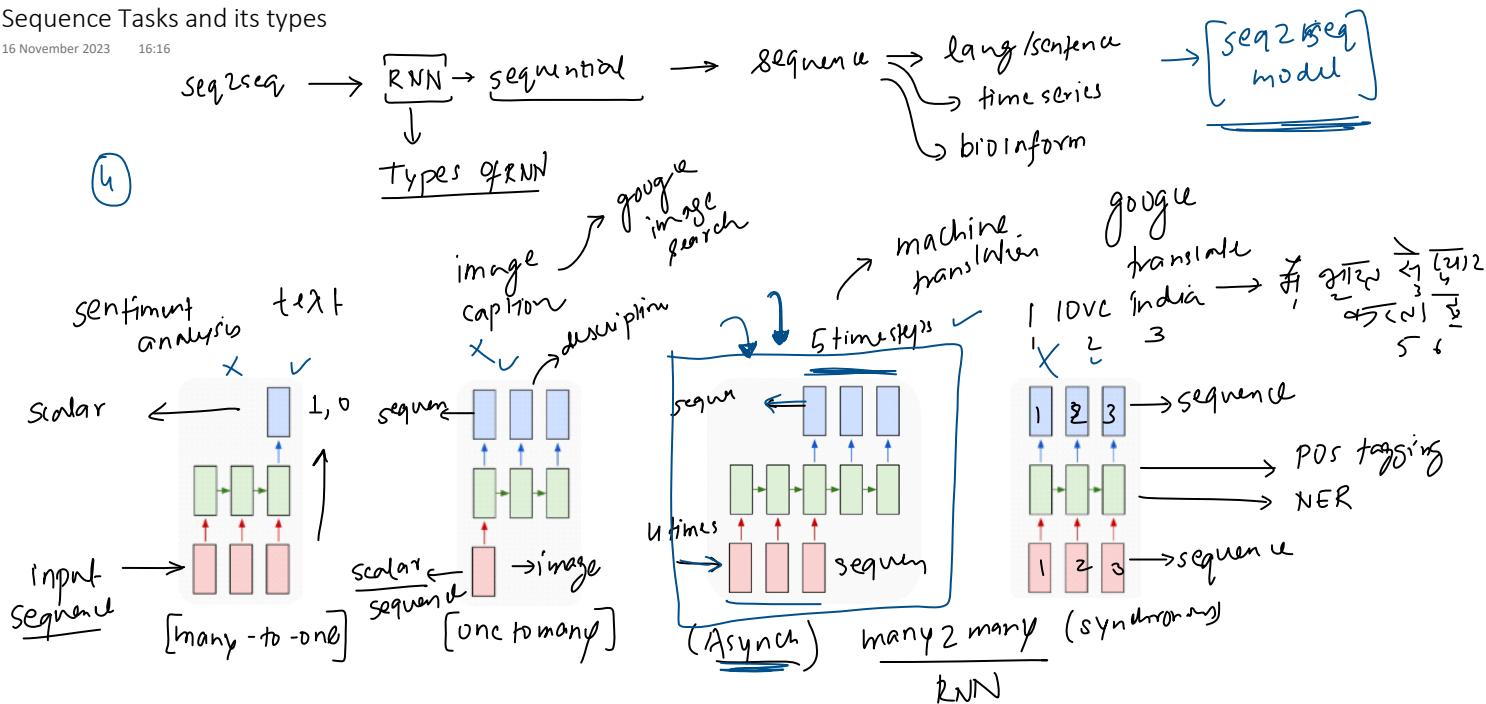


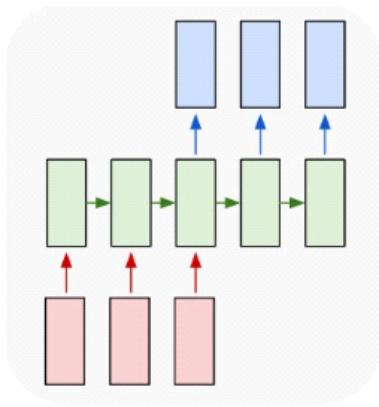
## Sequence Tasks and its types

16 November 2023 16:16



# Seq2Seq tasks

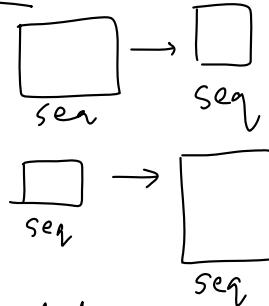
16 November 2023 16:16



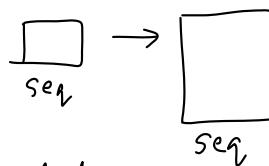
NLP

Seq2seq → machine trans

1) text summariz →



2) Question answer →



knowledge base

3) chatbot → input (text) → output (\*ex)

