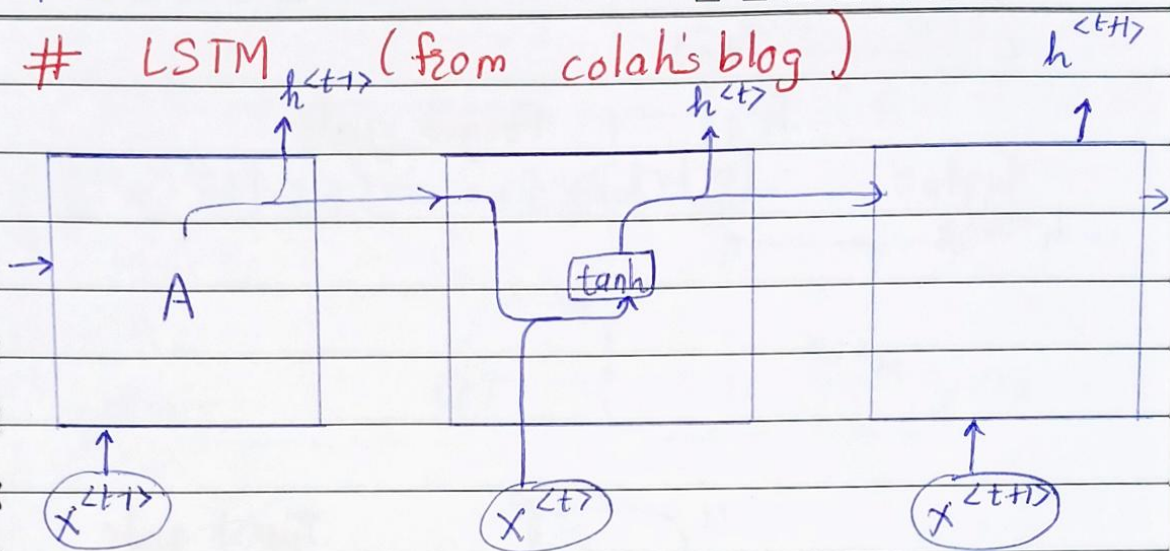


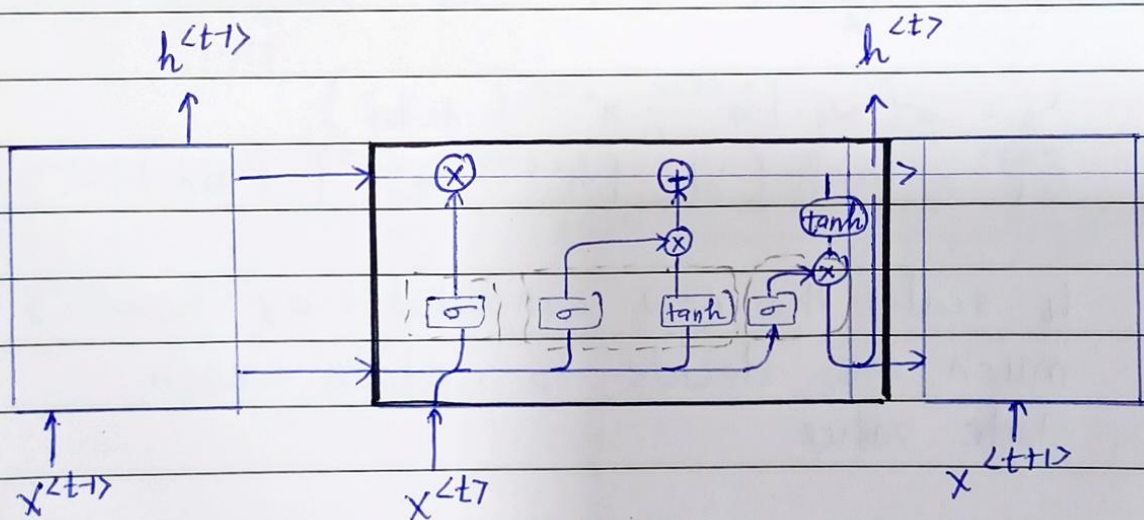
Basic RNN suffers from short-term memory problem

Date: 12/8/24
 MON TUE WED THU FRI SAT SUN
☒ ☐ ☐ ☐ ☐ ☐ ☐

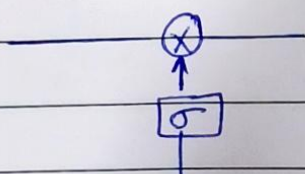
LSTM (from colah's blog)



Repeating module in a standard RNN contains single layer



each = neural net layer

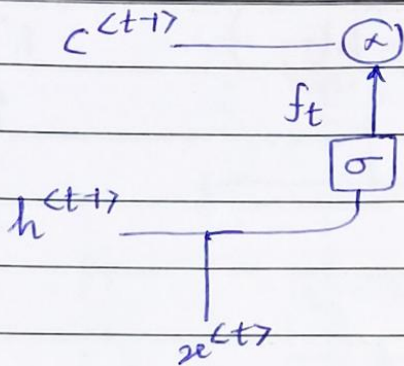


= gate
 sig \rightarrow how much of each component shd be let thru.

Date :

