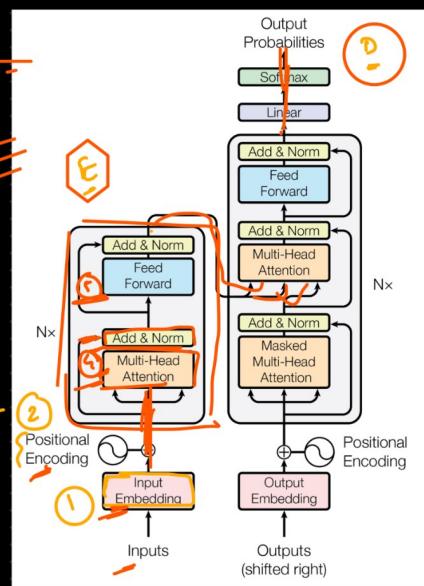
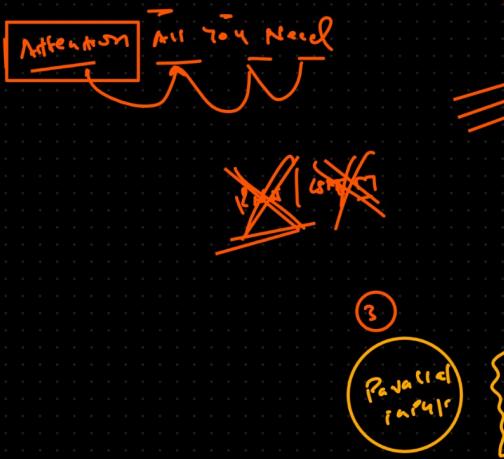
2018

Google \Rightarrow transformer

LSTM

GRU





Transformer Architecture. You may check your LinkedIn post for better understanding.

= break through

Sat/sun transformer

Encoder Decoder

Practice

→ Case Study

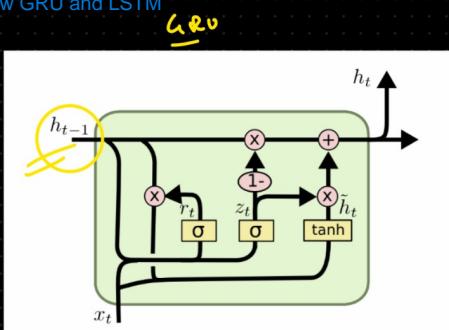
211 eval to end.

Note:
Please note that in the Encoder Decoder architecture we were still using RNN/LSTM/GRU as the hidden cell. There was one drawback of Encoder Decoder architecture, that is the context vector was able to hold the context of only 30-40 words at a time. But in real word scenario we can have words more than that. So, in 2018 Google published the research paper "Attention is all you need" where they simply added an Attention layer on top of Encoder Decoder architecture and were able to remediate this issue. This new architecture is what we call as transformer that does not use RNN/LSTM/GRU and is used to power up the modern day LLMs.

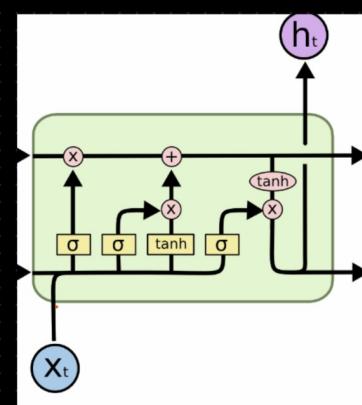
Decoding GRU architecture:



Difference b/w GRU and LSTM



GRU



LSTM

(1) Number of gates

2 gates

- (1) Reset gate
- (2) Update

(2) Memory Content

Single hidden state

From architecture we can clearly observe that GRU only has hidden state from previous timestamp and not any separate cell/memory state. Therefore, it is going to use hidden state for both passing the output of current cell to next cell and also as a memory unit to hold the context. Due to this change only GRU is able to hold more long term dependencies as compared to LSTM.