4) speech-to-text →

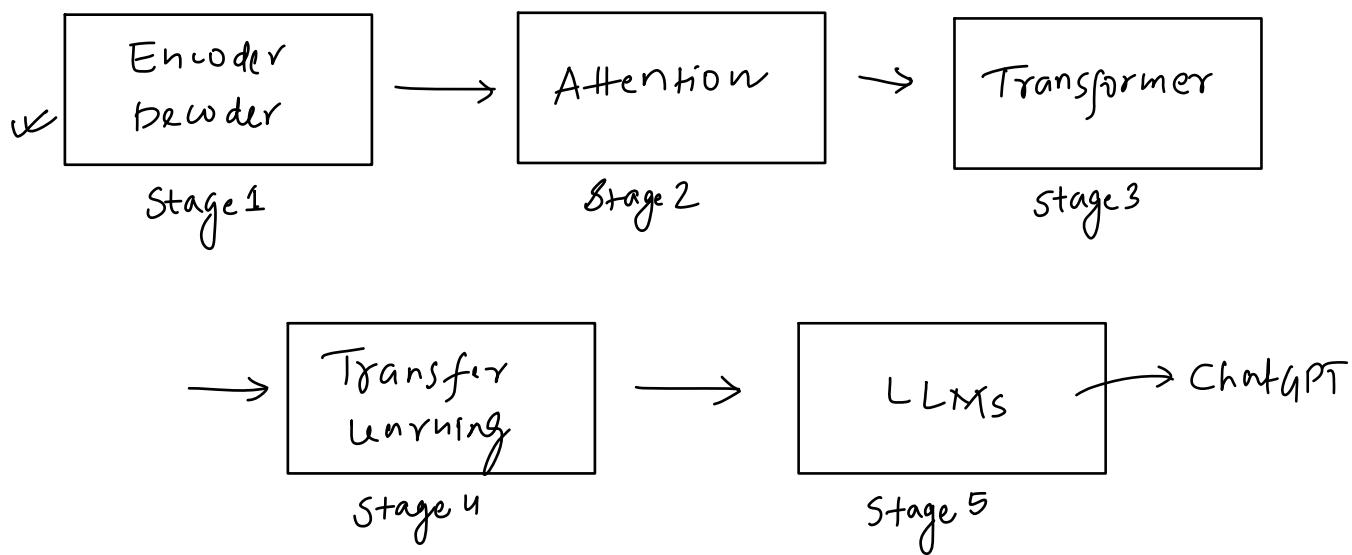
seq seq

→ seq2seq → chatgpt

# History of Seq2Seq Models

16 November 2023 16:16

ChatGPT



2014 Seminal

seq2seq  
 ↓  
 diff  
 ↳ encoder  
 decoder

## Sequence to Sequence Learning with Neural Networks

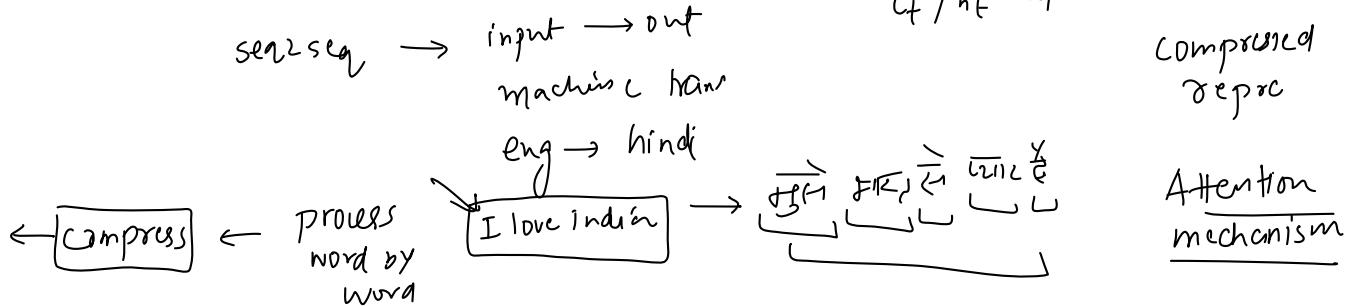
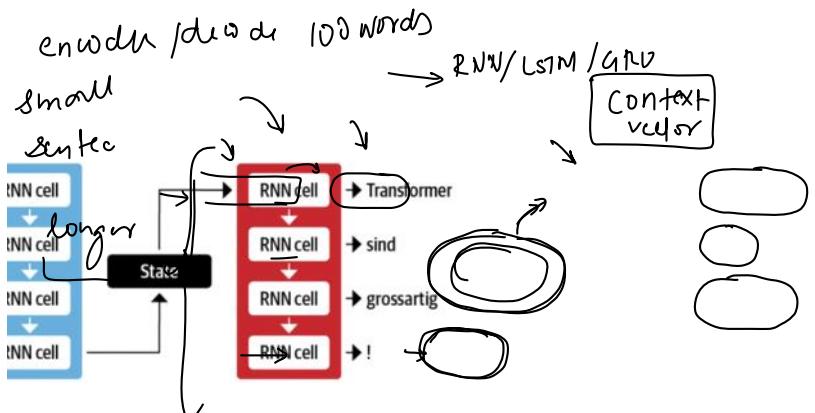
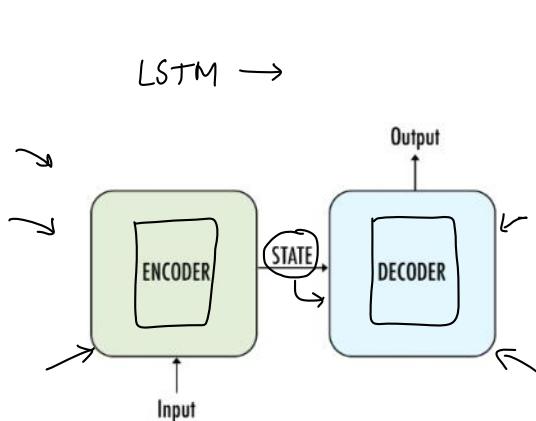
→ Ilya Sutskever  
 Google  
 ilyasu@google.com

[ Oriol Vinyals ]  
 Google  
 vinyals@google.com

[ Quoc V. Le ]  
 Google  
 qvl@google.com

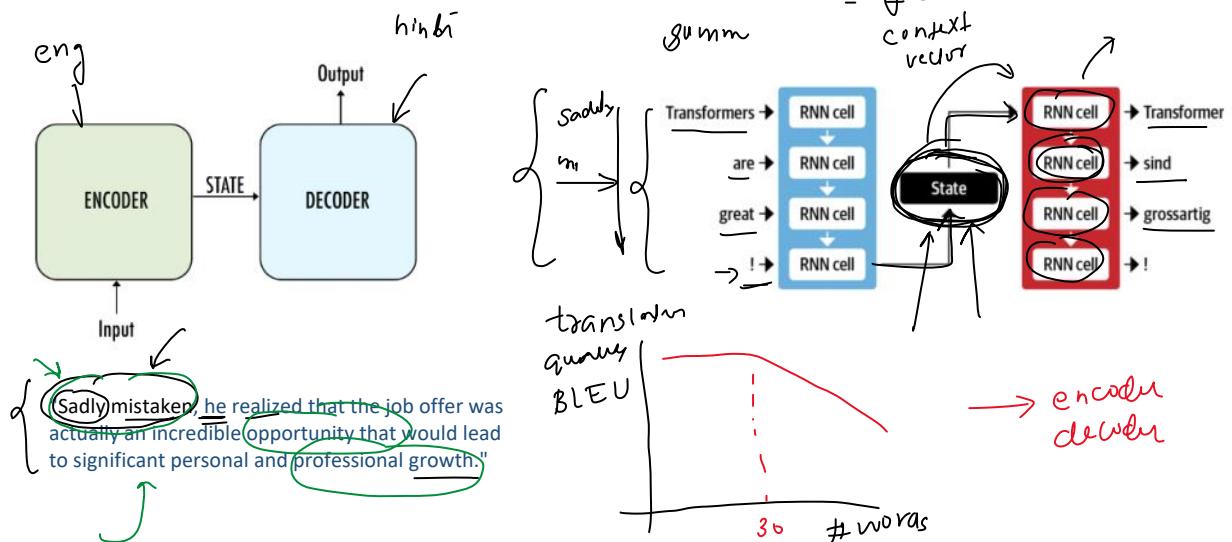
### Abstract

Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then another deep LSTM to decode the target sequence from the vector. Our main result is that on an English to French translation task from the WMT'14 dataset, the translations produced by the LSTM achieve a BLEU score of 34.8 on the entire test set, where the LSTM's BLEU score was penalized on out-of-vocabulary words. Additionally, the LSTM did not have difficulty on long sentences. For comparison, a phrase-based SMT system achieves a BLEU score of 33.3 on the same dataset. When we used the LSTM to rerank the 1000 hypotheses produced by the aforementioned SMT system, its BLEU score increases to 36.5, which is close to the previous best result on this task. The LSTM also learned sensible phrase and sentence representations that are sensitive to word order and are relatively invariant to the active and the passive voice. Finally, we found that reversing the order of the words in all source sentences (but not target sentences) improved the LSTM's performance markedly, because doing so introduced many short term dependencies between the source and the target sentence which made the optimization problem easier.



## Stage 2 - Attention Mechanism

20 November 2023 10:59



[2015] → A Henkim

## NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau  
Jacobs University Bremen, Germany

KyungHyun Cho [Yoshua Bengio]  
Université de Montréal

### ABSTRACT

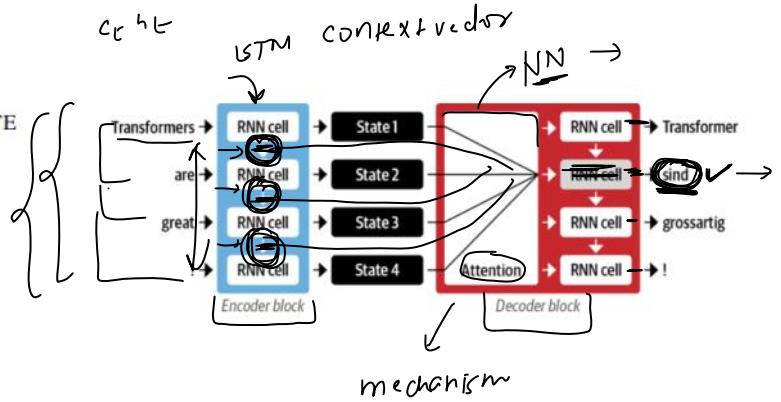
Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently for neural machine translation often belong to a family of encoder-decoders and encode a source sentence into a fixed-length vector from which a decoder generates a translation. In this paper, we conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder-decoder architecture, and propose to extend this by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly. With this new approach, we achieve a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-to-French translation. Furthermore, qualitative analysis reveals that the (soft-)alignments found by the model agree well with our intuition.

