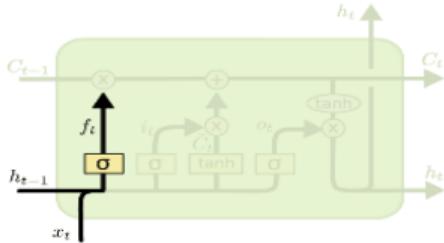




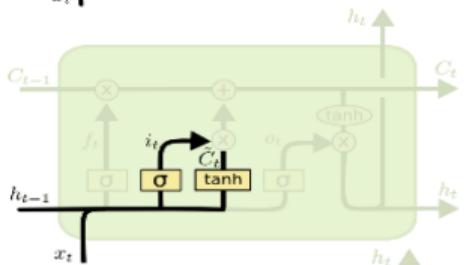
BITS Pilani
Pilani | Dubai | Goa | Hyderabad

DEEP LEARNING MODULE # 7 : Attention Mechanism

Dr YVK Ravi Kumar

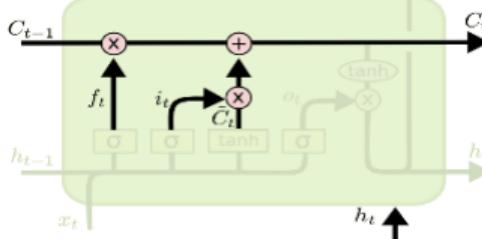


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

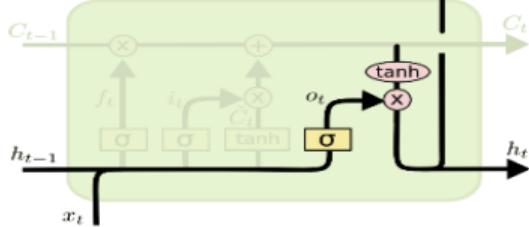


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

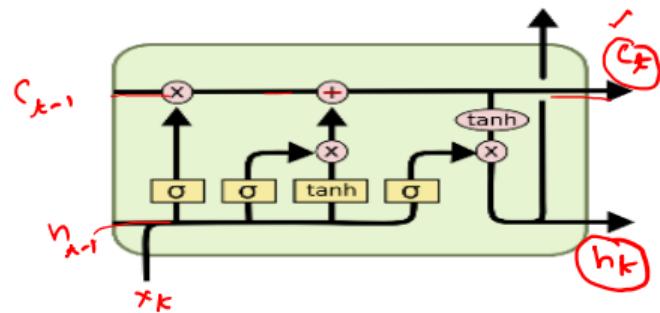


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$



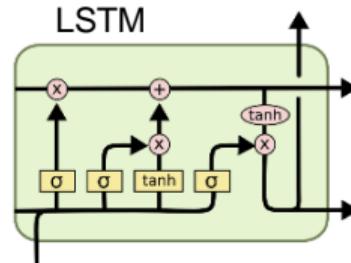
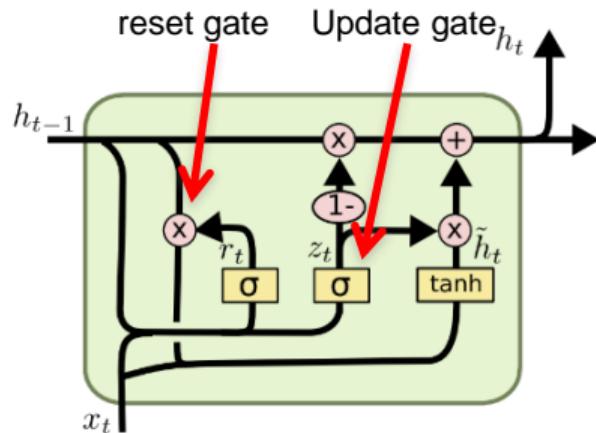
i_t decides what component is to be updated.
 C'_t provides change contents

Updating the cell state

Decide what part of the cell state to output

GRU – gated recurrent unit

(more compression)



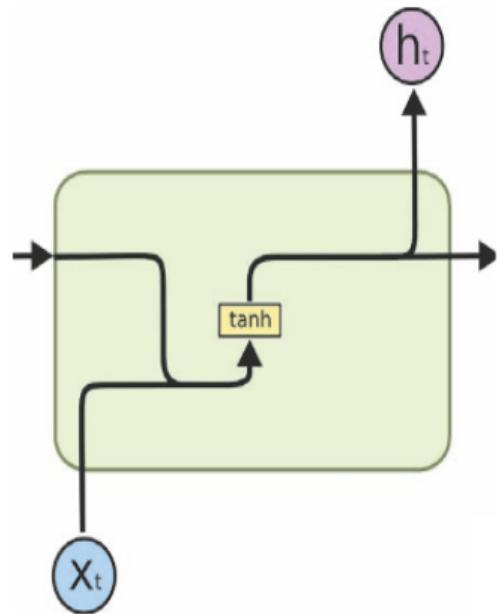
$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

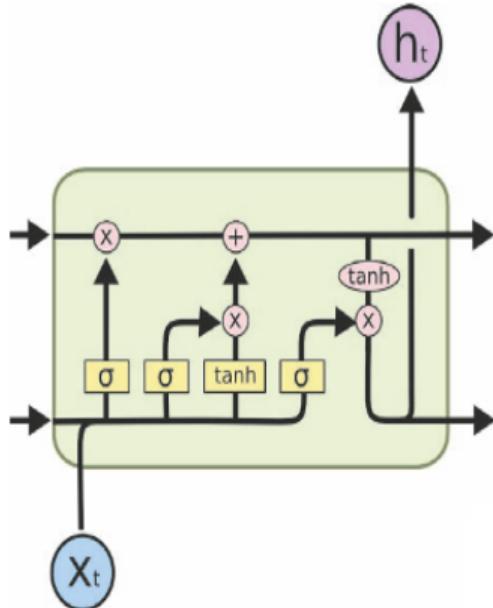
$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

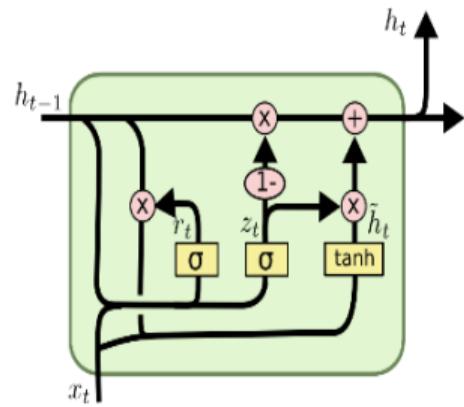
- It combines the **forget** and **input** into a single **update gate**.
- It also merges the cell state and hidden state.
- This is simpler than LSTM. There are many other variants too.



(a) RNN



(b) LSTM



Example: Name entity recognition

He said "Teddy bear is soft."