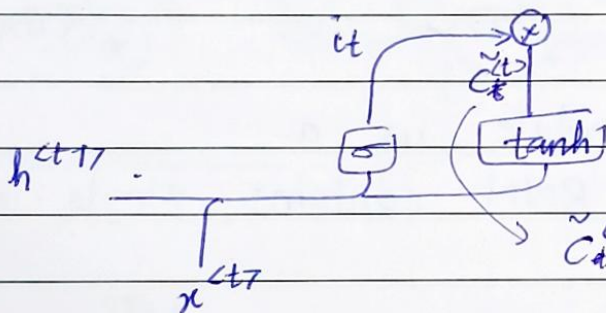
MON TUE WED THU FRI SAT SUN

☐ ☐ ☐ ☐ ☐ ☐ ☐



Forget gate

$$f_t = \sigma(W_f \cdot [h^{(t-1)}, x^{(t)}] + b_f)$$



Input gate

$\tilde{c}_t^{(t)}$ = vector of new candidate values

$$i_t = \sigma(W_i \cdot [h^{(t-1)}, x^{(t)}] + b_i)$$

$$\tilde{c}_t^{(t)} = \tanh(W_c \cdot [h^{(t-1)}, x^{(t)}] + b_c)$$

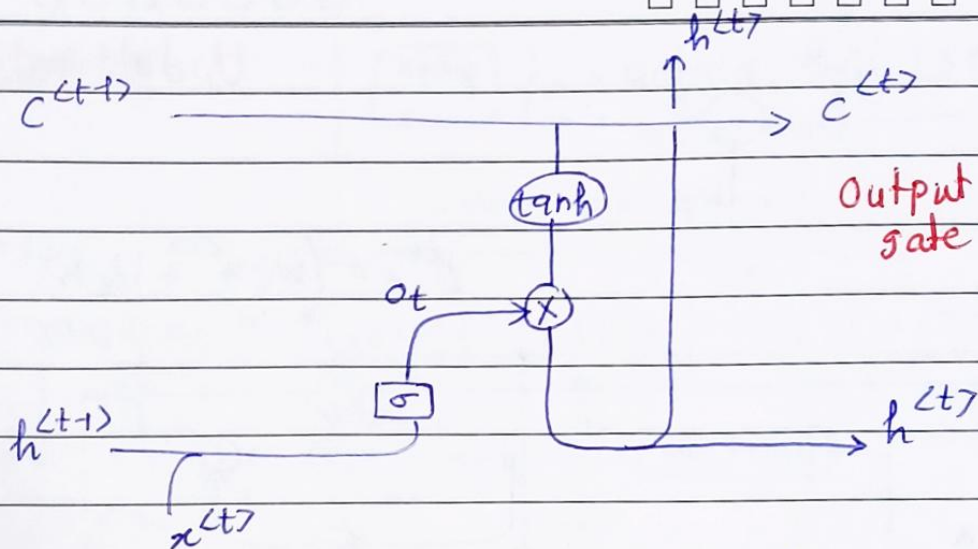
i_t scales the new candidate by how much we decide to update each state value.

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

Date :

MON TUE WED THU FRI SAT SUN

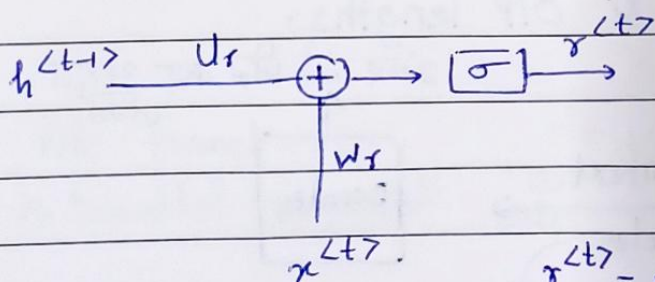
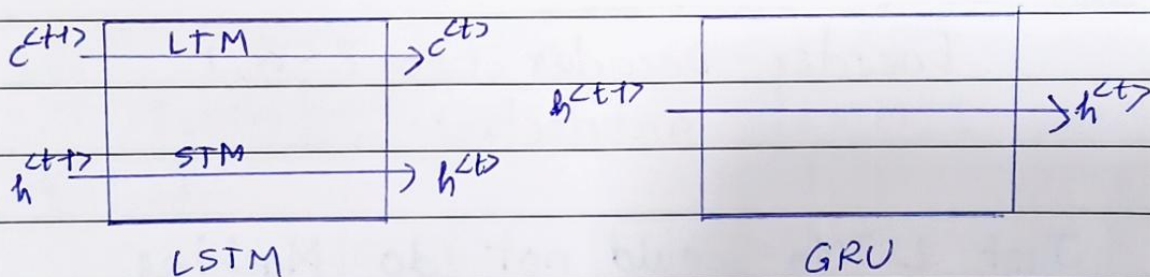
☐ ☐ ☐ ☐ ☐ ☐ ☐



$$o_t = \sigma(W_o [h^{(t-1)}, x^{(t)}] + b_o)$$

$$h^{(t)} = o_t * \tanh(c_t)$$

GRU



Reset gate

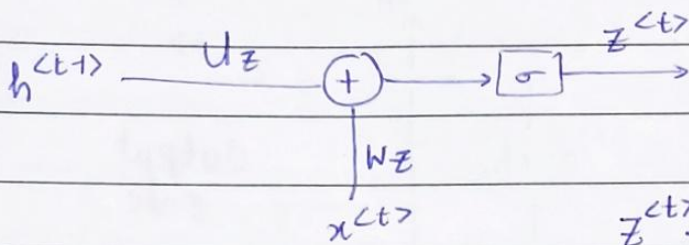
$$r^{(t)} = \sigma(W_r x^{(t)} + U_r h^{(t-1)})$$

Date :