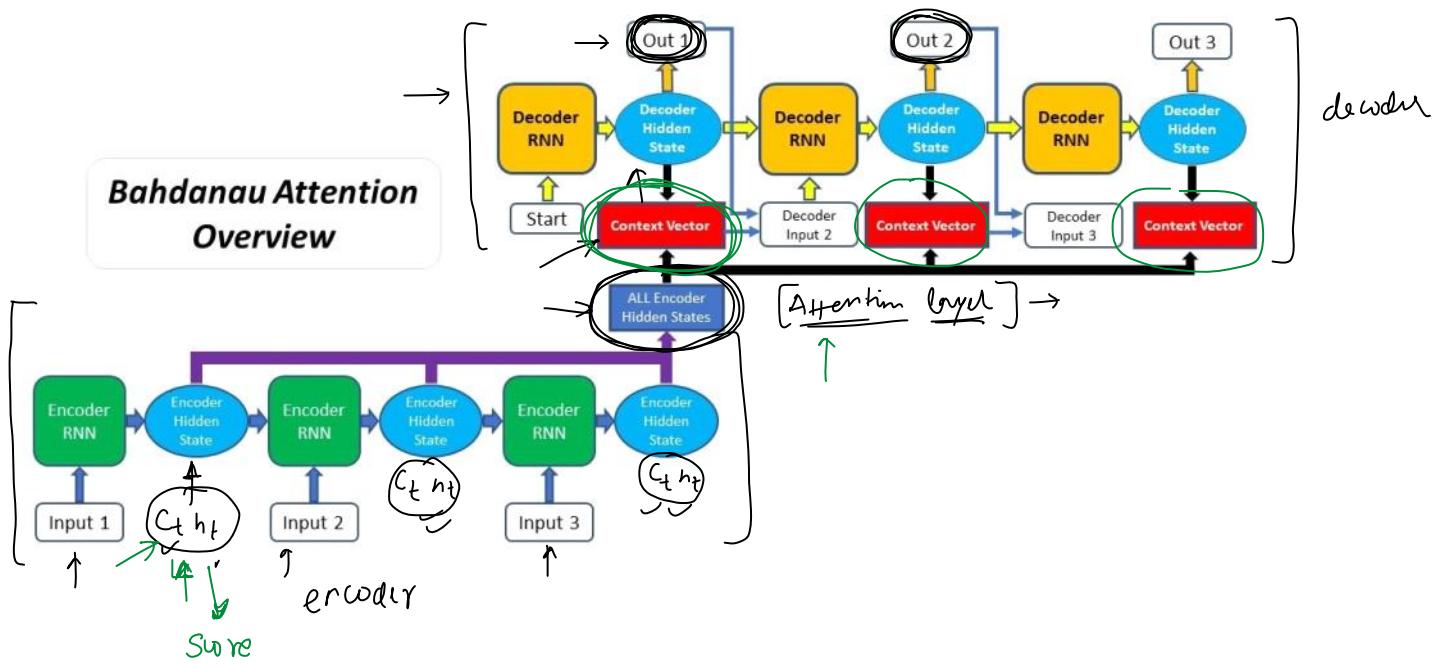
### 1 INTRODUCTION

*Neural machine translation* is a newly emerging approach to machine translation, recently proposed by Kalchbrenner and Blunsom (2013), Sutskever et al. (2014) and Cho et al. (2014b). Unlike the traditional phrase-based translation system (see, e.g., Koehn et al., 2003) which consists of many small sub-components that are tuned separately, neural machine translation attempts to build and train a single, large neural network that reads a sentence and outputs a correct translation.

Most of the proposed neural machine translation models belong to a family of *encoder-decoders* (Sutskever et al., 2014; Cho et al., 2014a), with an encoder and a decoder for each language, or involve a language-specific encoder applied to each sentence whose outputs are then compared (Hermann and Blunsom, 2014). An encoder neural network reads and encodes a source sentence into a fixed-length vector. A decoder then outputs a translation from the encoded vector. The whole encoder-decoder system, which consists of the encoder and the decoder for a language pair, is jointly trained to maximize the probability of a correct translation given a source sentence.

A potential issue with this encoder-decoder approach is that a neural network needs to be able to compress all the necessary information of a source sentence into a fixed-length vector. This may make it difficult for the neural network to cope with long sentences, especially those that are longer than the sentences in the training corpus. Cho et al. (2014b) showed that indeed the performance of a basic encoder-decoder deteriorates rapidly as the length of an input sentence increases.





Stage 3 - Transformers  
20 November 2023 12:18

{ computational complexity }

m words

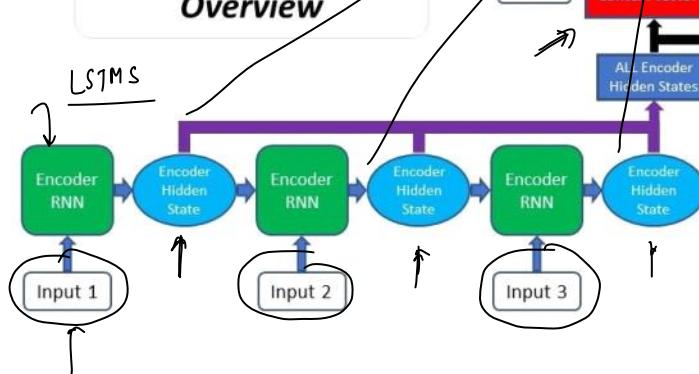
2015-2017

np

m words

after

## Bahdanau Attention Overview



n words

sequential order

en wieder diewel

parallel processing

2017

## Attention Is All You Need

Ashish Vaswani*	Noam Shazeer*	Niki Parmar*	Jakob Uszkoreit*
Google Brain	Google Brain	Google Research	Google Research
avaswani@google.com	noam@google.com	nikip@google.com	usz@google.com

Llion Jones*	Aidan N. Gomez* <sup>†</sup>	Lukasz Kaiser*
Google Research	University of Toronto	Google Brain
llion@google.com	aidan@cs.toronto.edu	lukaszkaiser@google.com

Ilia Polosukhin* <sup>‡</sup>
ilia.ilia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