NER

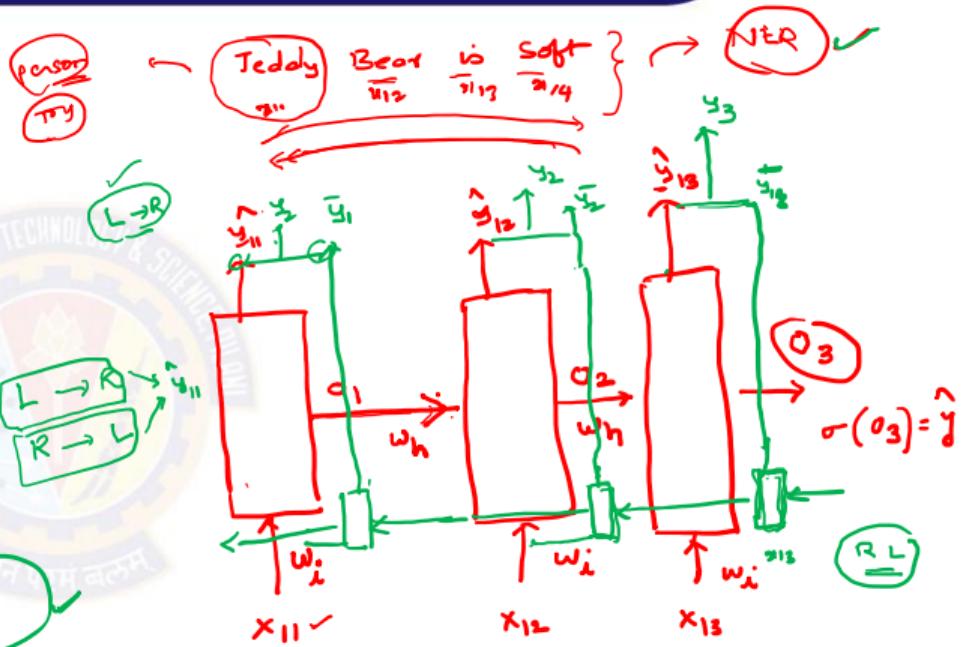
He said "Teddy Roosevelt was a President."

RNN → ?

Bi RNN ✓
Bi LSTM ✓

RNN ✓ $\vec{h}_t = \tanh(\vec{w} \vec{h}_{t-1} + \vec{u} \vec{x}_t + \vec{b})$
✓ $\overleftarrow{h}_t = \tanh(\overleftarrow{w} h_{t+1} + \overleftarrow{u} x_t + \overleftarrow{b})$

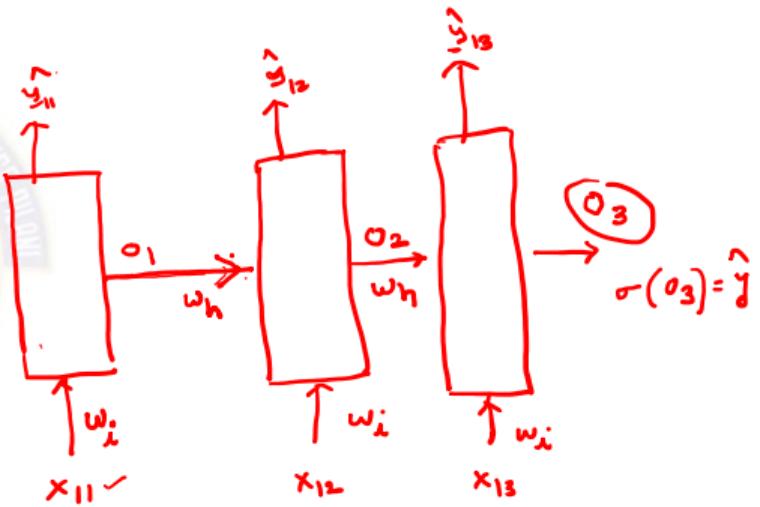
$y_t = \sigma [v [\vec{h}_t, \overleftarrow{h}_t] + b]$



$$f(x_{11} w_i) = o_1$$

$$f(x_{12} w_i + o_1 w_h) = o_2$$

$$f(x_{13} w_i + o_2 w_h) = o_3 \quad \sigma(o_3) = \text{output}$$



$$f(x_{11}w_i) = o_1$$

$$f(x_{12}w_i + o_1w_h) = o_2$$

$$f(x_{13}w_i + o_2w_h) = \underbrace{o_3}_{\text{output}} = \sigma(o_3)$$

BIDIRECTIONAL RNN (BRNN)

- Forward and backward connections.
- The blocks can be RNN, GRU, LSTM.
- Mostly used in the NLP.
- Acyclic graph
- Example: Name entity recognition
He said “Teddy bear is soft.”
He said “Teddy Roosevelt was a President.”

SUMMARY

- Use GRU, when dependency is short. Eg: Weather forecasting ✓
- Use LSTM, when dependency is long. Eg: NLP Translation ✓
- Use BRNN, dependency is in both direction. Eg: Stock prediction ✓

ANN → tabular data
CNN → images
RNN / LSTM, GRU → sequences
Machine translation ✓
Text Summarization ✓

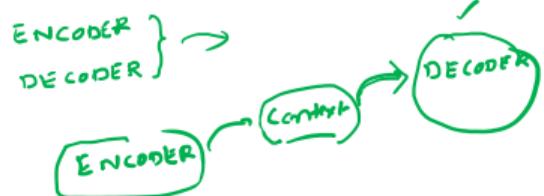
Attention Mechanism ✓

Teddy Roosevelt is the President.

Transformers



context
→ "Transformers"



Attention Is All You Need

2000 - 2014 → RNN / LSTM

2014 → "attention Mechanism"
2017 → Transformer
BERT / GPT

2017

Vision transfer
GPT2 → ChatGPT

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

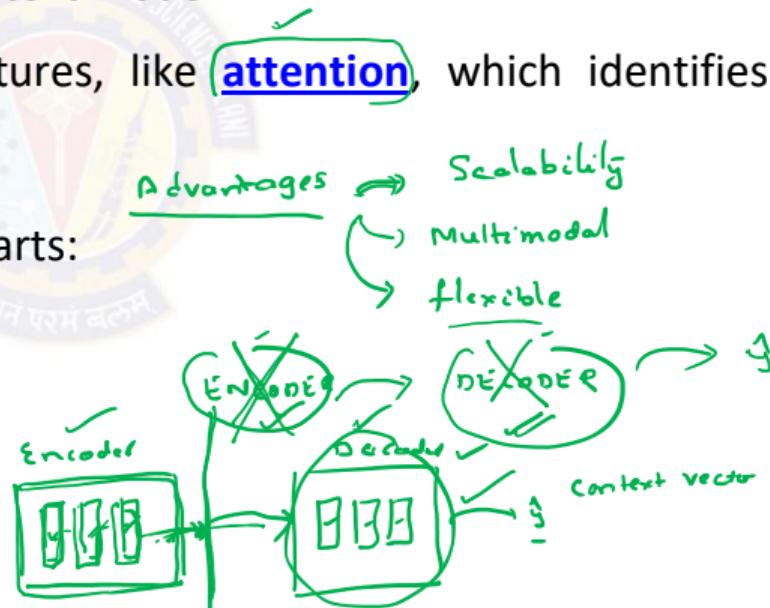
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Vision Transformers

Conv + Norm

- The **RNN** and **LSTM** neural models were designed to process language and perform tasks like classification, summarization, translation, and sentiment detection
 - RNN: Recurrent Neural Network
 - LSTM: Long Short Term Memory
- In both models, layers get the next input word and have access to some previous words, allowing it to use the word's left context
- They used word embeddings where each word was encoded as a vector of 100-300 real numbers representing its meaning