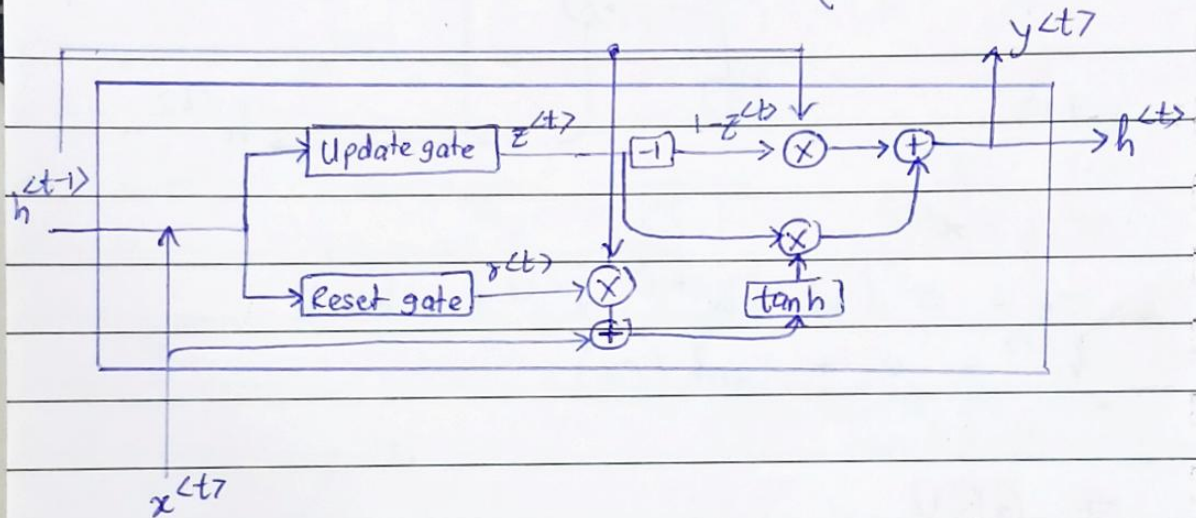
MON TUE WED THU FRI SAT SUN

☐ ☐ ☐ ☐ ☐ ☐ ☐

Update gate



$$z^{(t)} = \sigma(W_Z x^{(t)} + U_Z h^{(t-1)})$$

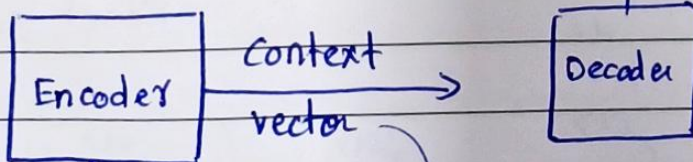


Encoder-Decoder (seq2seq) architecture

Just LSTMs could not do Machine translation bcz they cannot handle variable I/P & O/P lengths.

आप से मिल कर अच्छा लगता


=

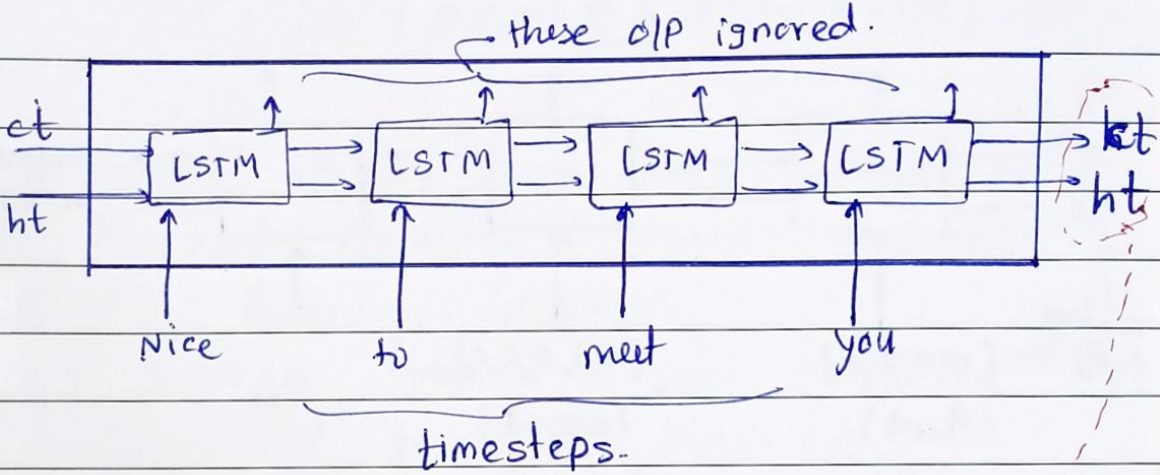


Nice to meet you

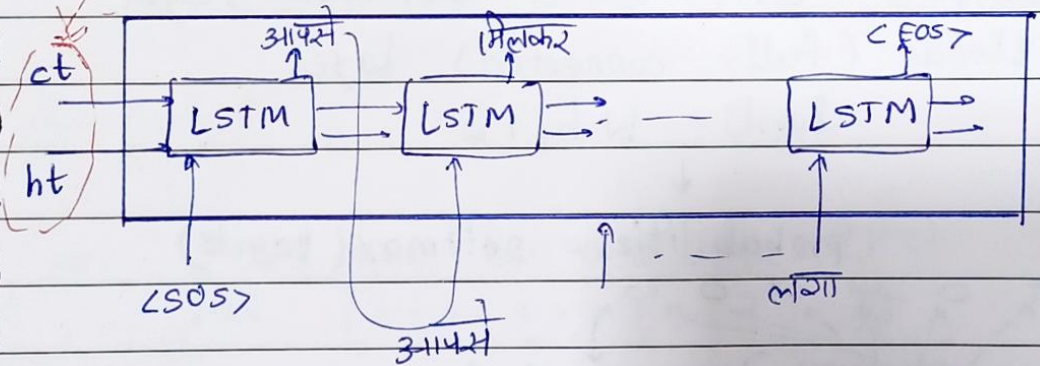
stores the summary of I/P.

Encoder:

 → unfold this LSTM over time



Decoder



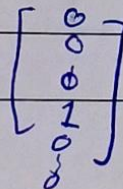
M/c transl:

1) Parallel Dataset

	ENG	HIN
①	—	①
②	—	②

2) tokenization.

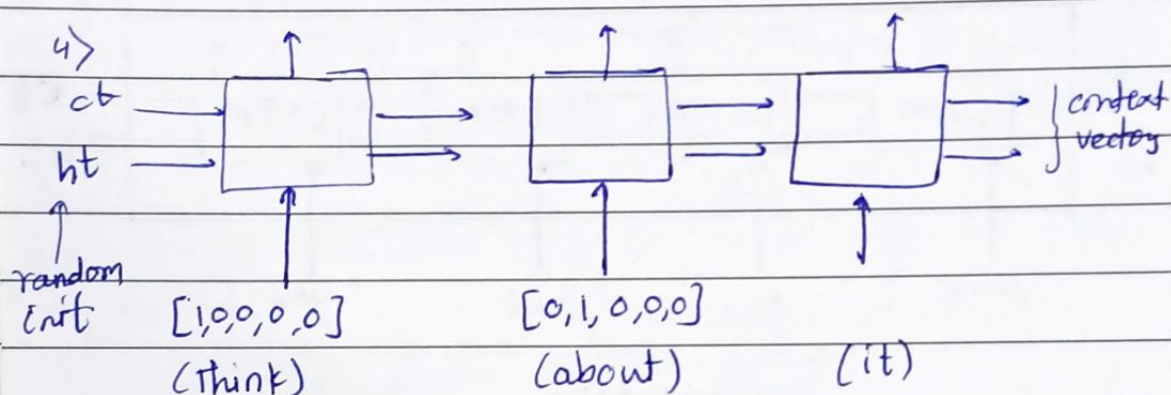
3) One-hot enc.