3 gates

- (1) Forget
- (2) iP or Update
- (3) OlP

- (1) Cell State
- (2) Hidden State

Here we have both cell state Ct-1 and hidden state ht-1

③

Parameter complexity

Less Parameter

Less trainable parameter
since we have only 2 gates

③

More Parameter

More trainable parameter since
we have 3 gates

④

Computation efficiency

less intensive

④

more powerful and more intensive

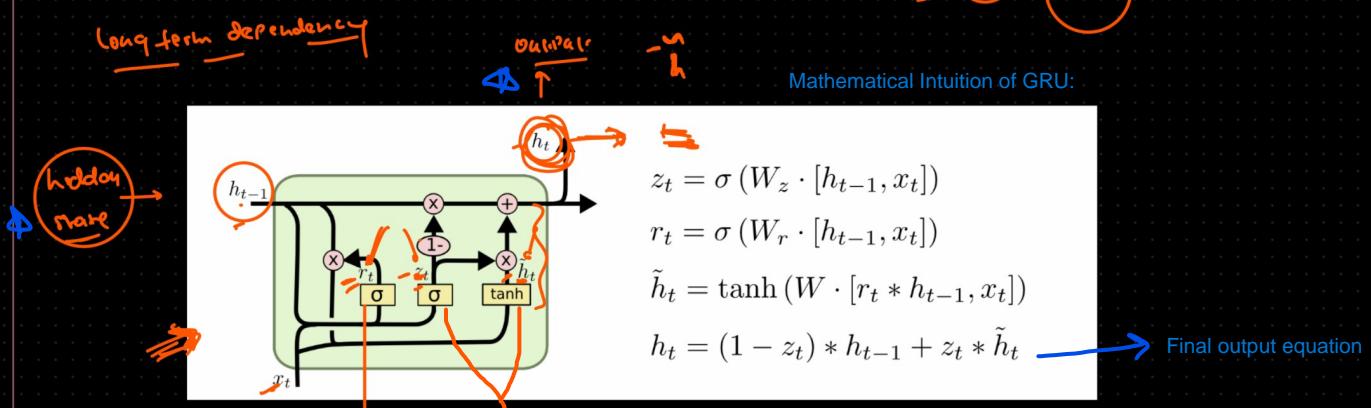
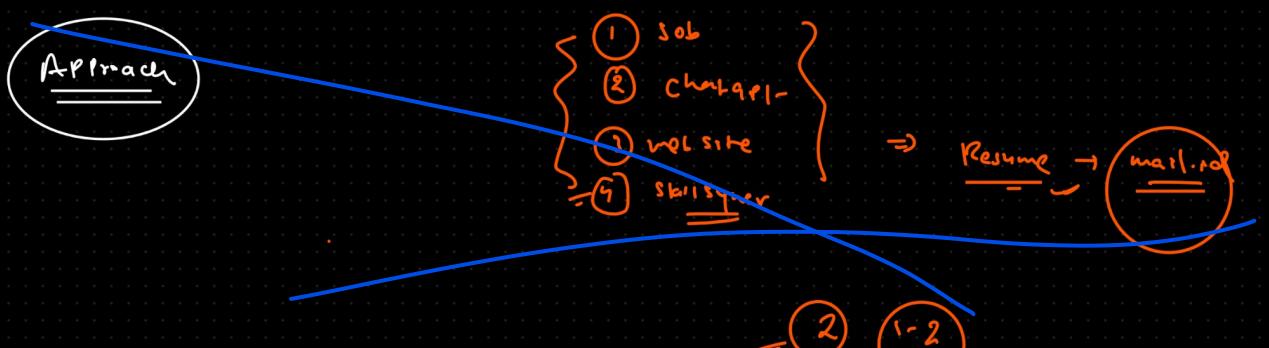
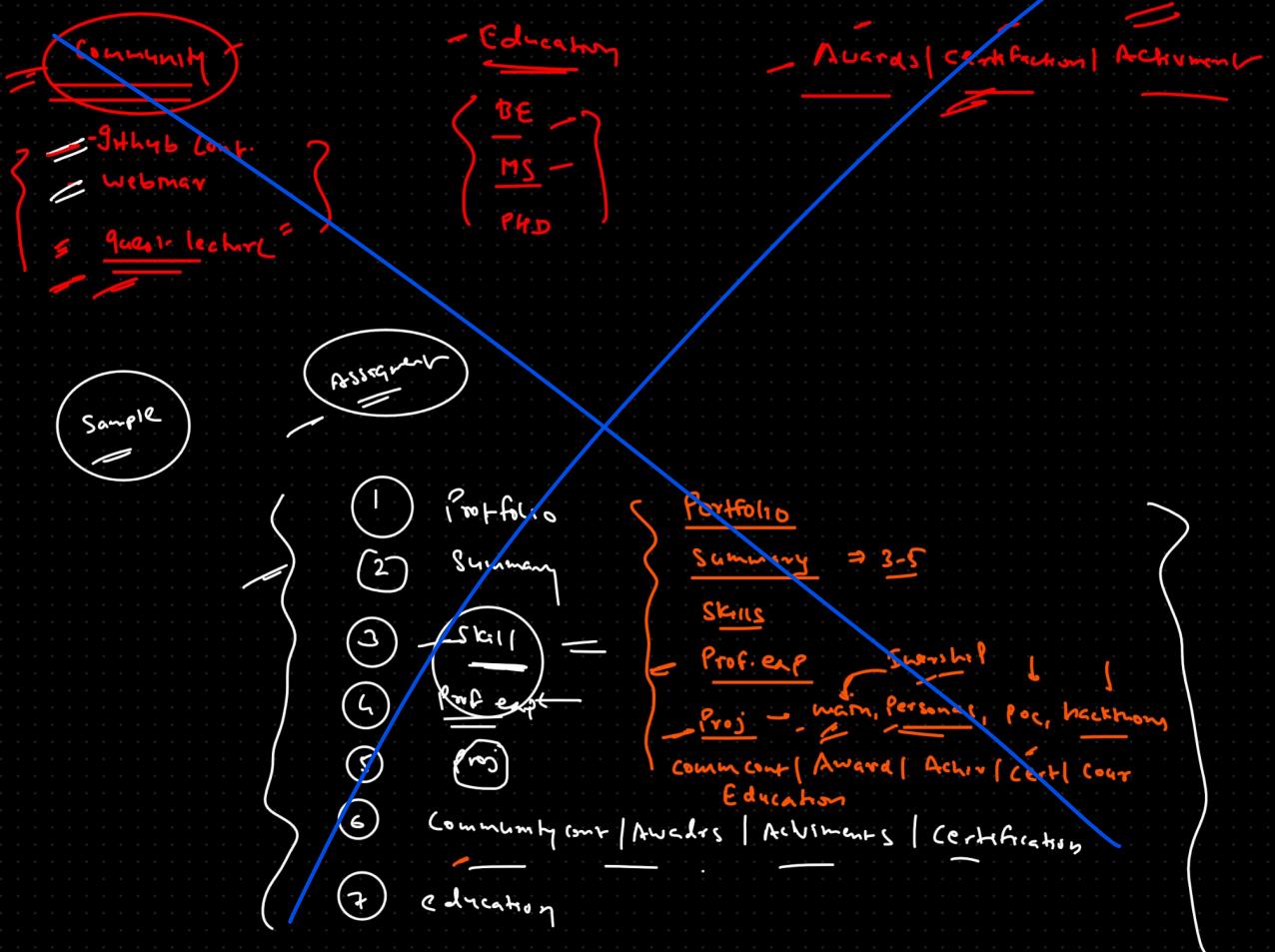
⑤