- Transformers extend this to allow the network to process a word input knowing the words in both its left and right context
- This provides a more powerful context model
- Transformers add additional features, like attention, which identifies the important words in this context
- And break the problem into two parts:
 - An encoder (e.g., Bert)
 - A decoder (e.g., GPT)



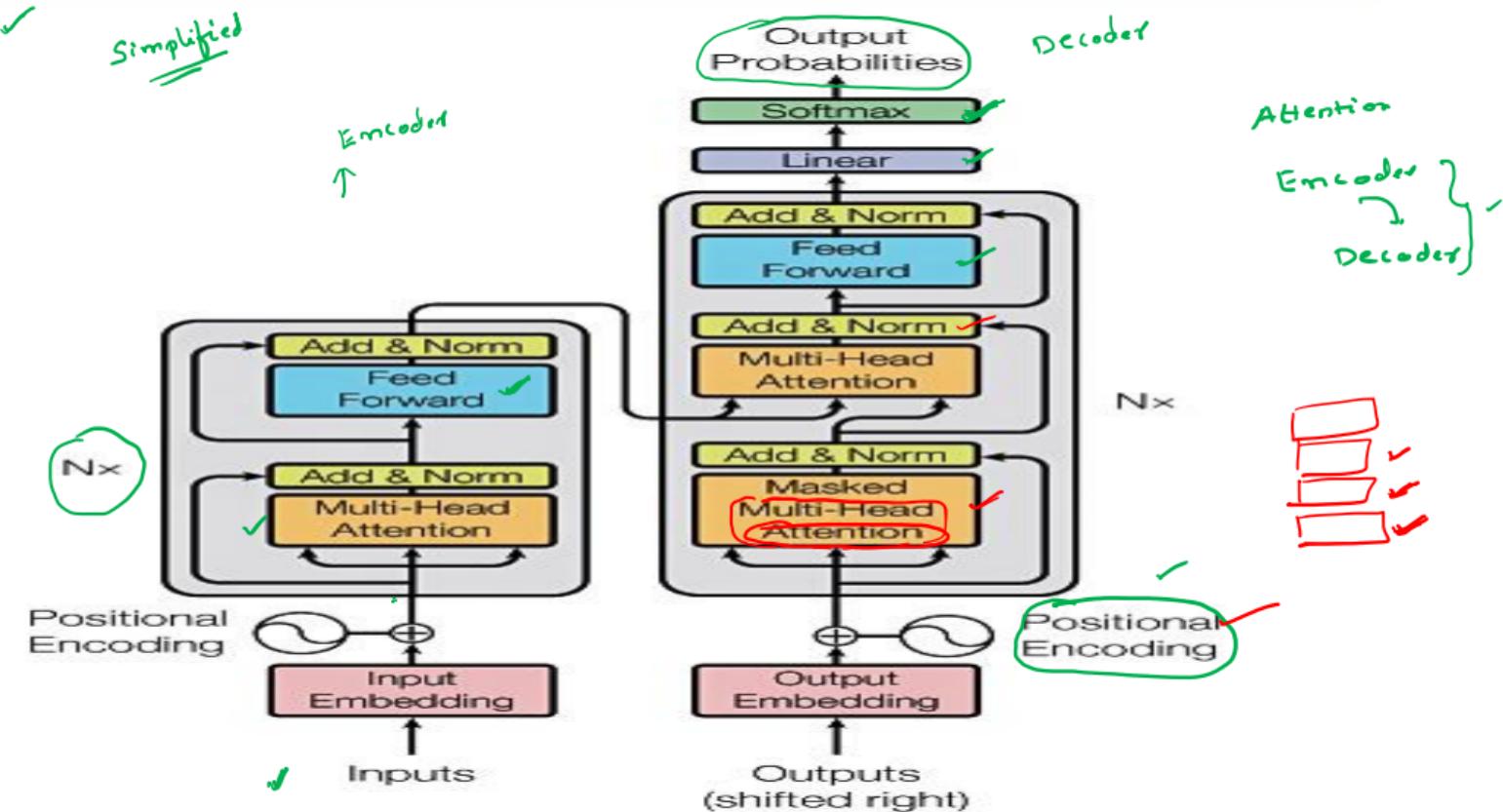
- A transformer uses an **encoder stack** to model input, and uses **decoder stack** to model output (using input information from encoder side)
- If we do not have input, we just want to model the "next word", we can get rid of the encoder side of a transformer and output "next word" one by one. This gives us **GPT**
- ✓ If we are only interested in training a language model for the input for some other tasks, then we do not need the decoder of the transformer, that gives us **BERT**

Transfer learning → "Hugging face"

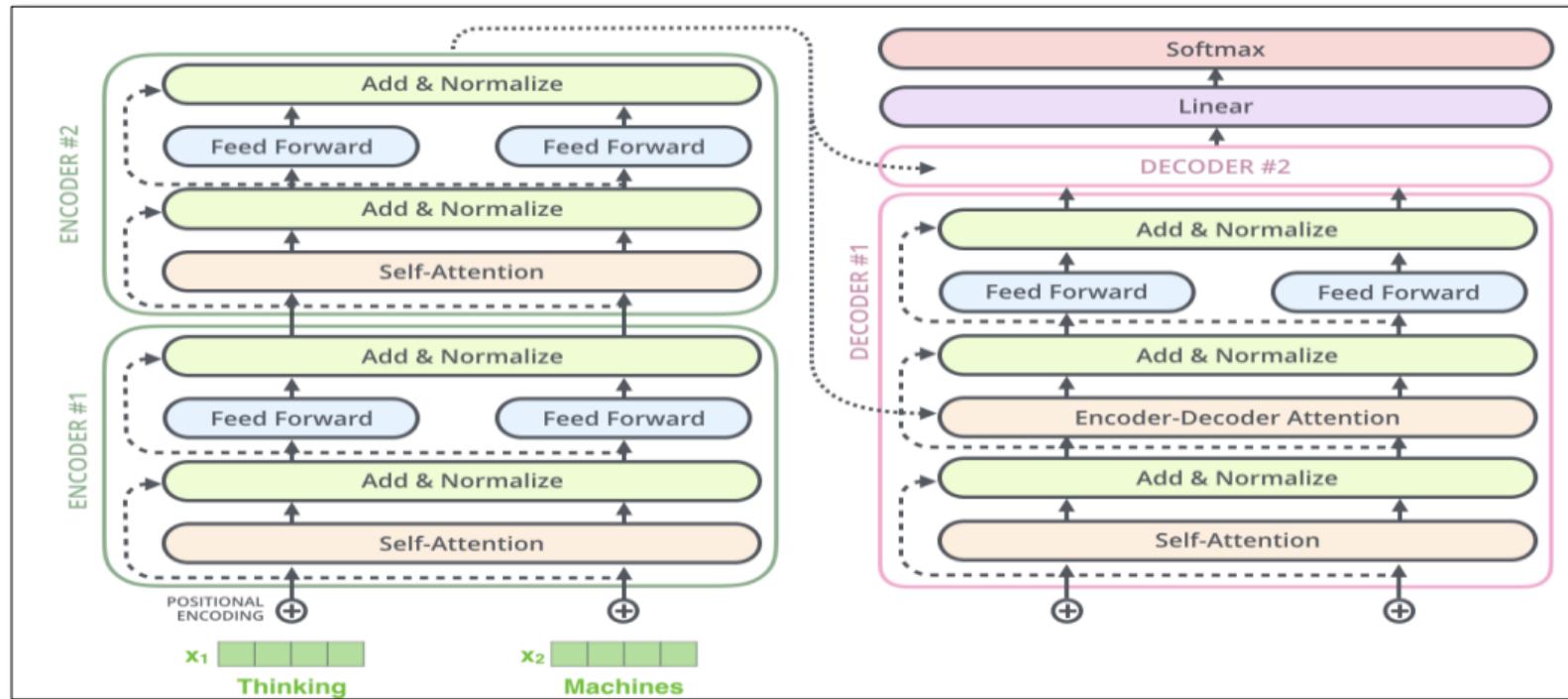
Disadvantages → ↗ interpretable,
→ computational algo