Date :

MON TUE WED THU FRI SAT SUN
☐ ☐ ☐ ☐ ☐ ☐ ☐

For Hindi vocab: add <SOS> & <Eos> tokens too

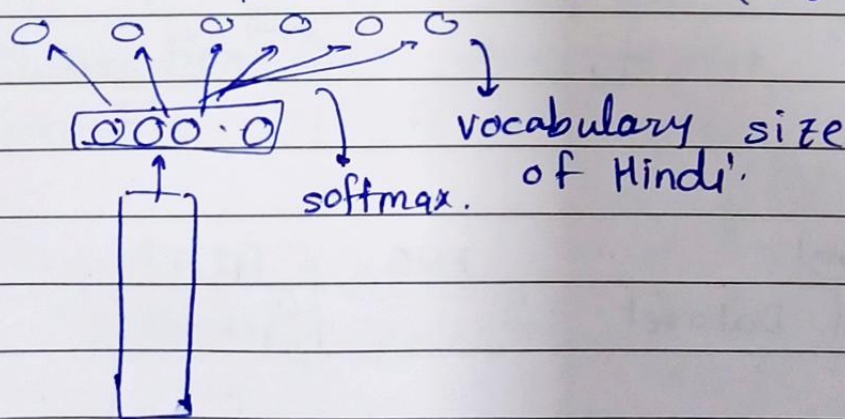


5) To output vectors (1-hot wale) from the decoder, we pass the LSTM's OLP to a ~~softmax~~ layer. Linear (Fully connected) layer.

$$\text{logits} = W \cdot h_t + b$$



$$\text{probabilities} = \text{softmax}(\text{logits})$$



Date :

MON TUE WED THU FRI SAT SUN
☐ ☐ ☐ ☐ ☐ ☐ ☐

y_{true} [0, 1, 0, 0, 0, 0, 0] सही

for y_{pred} [0.2, 0.1, 0.3, 0.2, 0.1, 0.1, 0] ली

6> Now, if we send this wrong O/P ली to next timestep, it will decrease the learning process.

∴ During Training, we supply सही (correct) to next timestep.
 this is called "Teacher Forcing"

Now, suppose O/Ps from decoder were:

ली सही $\langle \text{EOS} \rangle \rightarrow$ Forward Prop

7> Loss

eg. y_{true} [0, 1, 0, 0, 0, 0, 0] सही
 y_{pred} [0.2, 0.1, 0.3, 0.2, ...] ली

Loss = Categorical Cross Entropy.
 n = Hindi vocab size

$$L = - \sum_{i=1}^n y_i^{true} \log(y_i^{pred})$$

$$L_{t=1} = - 1 \times \log(0.1) = 1$$

Simply calc. losses for all timesteps

8) Gradient

calc. the gradient of loss wrt each of the trainable parameters (LSTM, softmax, linear layer)

$$\frac{\partial \text{loss}}{\partial \text{logits}_t} = \hat{y}_t - y_t$$

For each hidden state h_t ,

$$\frac{\partial \text{loss}_t}{\partial h_t} = \frac{\partial \text{loss}_t}{\partial \text{logits}_t} \cdot \frac{\partial \text{logits}_t}{\partial h_t}$$

$$\frac{\partial \text{loss}}{\partial W_h} = \sum_{t=1}^T \frac{\partial \text{loss}_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_h}$$

Backprop thru encoder

$$\frac{\partial \text{loss}}{\partial h_{enc}} = \sum_{t=1}^{T_{dec}} \frac{\partial \text{loss}_t}{\partial h_{dec,t}} \cdot \frac{\partial h_{dec,t}}{\partial h_{enc}}$$

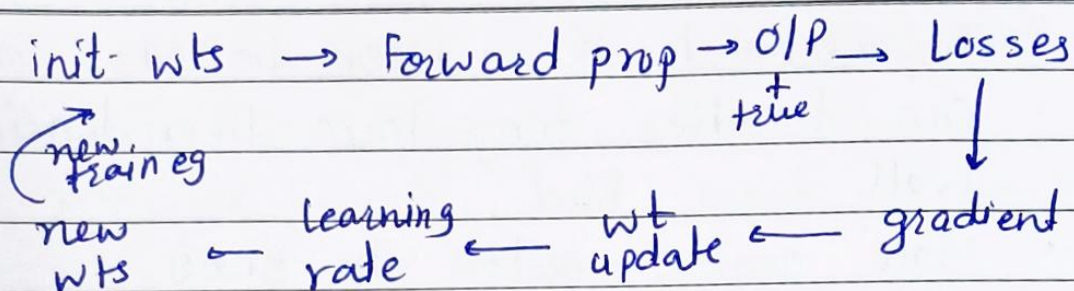
9) Update params

$$W_{new} = W_{old} - \alpha \frac{\partial \text{loss}}{\partial W}$$

Date :

MON TUE WED THU FRI SAT SUN

☐ ☐ ☐ ☐ ☐ ☐ ☐



10) Prediction

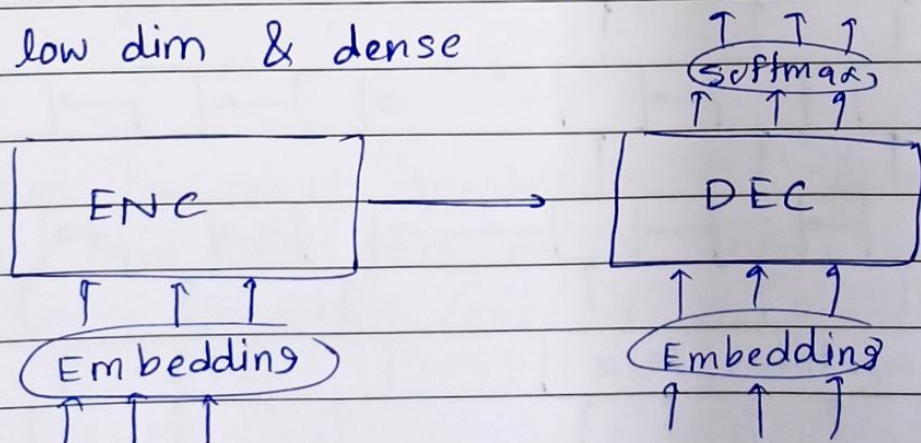
When the wts are freezed, we don't do teacher forcing & supply a decoder LSTM o/p as I/P for next timestep