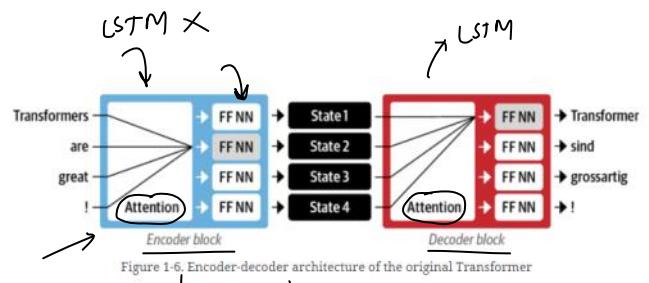


Figure 1-6: Encoder-decoder architecture of the original Transformer

LSTM / RNN cell

Attention

Self-attention

stage

arch

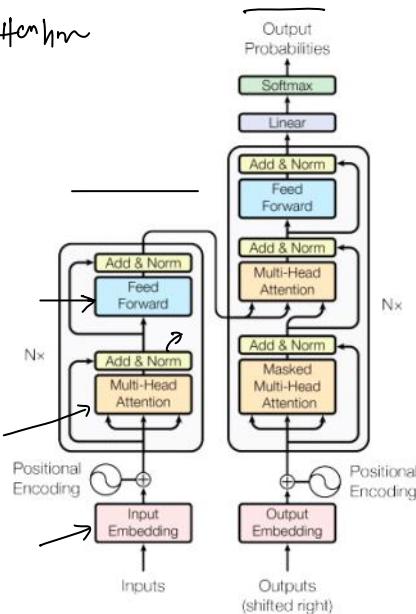
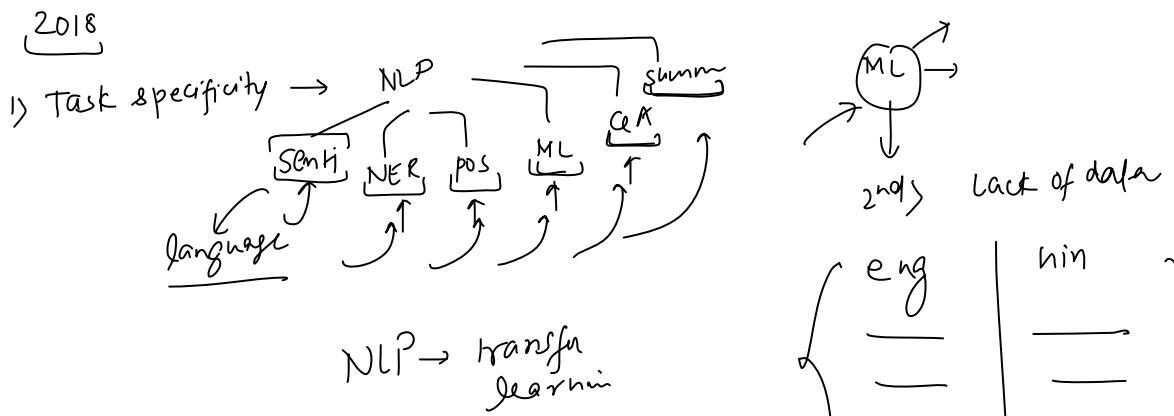
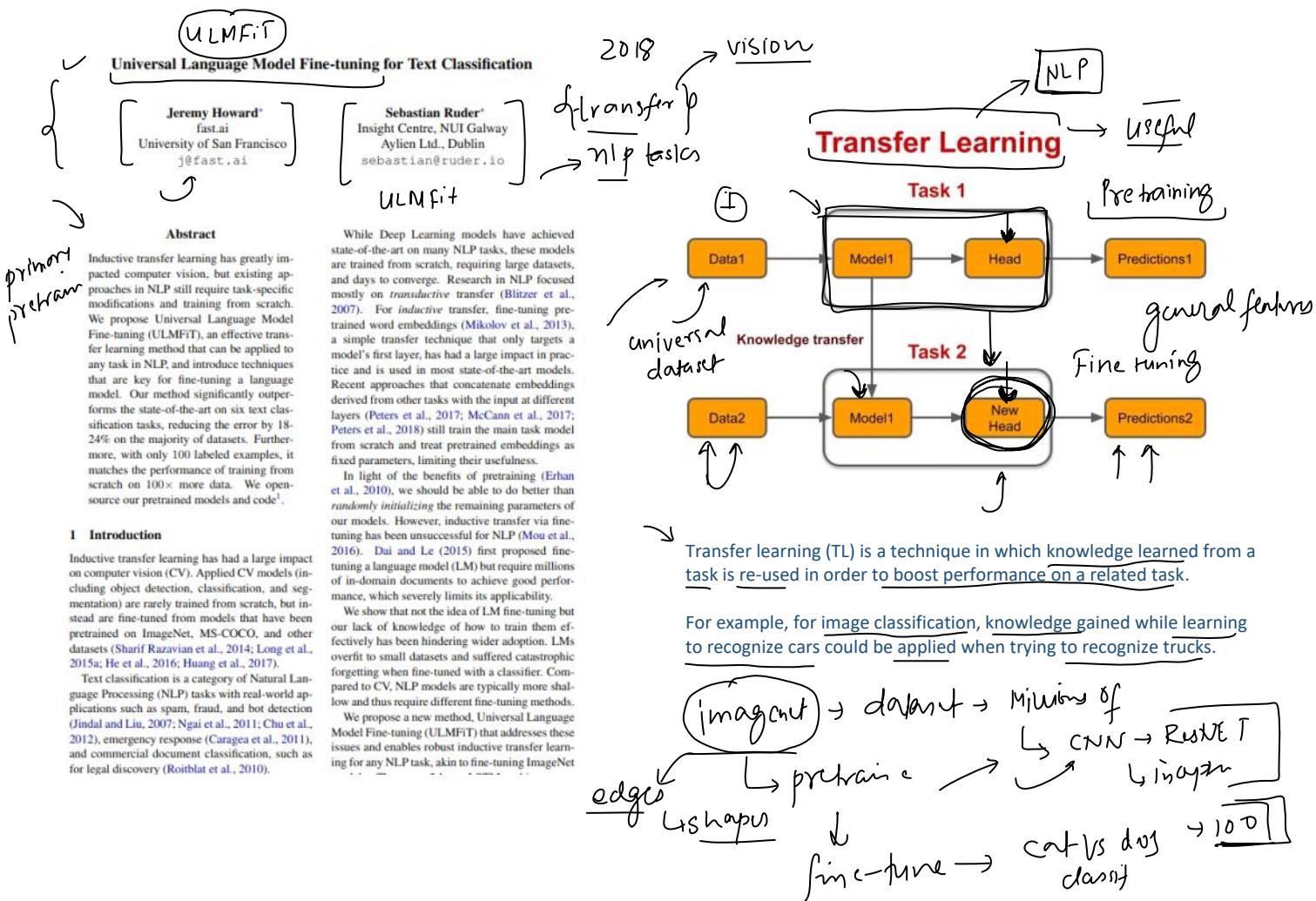
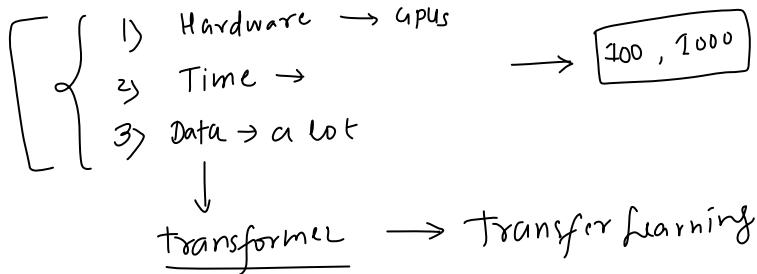
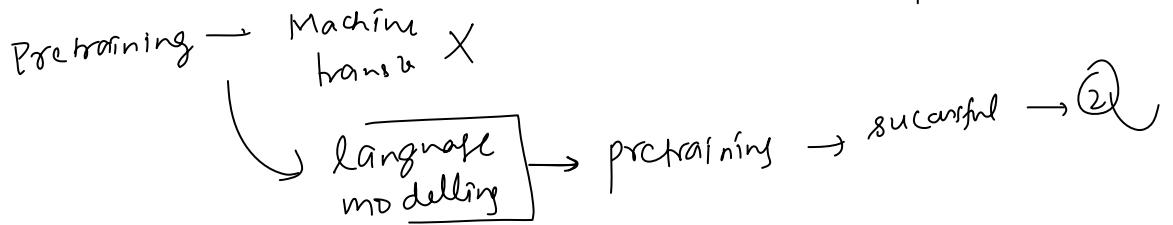
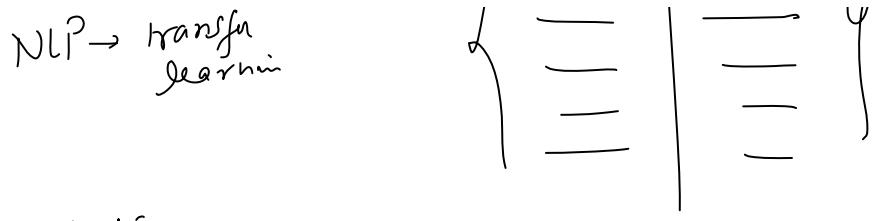


Figure 1: The Transformer - model architecture.