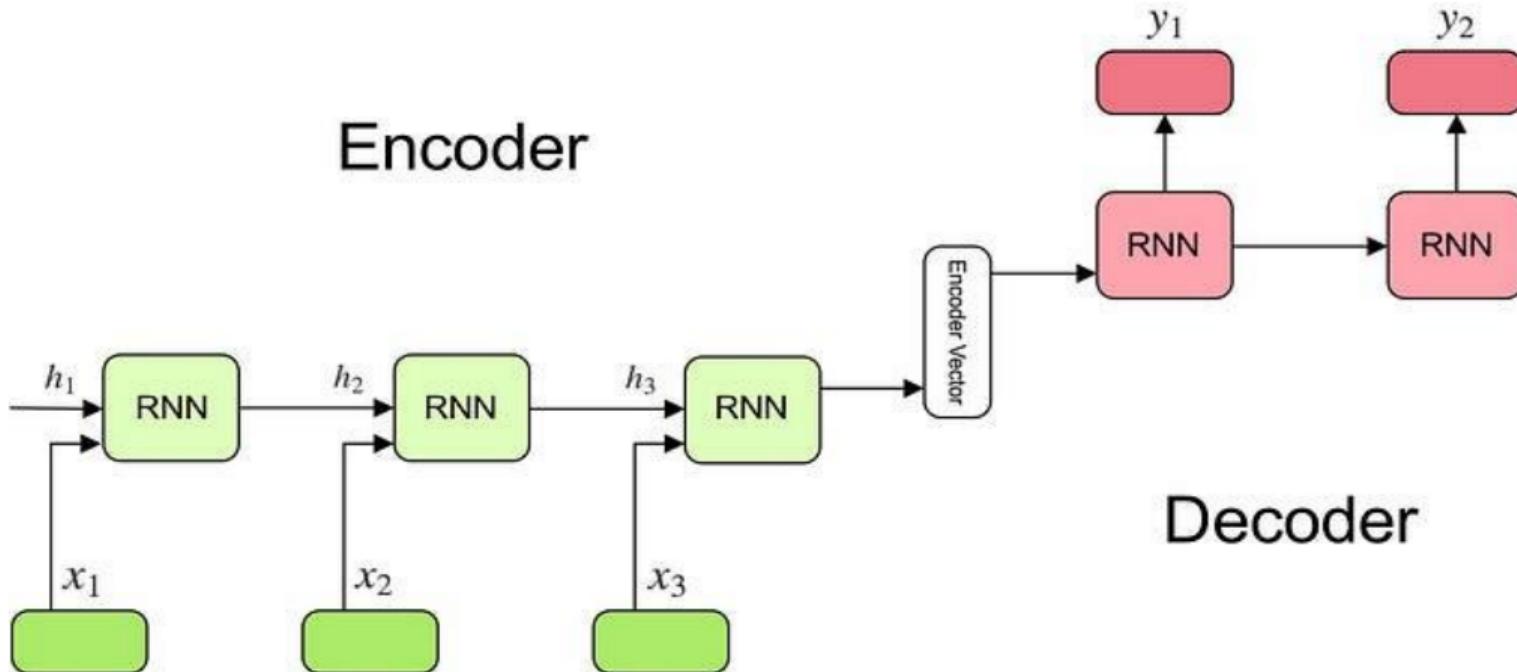
Transformer



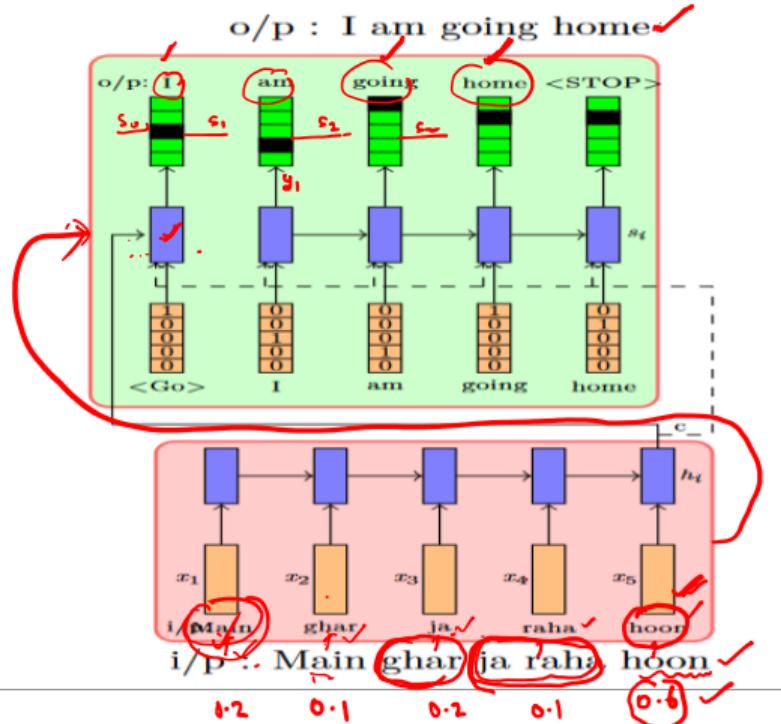
Transformer



Encoder - Decoder



Attention



✓ $t_1 : [\underline{1} \ 0 \ 0 \ 0 \ 0]$ Main

✓ $t_2 : [0 \ 0 \ 0 \ 0 \underline{1}]$ Hoon

✓ $t_3 : [0 \ 0 \underline{0.5} \ 0.5 \ 0]$ ja raha

✓ $t_4 : [0 \ 1 \ 0 \ 0 \ 0]$ ghar

i/p : Main ghar ja raha hoon

Vanilla Encoder

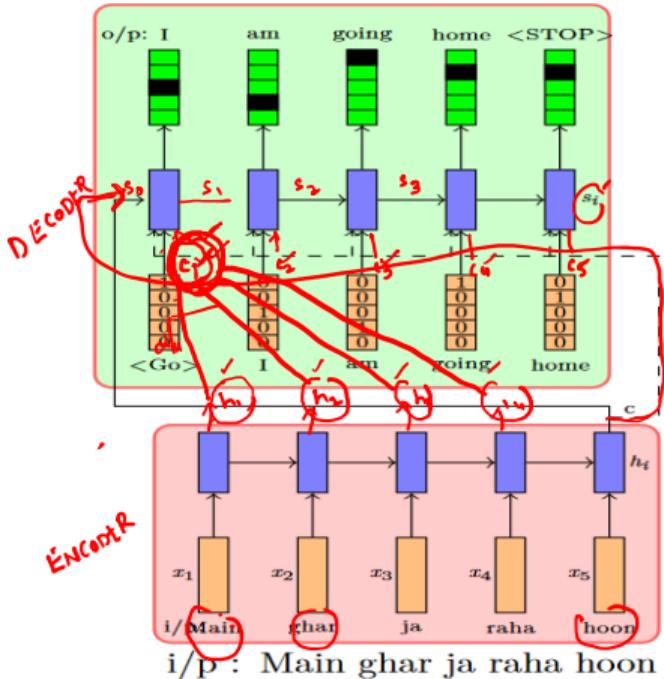
Decoder

i $[y_{i-1}, s_{i-1}]$

attention ✓
attention input
 $[y_{i-1}, s_{i-1}, c_i]$