Improvement 1

Embedding Layer \rightarrow for both ENC & DEC

\downarrow
it contains a summary for each word

- low dim & dense

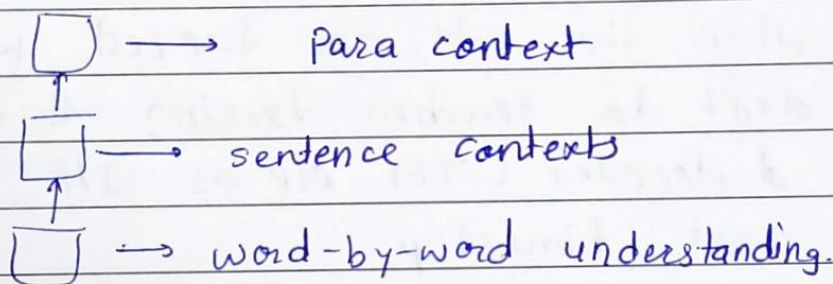


Date :

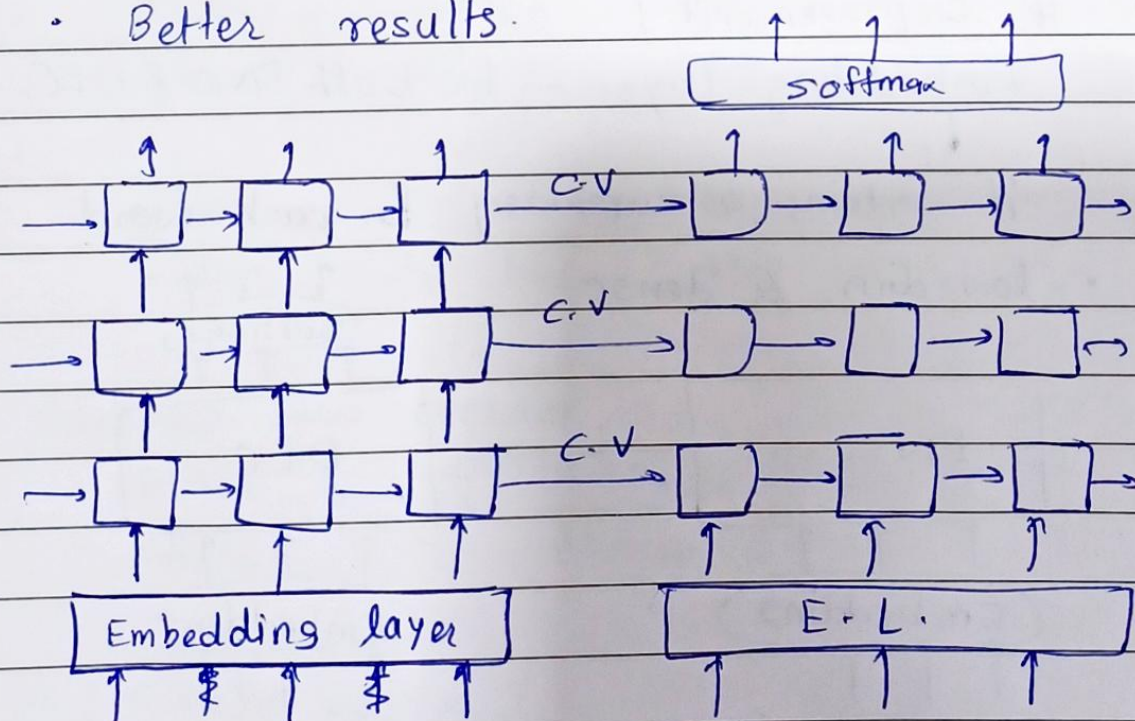
MON TUE WED THU FRI SAT SUN
☐ ☐ ☐ ☐ ☐ ☐ ☐

Improvement 2 - Deep LSTMs

- Can handle Long Term Dependencies well.
- More context vectors to store summaries.
- Can understand hierarchical levels -



- Better results.



3 layer LSTM

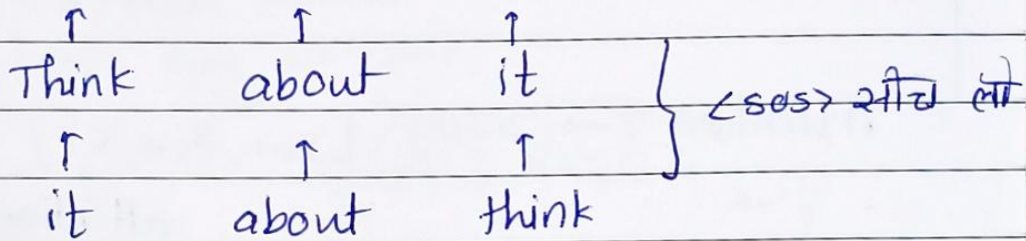
Date :

MON TUE WED THU FRI SAT SUN

☐ ☐ ☐ ☐ ☐ ☐ ☐

Improvement 3 - Reversing the I/P.