## Stage 4 - Transfer Learning

20 November 2023 15:39





NLP task  $\rightarrow$  NLP/DL model next word pred  
 I live in India. and the capital is  $\frac{\text{is New Delhi}}{\text{J}}$

Language modeling as a Pretraining task  
 1) Rich feature learning  
 The hotel was exceptionally clean, yet the service was  $\frac{\text{bad}}{\text{pathetic}}$

$\rightarrow$  know trans  
 ↓  
 text classif / ques. | textsum | NLP / PGM  
 pdf  $\rightarrow$  dataset labeling

mt (new  $\rightarrow$  supervised labeled)  
 eng | hin  
 $\rightarrow$  unsupervised task

2) huge avail of data

[ULMFiT]

X transformer

AND LSTM  $\rightarrow$  wikipedia

Unsupervised  
pretrain  
language  
modeling

finetuning

classifier  $\rightarrow$  imdb

$\rightarrow$  yelp

$\rightarrow$  new domain

$\rightarrow$  model

test

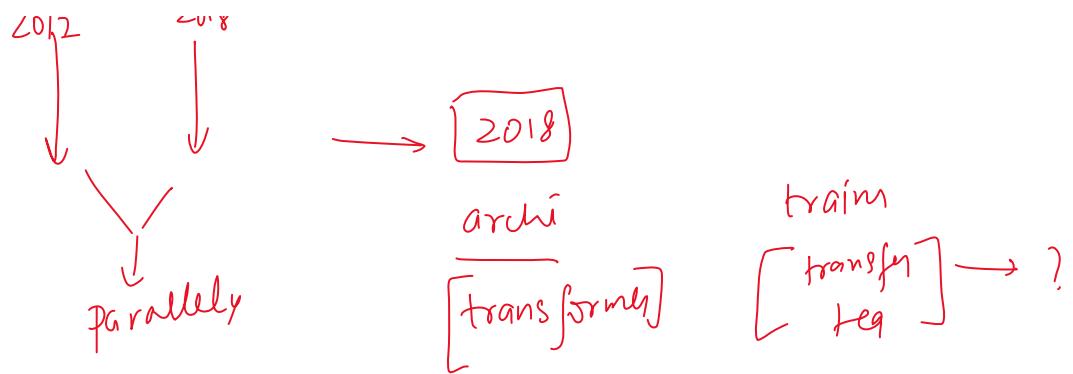
scratch  $\rightarrow$  1000 rows

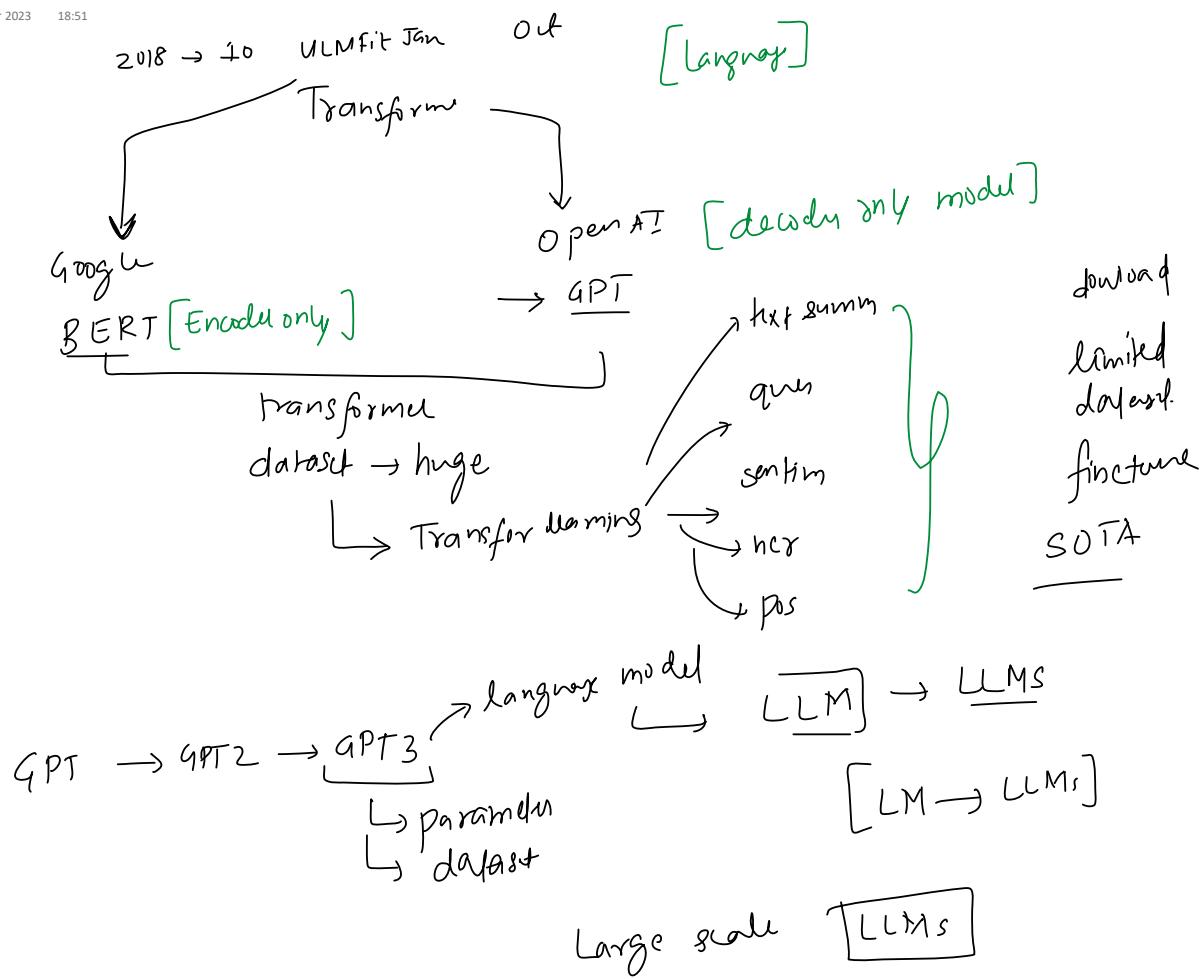
100 rows  $\rightarrow$  better  $\rightarrow$

State of the art

2012

2018



Qualities of LLMs

1) Data → billions → GPT3 → 45 TBs  
 ↗ book, websites, internet  
 ↗ diversity → bias