attention
input

Attention

O/P : I am going home



$[y_{i-1}, s_{i-1}] \leftarrow \text{Encoder / Decoder}$

$[y_{i-1}, s_{i-1}, c_i] \leftarrow \text{attention mech}$

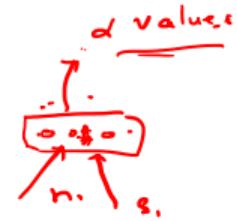
$$c_i = \sum_j \alpha_{ij} h_j$$

$$\alpha_{ij} = \sum \omega_{ij} h_j \quad \text{attention}$$

$\alpha \rightarrow \text{how to get}$

$$\alpha_{ij} = f(h_j, s_{i-1})$$

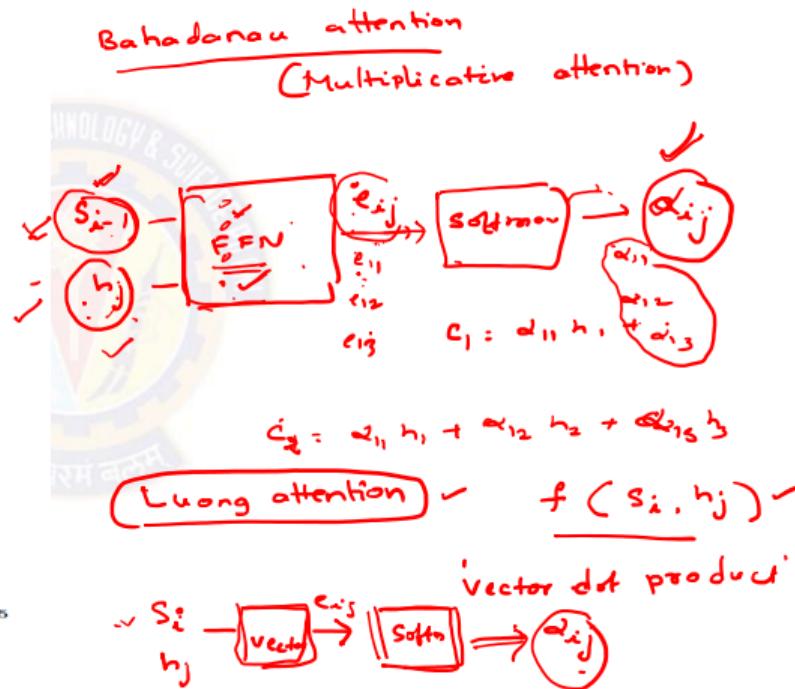
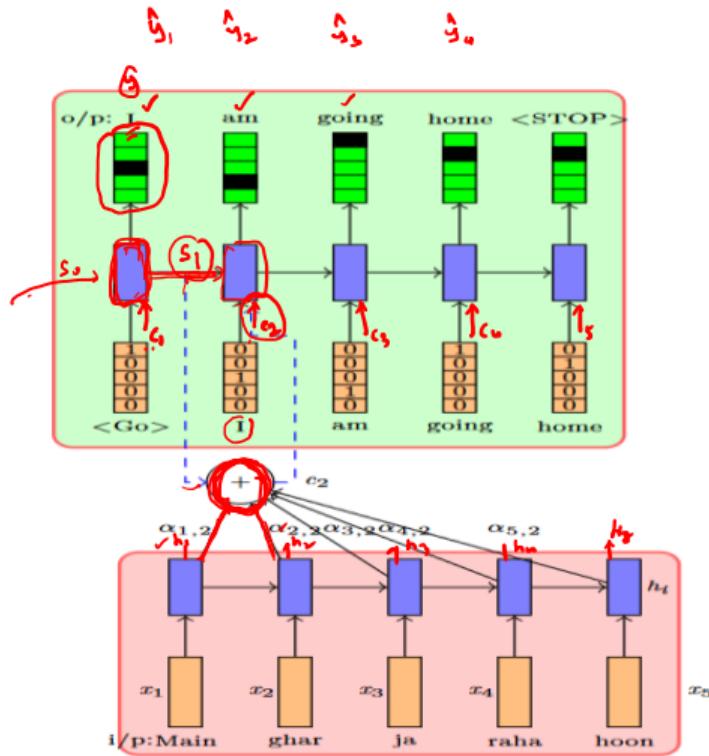
$$\alpha_{ij} = f(h_j, s_{i-1})$$