- Works for languages where the starting words hold much of context



#

ATTENTION

#

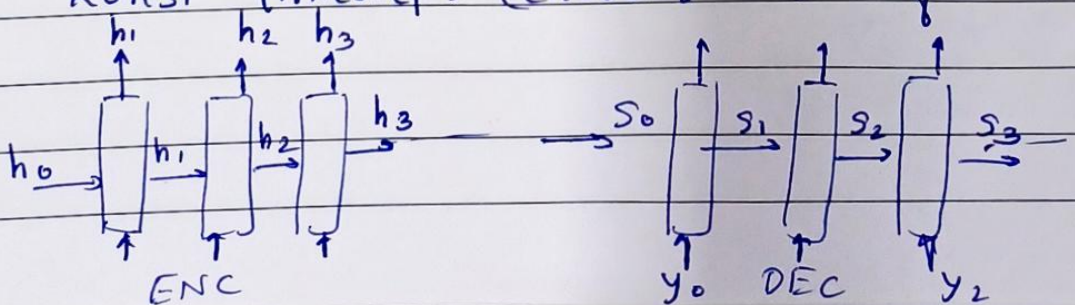
Problems with Encoder-Decoder:

- 1) Cannot summarize sentences with $\sim > 25$ words.
- 2) Turn off the lights

↳ <sos> मिस्ट जय श्री <Eos>

We do not need the whole I/P sentence to generate only मिस्ट

Thus for every timestep in the decoder, we have to pass on an info as to 'konsi timesteps (encoder ki) useful hai'.



Date :

MON TUE WED THU FRI SAT SUN
☐ ☐ ☐ ☐ ☐ ☐ ☐

Vanilla enc-dec

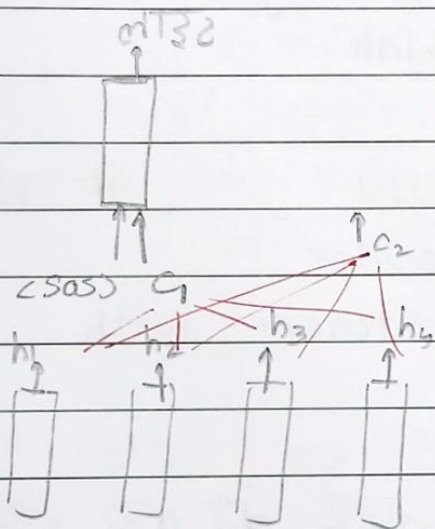


for every timestep of DEC → 2 I/Ps

$[y_{i-1} \& s_{i-1}]$

Attention → 3 I/Ps $[y_{i-1}, s_{i-1}, c_i]$

↑
attention I/P.



For each timestep in DEC, a $w_t(\alpha_i)$ is assigned for each of the ENC timestep.

$$c_1 = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3 + \alpha_4 h_4$$

↑
scalars vectors (4 dim here)

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

↑ ↑
i, j
timestep of DEC " ENC

Date :

MON TUE WED THU FRI SAT SUN

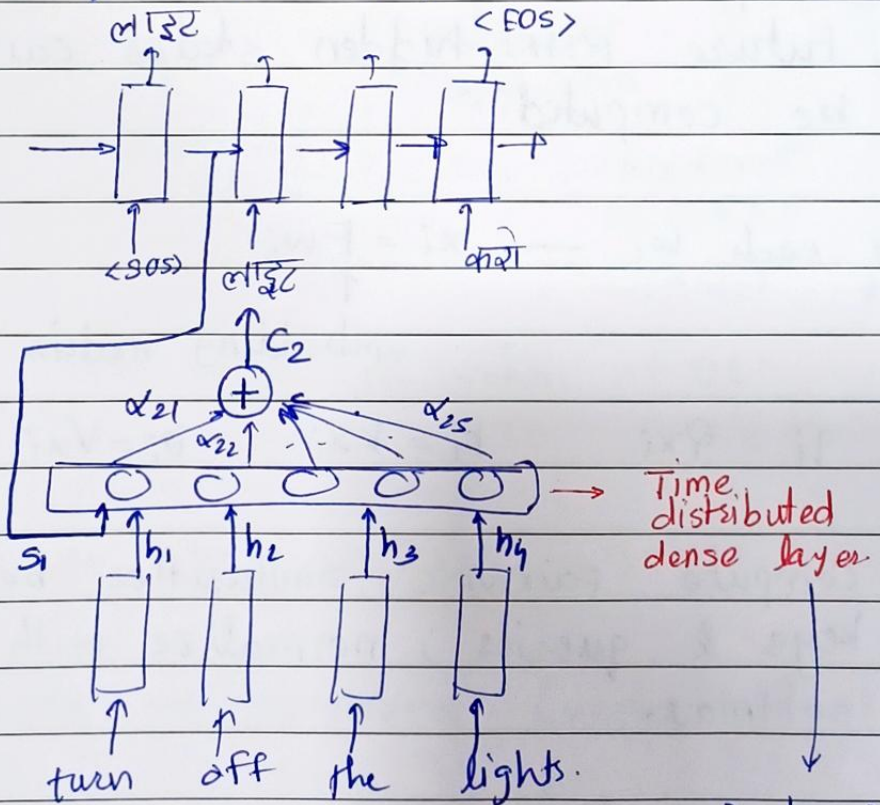
☐ ☐ ☐ ☐ ☐ ☐ ☐

α_{21} depends on h_1 as well as s_1

this is bcz:
already kitna translⁿ ho chuka hai is IMP. ← prev hidden state of DEC.

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

instead of finding any function here, we use ANN



Time distributed dense layer.

Backprop

ke upar bhi me ANN wts bhi adjust hote hai

→ Lecture 8 - Self Attn & Transf ←

- * Lack of parallelizability: in recurrent models!
- For & back prop have $O(\text{seq. length})$ unparallelizable ops.
- GPU can perform independent comps.
- Future RNN hidden states cannot be computed

For each $w_i \rightarrow x_i = E w_i$
 ↑
 embedding matrix

$$1) \quad q_i = Q x_i \quad k_i = K x_i \quad v_i = V x_i$$

2) Compute pairwise similarities betⁿ keys & queries; normalize with softmax.

$$e_{ij} = q_i^T k_j$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_j \exp(e_{ij})}$$

Date :