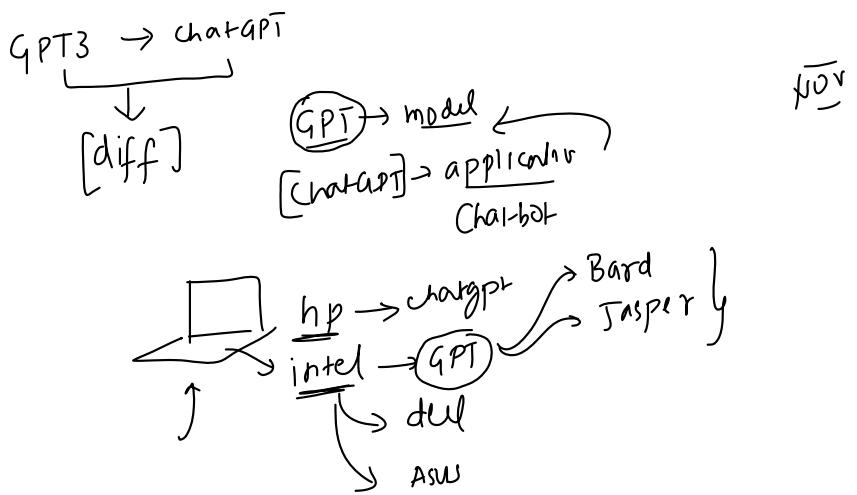
2) Hardware → Cluster of GPU → GPT3 → Supercomputer → 100s NVIDIA GPUs  
 ↗ individual companies, govt, institutes

3) Training → days to wccs

4) Cost → hardware + elec + infra + experiments → millions  
 ↗ individual companies, govt, institutes

4) energy consumption  
 ↗ GPT3 → ...

↳ energy consup  
↳ q p 13  
↳ small town  
↳ month



GPT3 → [ChatGPT]

1) RLHF → Reinforcement learning from human feedback

- + 1 Supervised finetuning → dataset
- + 2 reinforce → prompt production
- + responses
- + human → response bank

===== y labeled

2) Incorporate safety and ethical guideline

- + minimize bias

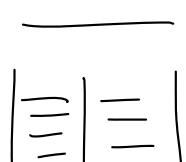
3) improvement in contextual point

===== context retain → maintain context } dialogue convers

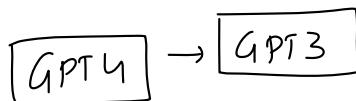
4) Dialogue specific training

- + conversation
- + better understanding → dialogue lang → partitions

5) ChatGPT continuous imp → human feedback  
↳ usu



train → refining



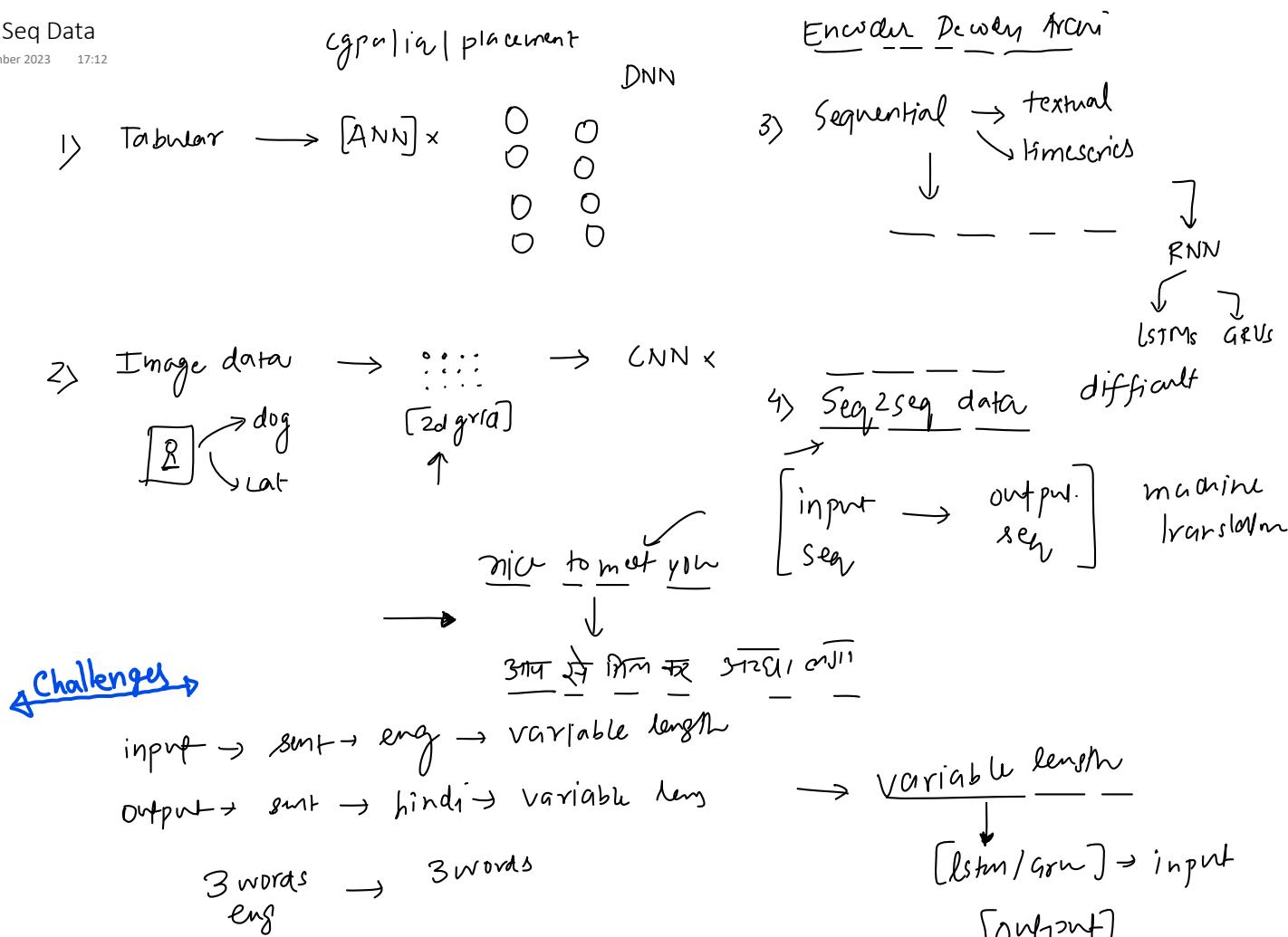
$\left| \begin{array}{c} \diagup \\ \diagdown \end{array} \right| = \left| \begin{array}{c} \diagdown \\ \diagup \end{array} \right|$

train  $\rightarrow$   $y^{true}$

$\hookrightarrow$   $\boxed{\text{GPT4}}$   $\rightarrow$   $\boxed{\dots}$

# Lecture #68 (Encoder-Decoder)

Seq2Seq Data  
08 December 2023 17:12



- input એ sentence જાહેર કરી શકતું હોય (variable length)
- output " " " " hindi (" " ")
- input એ 3 ટે word -નું output એ જેવી નાથ બાળ પણ હોય

Variable length વાસ્તવી કિન્તુ handle કરુણે ચિન્હિત કરી શકતું હોય।

input એ, GRU/LSTM કર્યાન્નું હાણ્ણું। Output એ કરીજાનું handle કર્યાન્ણું હોય।

Sequence-sequence prob solve કર્યાન્ણું છેન્સ આપ્યાન્ણું Encoder & Decoder Architecture