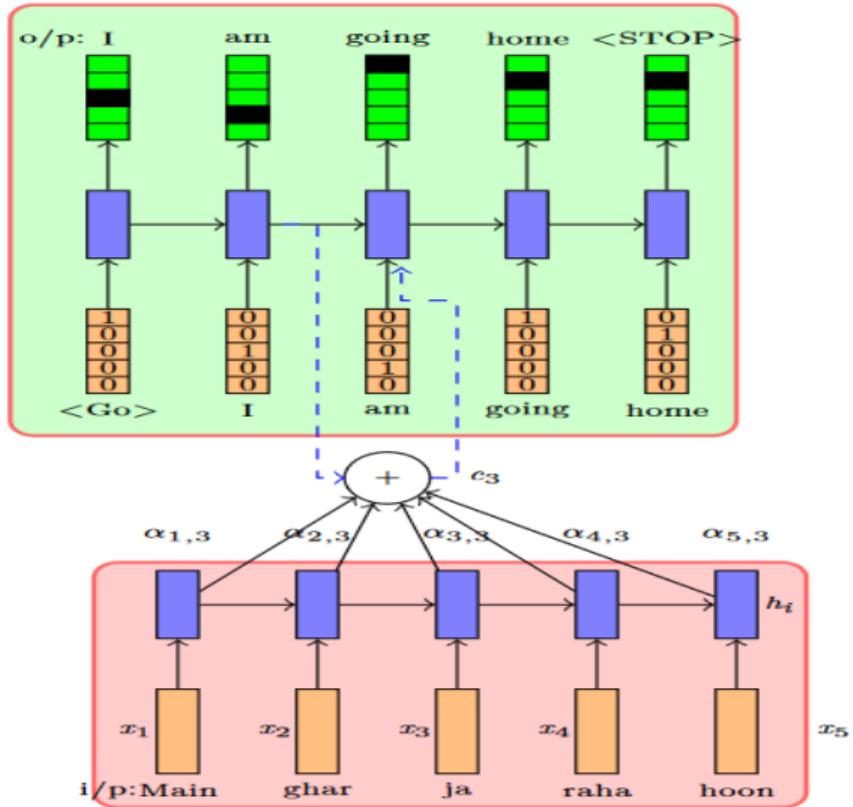


✓ Bahadaru
attention
ANN

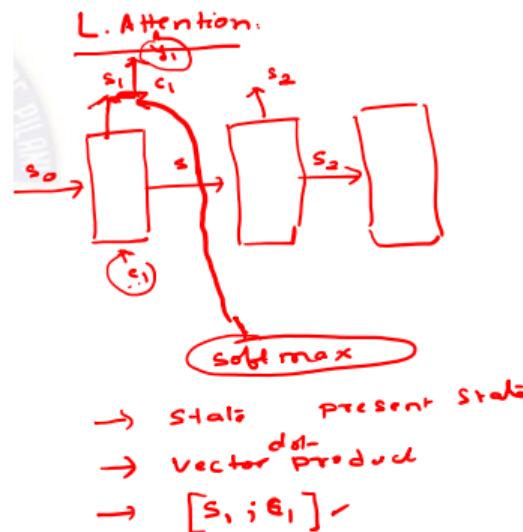
Luong attention
vector dot product

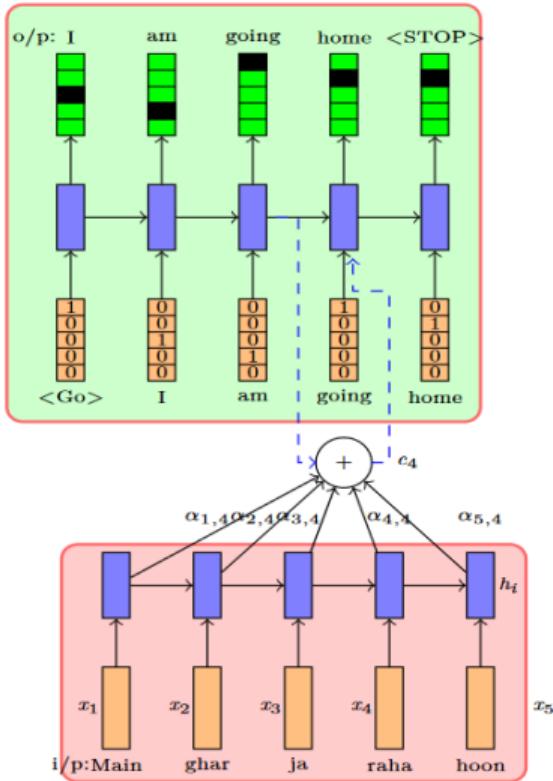
Attention



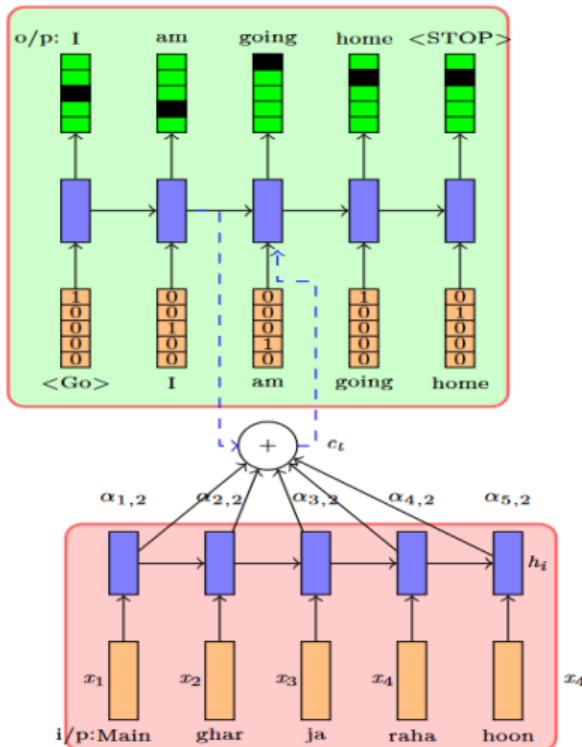


For example at timestep 3, we can just take a weighted average of the representations of 'ja' and 'raha'





- We could just take a weighted average of the corresponding word representations and feed it to the decoder
- For example at timestep 3, we can just take a weighted average of the representations of 'ja' and 'raha'
- Intuitively this should work better because we are not overloading the decoder with irrelevant information (about words that do not matter at this time step)



oracle

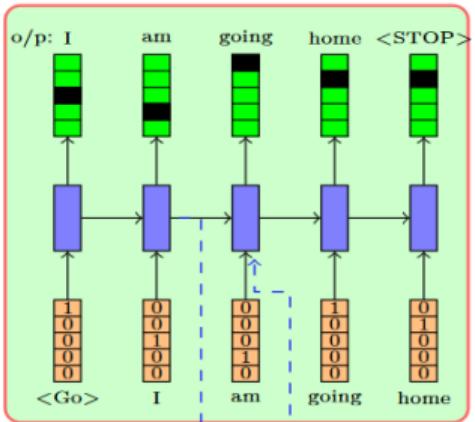
- The machine will have to learn this from the data
- To enable this we define a function

$$e_{jt} = f_{ATT}(s_{t-1}, c_j)$$

This quantity captures the importance of the j th input word for decoding the t th output word

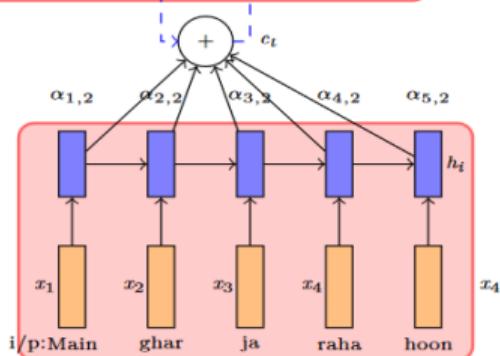
We can normalize these weights by using the softmax function