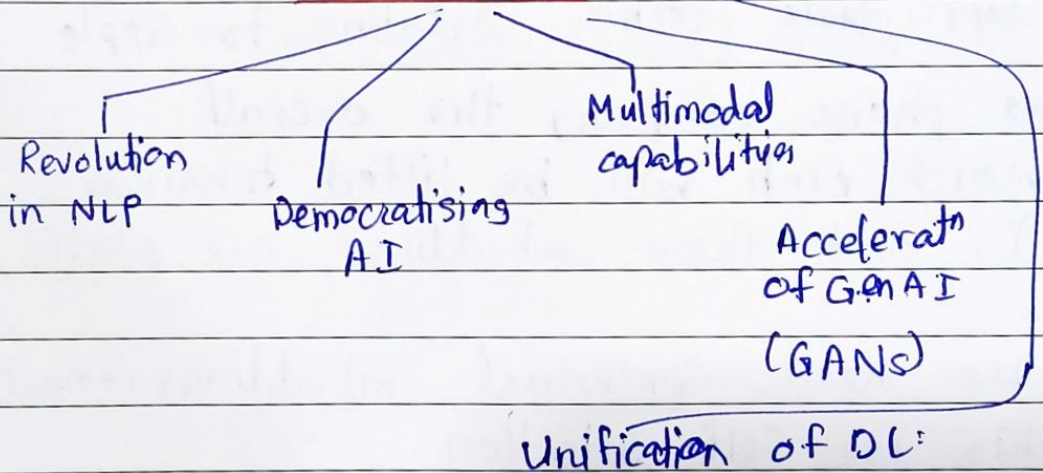
MON TUE WED THU FRI SAT SUN

☐ ☐ ☐ ☐ ☐ ☐ ☐

3) Compute o/p for each word as weighted sum of values.

$$a_i = \sum_j \alpha_{ij} v_j$$

Transformers



Self - Attention

words \rightarrow numbers (vectorization)

1) One Hot Encoding

2) Bag of words

	mat	rat	cat
s1	2	0	1
s2	0	2	1

3) Word embeddings - semantic meaning

this does not represent the "meaning"
but "avg meaning" of the word

eg. $[x, y]$ If there are
 $\uparrow \quad \uparrow$ 9000 ~~words~~ sentences
 repr: taste tech. relating to Apple
 as phone company the overall
 word emb will be tilted towards
 Y. And these embeddings are static

\therefore We need Contextual embeddings (smart)
 This is Self-Attention.

eg. Money bank grows
 River bank flows

$$\text{money} = 0.7 \text{money} + 0.2 \text{bank} + 0.1 \text{grows}$$

$$\text{bank} = 0.25 \text{money} + 0.7 \text{bank} + 0.05 \text{grows}$$

these are individual old
 embedding (n dim)

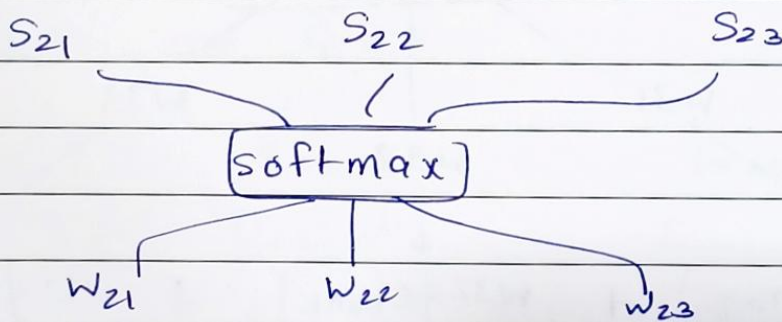
Date :

MON TUE WED THU FRI SAT SUN
☐ ☐ ☐ ☐ ☐ ☐ ☐

* Similarity between 2 vectors is given by dot product

∴ The Ws in prev eqns are similarities.

$$e_{bank}^{(new)} = \overbrace{[e_{bank} \cdot e_{money}]}^{S_{21}} e_{money} + \overbrace{[e_{bank} \cdot e_{bank}]}^{S_{22}} e_{bank} + \overbrace{[e_{bank} \cdot e_{grows}]}^{S_{23}} e_{grows}$$



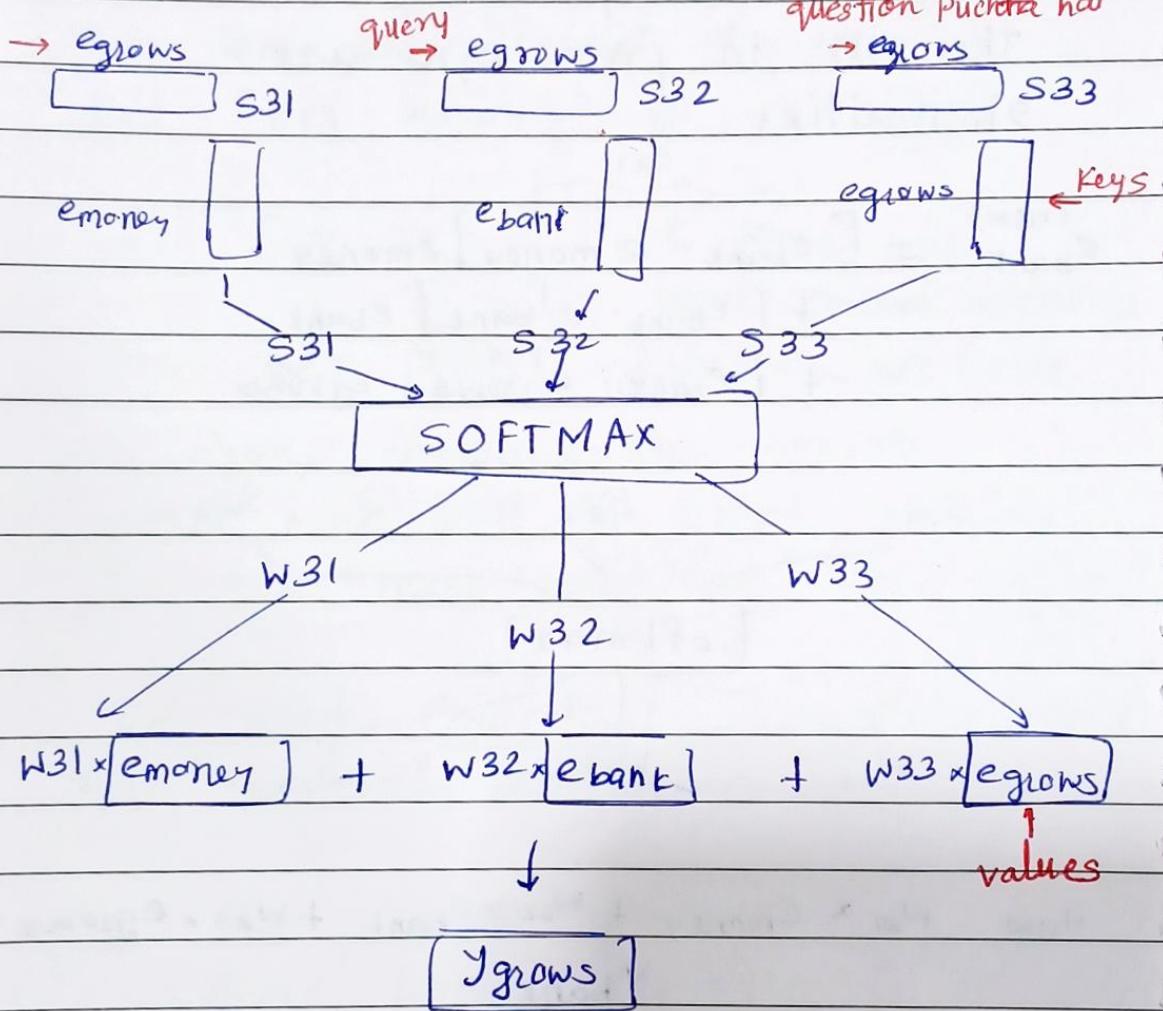
$$\text{Now } W_{21} \times e_{money} + W_{22} \times e_{bank} + W_{23} \times e_{grows} = y_{bank}$$