→ TFLearn 1

# Before Starting

08 December 2023 19:24

## Prerequisite

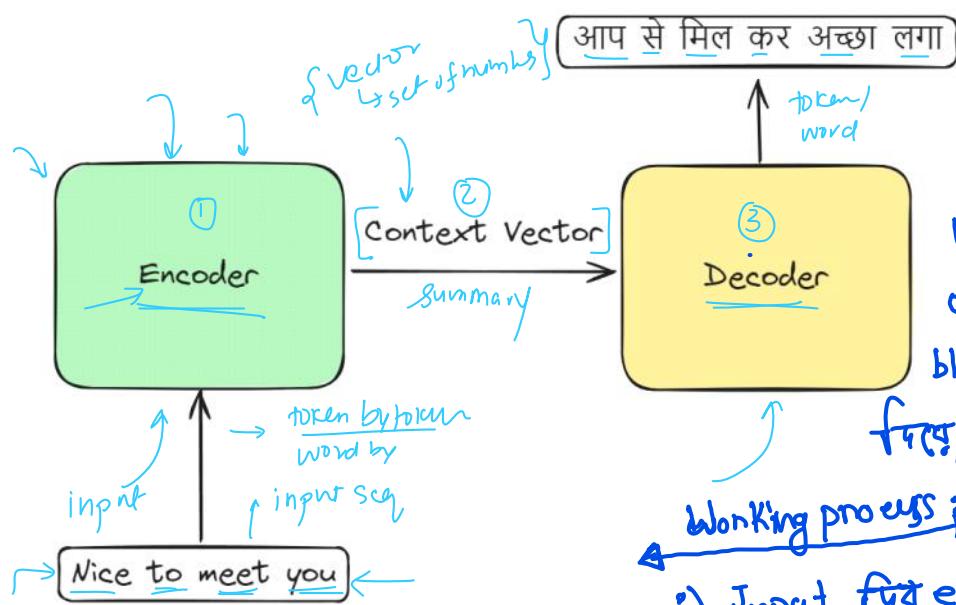
- RNN / UTM

## Plan of attack

- simple version
- deep
- improvements

## Translation as an Example

Machine tran  $\rightarrow$  exam



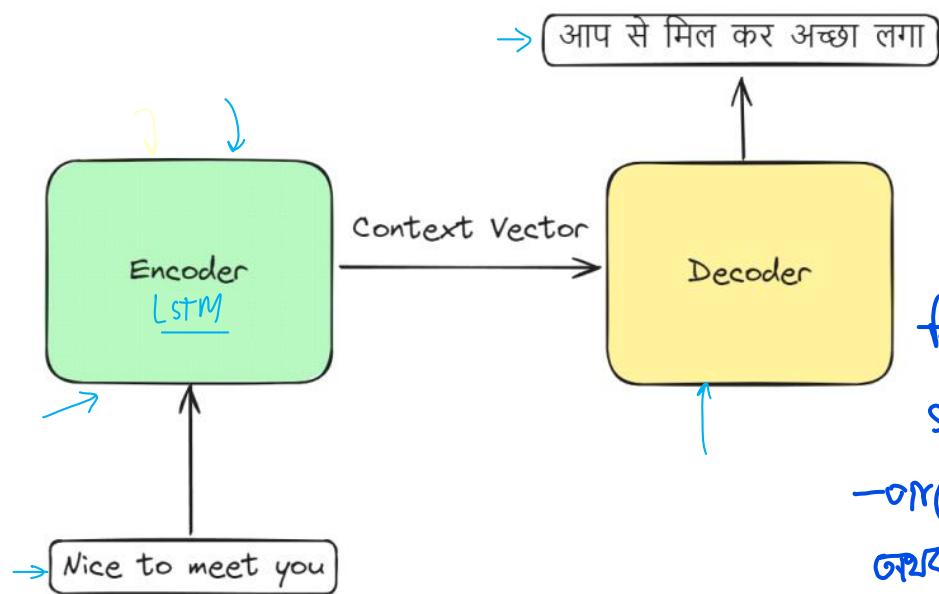
Encoder-Decoder (एडिटर)

block याकू एटे Encoder  
याकू एटे Decoders। प्रति रोटे  
block एटे context vector