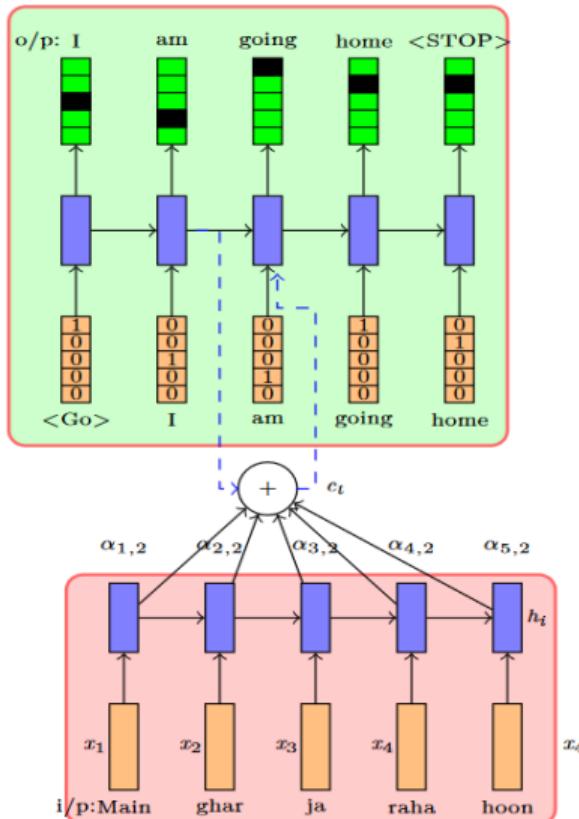
$$\alpha_{jt} = \frac{\exp(e_{jt})}{\sum_{i=1}^M \exp(e_{jt})}$$



$$\alpha_{jt} = \frac{\exp(e_{jt})}{\sum_{j=1}^M \exp(e_{jt})}$$

- α_{jt} denotes the probability of focusing on the j^{th} word to produce the t^{th} output word





This model would make a lot of sense if we were given the true α 's at training time

$$\alpha_{tj}^{true} = [0, 0, 0.5, 0.5, 0]$$

$$\alpha_{tj}^{pred} = [0.1, 0.1, 0.35, 0.35, 0.1]$$

We could then minimize $\mathcal{L}(\alpha^{true}, \alpha^{pred})$

Bahadanau Attention



Bahadanau Attention



Handwritten notes in red ink:

- Sunday: 8.10 AM to 10.45 AM
- Friday: 3.50 - 5.50
- Saturday: 10.45 - 12.45

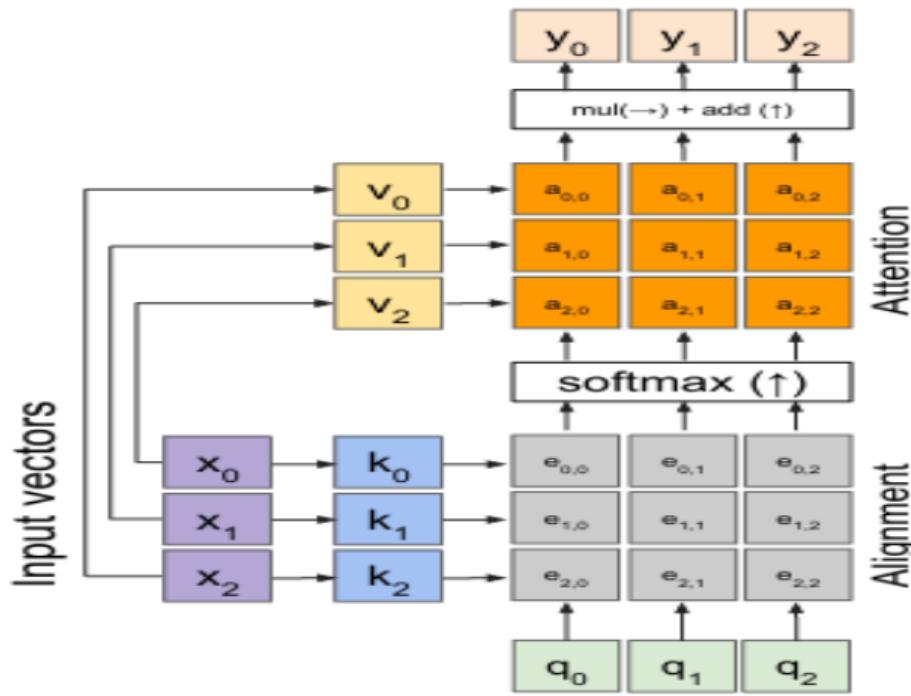
Luong Attention



Luong Attention



General attention layer



Outputs:

context vectors: \mathbf{y} (shape: D_v)

Operations:

Key vectors: $\mathbf{k} = \mathbf{x}\mathbf{W}_k$

Value vectors: $\mathbf{v} = \mathbf{x}\mathbf{W}_v$

Alignment: $\mathbf{e}_{i,j} = \mathbf{q}_j \cdot \mathbf{k}_i / \sqrt{D}$

Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$

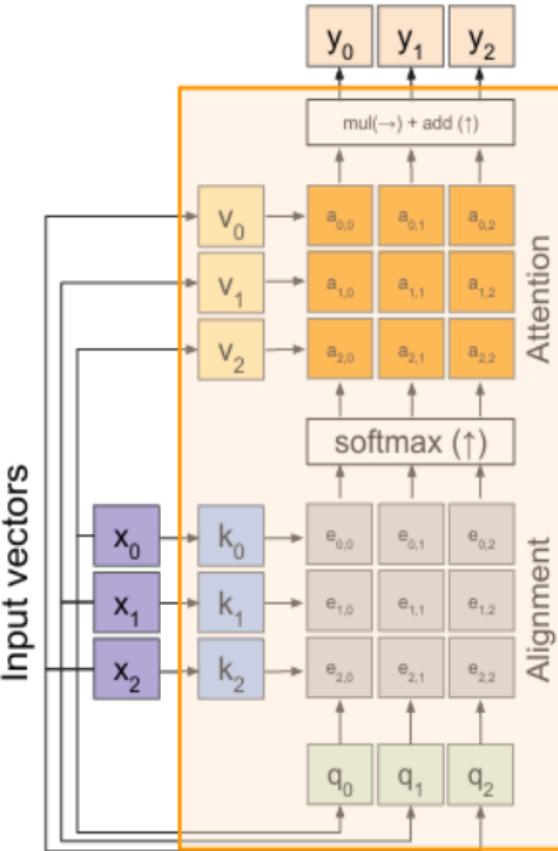
Output: $\mathbf{y}_j = \sum_i \mathbf{a}_{i,j} \mathbf{v}_i$

Inputs:

Input vectors: \mathbf{x} (shape: $N \times D$)

Queries: \mathbf{q} (shape: $M \times D_k$)

Self attention layer - attends over sets of inputs



Outputs:

context vectors: \mathbf{y} (shape: D_y)

Operations:

Key vectors: $\mathbf{k} = \mathbf{x}W_k$

Value vectors: $\mathbf{v} = \mathbf{x}W_v$

Query vectors: $\mathbf{q} = \mathbf{x}W_q$

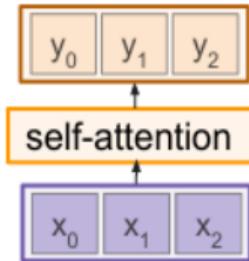
Alignment: $\mathbf{e}_{i,j} = \mathbf{q}_j \cdot \mathbf{k}_i / \sqrt{D}$

Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$

Output: $\mathbf{y}_j = \sum_i \mathbf{a}_{i,j} \mathbf{v}_i$

Inputs:

Input vectors: \mathbf{x} (shape: $N \times D$)