↑
contextual embedding.

Date :

MON TUE WED THU FRI SAT SUN
☐ ☐ ☐ ☐ ☐ ☐ ☐

money bank grows
 ↓ ↓ ↓
 [] [] []



→ All are parallel operations - (for GPU)

→ No Learning Parameters involved

Model ko data se pharak nahi padta.
 Just depends on individual sentence.

eg. 'piece of cake' → कक की दुकान
 ypred → बहुत आश्चर्य की है
 ytrue

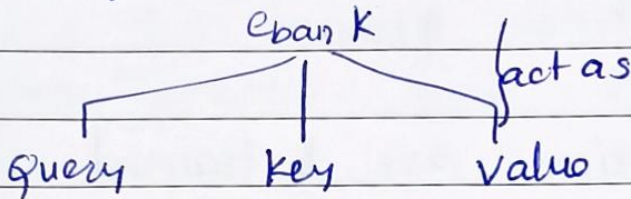
Date :

MON TUE WED THU FRI SAT SUN

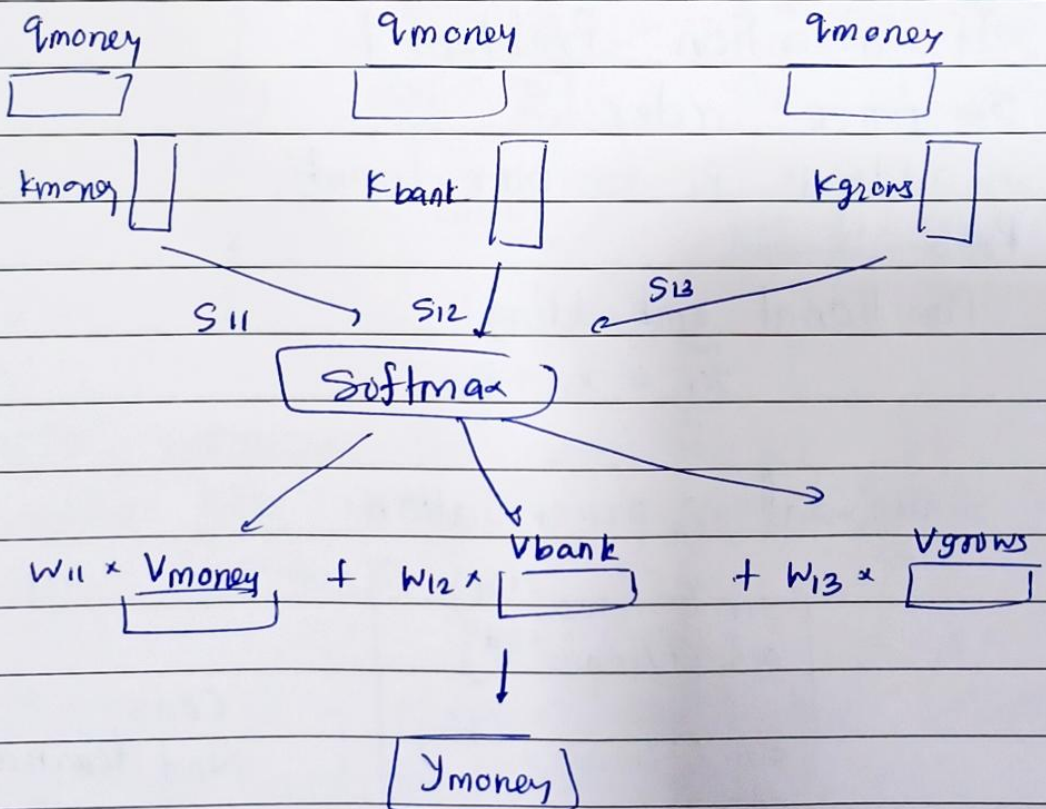
☐ ☐ ☐ ☐ ☐ ☐ ☐

Thus we need **Task-specific** contextual emb. and not general con. emb.

Query, key, Value



Now, we want to make out 3 diff vectors from single e_bank to function as Q, k, V



Date :

MON TUE WED THU FRI SAT SUN

☐ ☐ ☐ ☐ ☐ ☐ ☐

To generate these 3 different vectors,
we take 3 matrices W_q, W_k, W_v

$$e_{bank} \times W_q^T = q_{bank}$$

These matrices are trained with
data using backpropagation.

Stanford Lec 8:

Self-attention Problem 1:

Sequence order:

adding P_i to our inputs.

~~Prob 2: seq~~

Positional embedding:

$$\tilde{x}_i = x_i + P_i$$

→ Sinusoidal representation:

$$P_i = \begin{bmatrix} \sin(i/10000^{2/d}) \\ \cos(i/10000^{2/d}) \\ \sin(i/10000^{2/d/2}) \\ \cos(i/10000^{2/d/2}) \end{bmatrix}$$

Cons:
Not learnable