Computation time

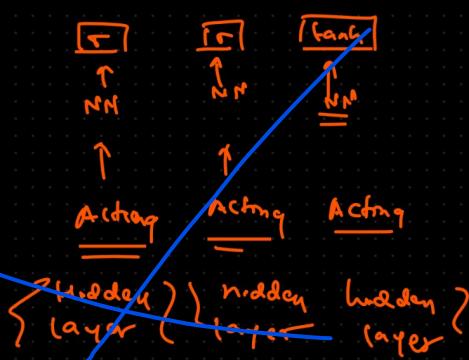
less

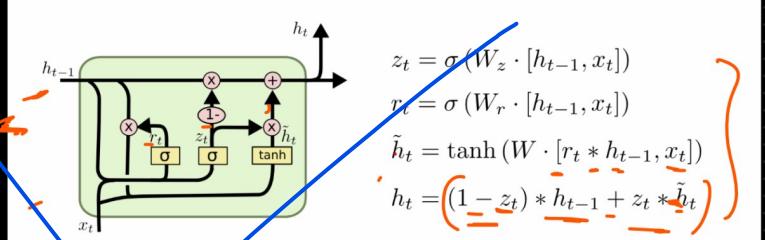
⑤

more time

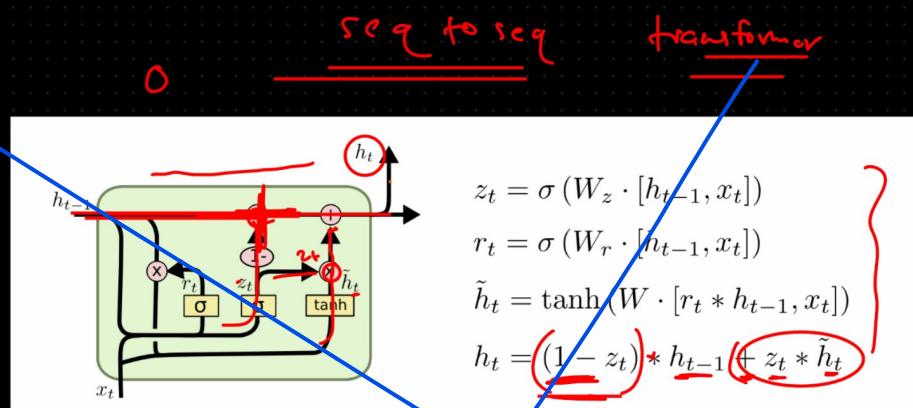
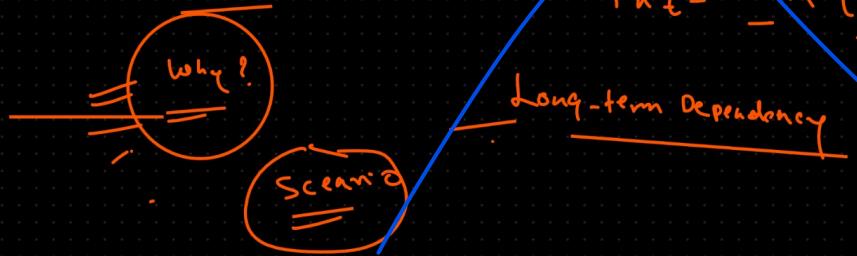


word by word
= on a diff timestamp





LSTM v/s GRU
this is happening



(How) Diffr Long v/s GRU

why? \Rightarrow context

theoretical \Rightarrow project Create Open

Conclusion

Deep RNN | LSTM | GRU

Backprop RNN LSTM | GRU

= seq to seq mapping \Rightarrow RNN

= Encoder to decoder seq to seq

= transformer

= GRU v/s LSTM Mathematical

= Why?