Hierarchical Attention

Context

U: Can you suggest a good movie?✓

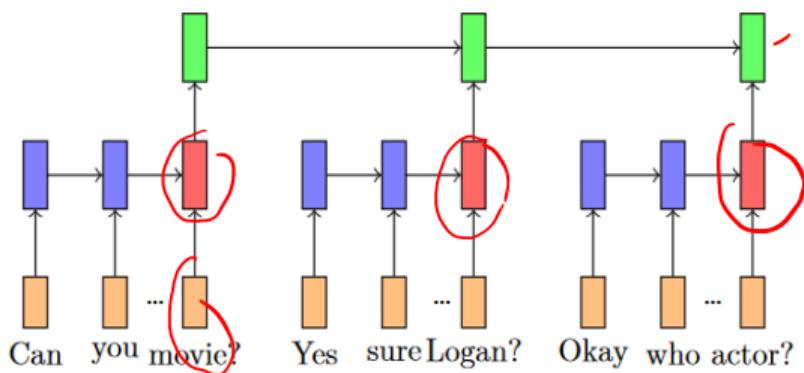
B: Yes, sure. How about Logan?✓

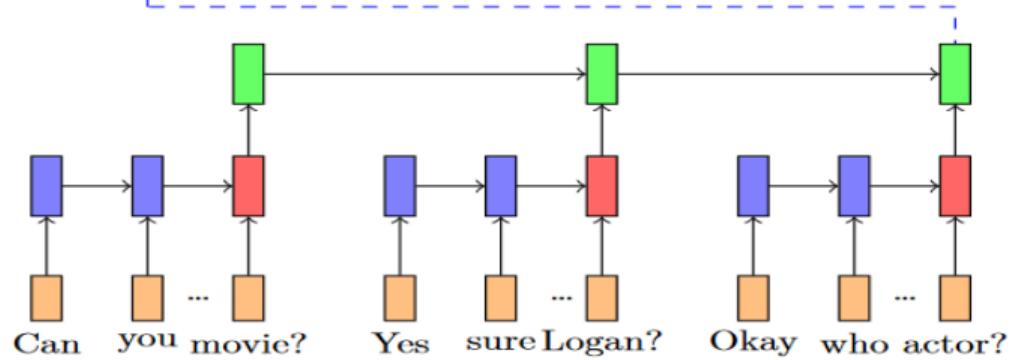
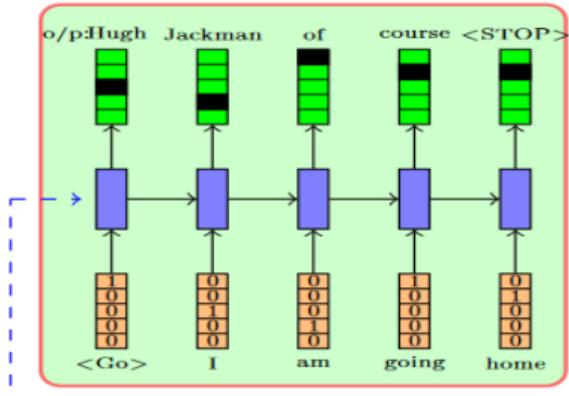
U: Okay, who is the lead actor?✓

Response

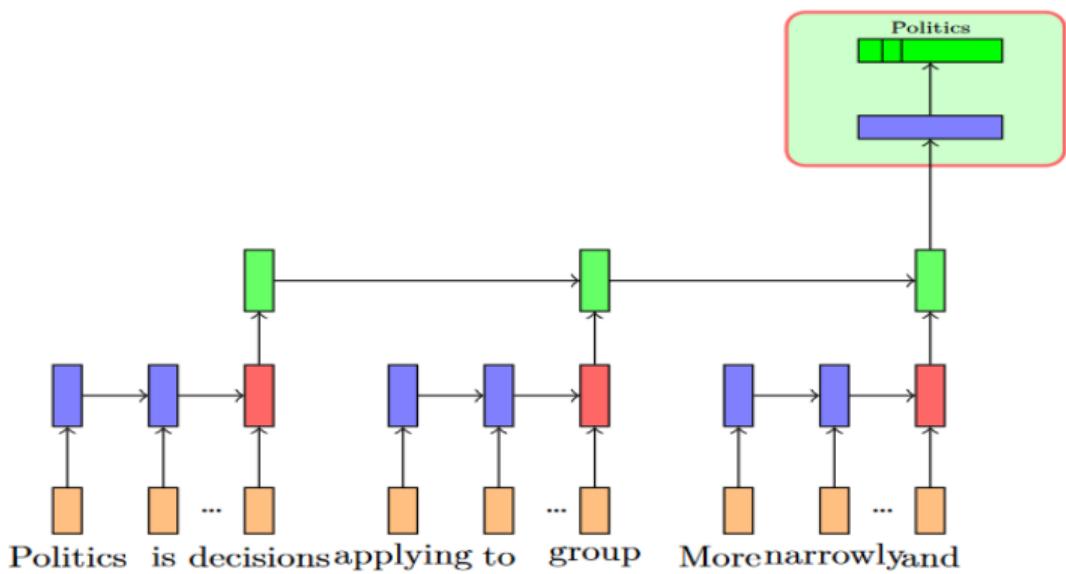
B: Hugh Jackman, of course✓

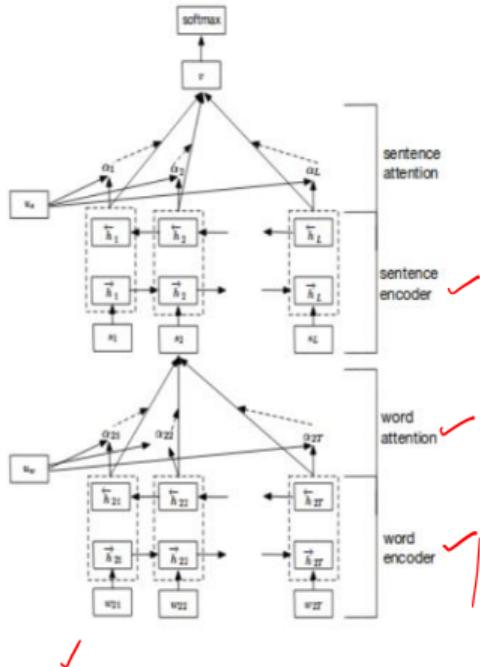
- We could think of a two level hierarchical RNN encoder
- The first level RNN operates on the sequence of words in each utterance and gives us a representation
- We now have a sequence of utterance representations (red vectors in the image)
- We can now have another RNN which encodes this sequence and gives a single representations for the sequences of utterances





Politics is the process of making decisions applying to all members of each group.
More narrowly, it refers to achieving and ...





- Data: $\{Document_i, class_i\}_{i=1}^N$
- Word level (1) encoder:

$$h_{ij} = RNN(h_{ij-1}, w_{ij})$$

$$u_{ij} = \tanh(W_w h_{ij} + b_w)$$

$$\alpha_{ij} = \frac{\exp(u_{ij}^T u_w)}{\sum_t \exp(u_{it}^T u_w)}$$

$$s_i = \sum_j \alpha_{ij} h_{ij}$$

- Sentence level (2) encoder:

$$h_i = RNN(h_{i-1}, s_i)$$

$$u_i = \tanh(W_s h_i + b_s)$$

$$\alpha_i = \frac{\exp(u_i^T u_s)}{\sum_i \exp(u_i^T u_s)}$$

$$s = \sum_i \alpha_i h_i$$

Figure: Hierarchical Attention Network
[Yang et al.]