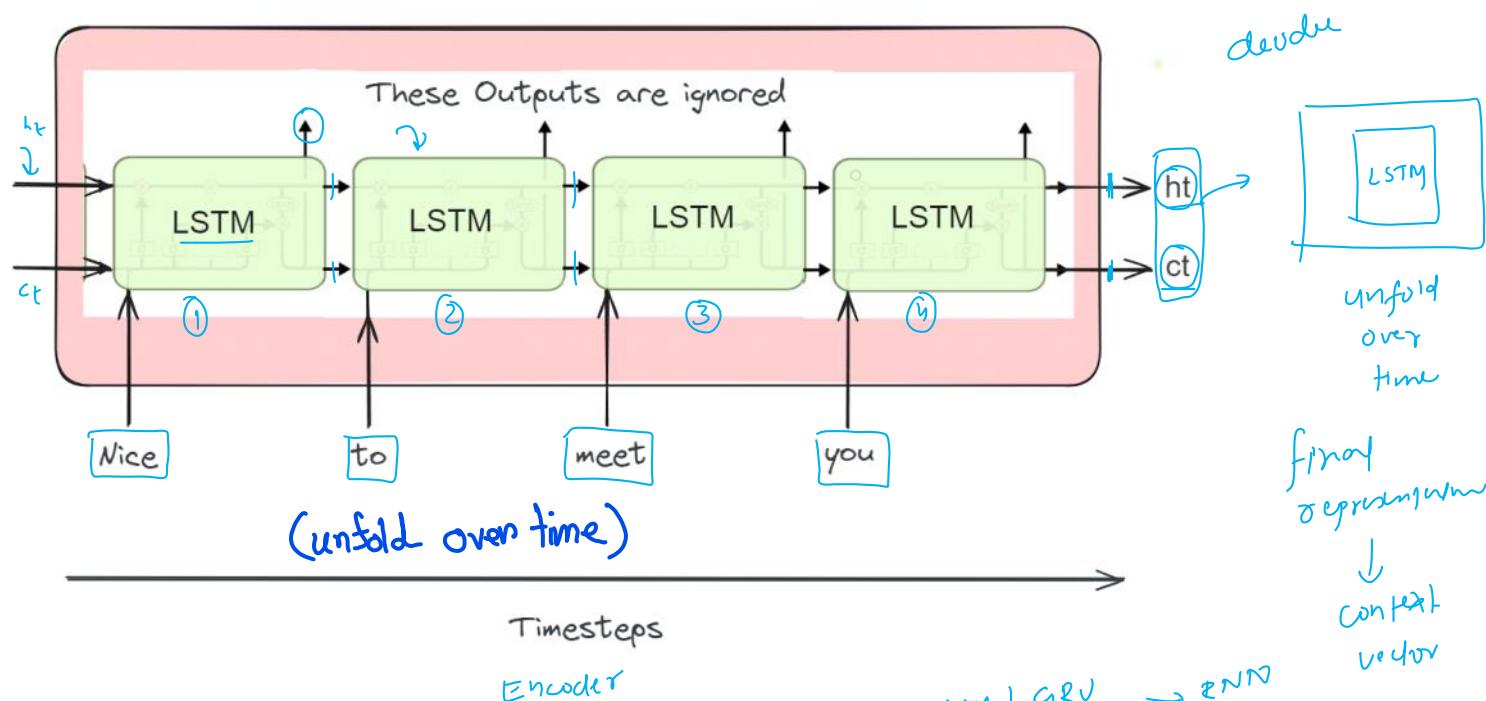
फिल्म याकू थाकू।

Working process :-

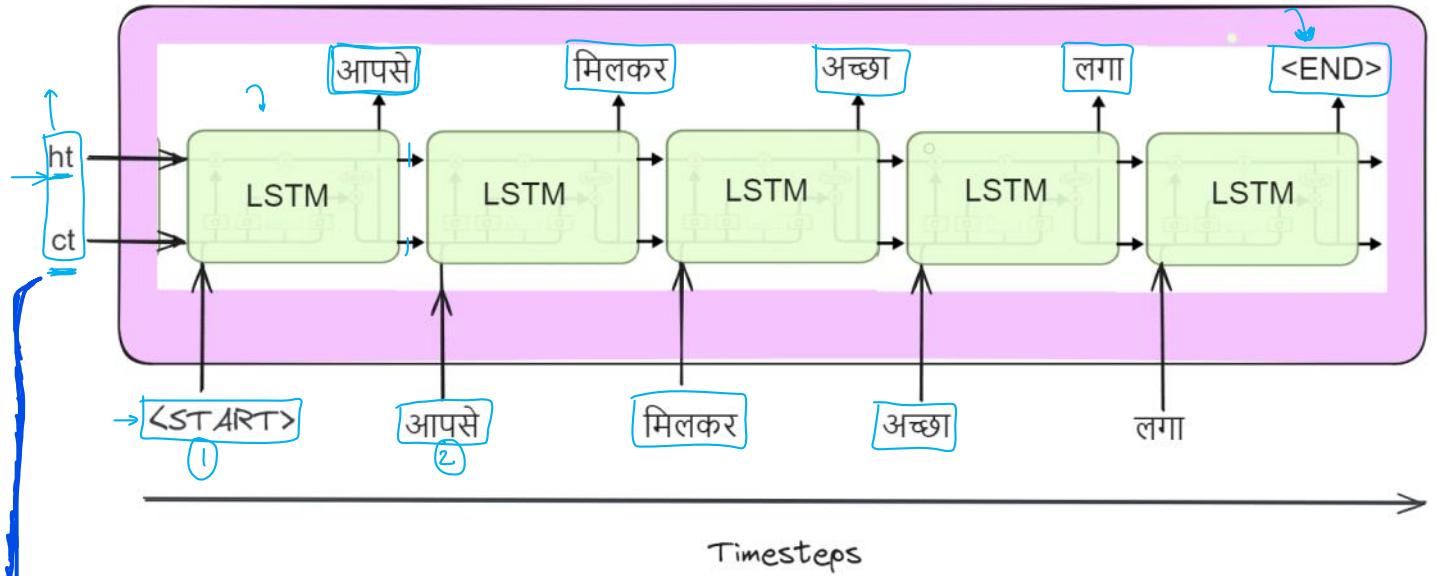
- i) Input for encoder (a token by token or word by word).
- ii) Encoder input vector रोक्तु रोक्तु summarise करकूवे। यसु summarization output रोक्तु एटे context vector एक्सु आमदु context vector यामुठि।
- iii) Decoder एटे context vector-को रुक्तु word by word or token by token उन्हु language (ए त्रान्शित फिर्ते।



अर्थात्, जटाएँ एडेकोडर  
& डेकोडर द्वा गर्फु टेम्प  
किंवा यहाँ पढ़ना है ऐटो  
सीक्यूएल के प्रोसेस करूँ  
—एक्ट्रो पाठ् एजरन्स् एल्टी  
जटाए ग्रु। Research papers  
एल्टी-युवश्वर् कर्या रखेह।

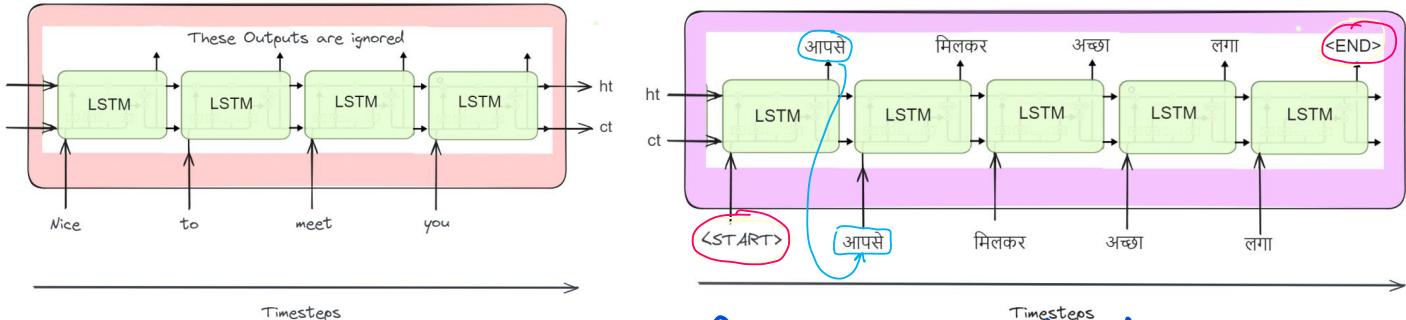


⊗ —जटाए एल्टी और ग्रु युवश्वर् करति। किंवा, RNN ना याहुर् RNN के vanishing  
Gradient एवल्म देख देह। ⊗



- Encoder का output एवं Decoder का input रिसेप्टर हैं।
- प्रारंभिक आवश्यक स्पेशल सिम्बल <start> एवं encoder & decoder का starting & ending गुणात्मक है।
- timestep-1 का output जो है timestep-1 का पहला input रिसेप्टर है। अब, उसका अनुदान देना पड़ता है 2nd special symbol <end>।
- पार्ट 2।

encoder & decoder ଏକାଇ ମାତ୍ର train ହୁଏ । ତେଣୁ training  
ଅନ୍ୟ ସୂଚାରୁ ଧ୍ୟାନ ଦିଆଯାଇଲେ diagram ଏକମାତ୍ର ହୁଅଥାବା ।



Machine translation (मशीन अनुवाद) dataset हिस्तेन व्यवहार तरीके training process तुल्यता।

Dataset for training  $\Rightarrow$

Eng → Bangla  
Eng → Hindi

i) Think about it → পেরি মানুষের ব্যবহাৰ  
ii) come in → আসিব আসিব

→ DOHE (total unique words)

→ numbers  
↳ NLP

Tokenize the data set:

row1 → [Think, about, it] → [পেরি, মানুষের, ব্যবহাৰ]

row2 → [come, in] → [আসিব, আসিব]

eng      one box

bangla

**think** [10000]

**art** [01000]       

→ [0,1,0,0,0]

<start> <end>

[10000000]

      |  
      |

**Kend** [010000]

[0,1,0,0,0]

प्राप्त probability  
दूसरी नियोजित रूप  
output i

Encoder पे जानेगी अबत्रा row input first timestep  
पर्वत्रा feauture word by word encoder पे ध्याप (think about it) 

Decoder-এ অপরে special symbol start থাকে। Decoder  
এর output-একটি softmax layer থাকে যেখানে  
input (গুরুত্ব) for node থাকবে।

Diagram illustrating a sequence-to-sequence model with an attention mechanism:

- Input Sequence:** Think
- Hidden States ( $h_t$ ):** Three vertical rectangles representing hidden states.
- Initial Vector:**  $0$
- Context Vector:**  $[random]$  (highlighted by a green bracket)
- Intermediate States:**
  - First hidden state receives  $0$  and produces  $(random)$ .
  - Second hidden state receives  $(random)$  and produces  $[random]$ .
  - Third hidden state receives  $[random]$  and produces  $[random]$ .
- Attention Mechanism:** Arrows point from the first hidden state to the second, and from the second to the third, indicating information flow or context passing.

The diagram illustrates a sequence-to-sequence model architecture with an attention mechanism, showing the forward pass and the backward pass (gradient flow) for calculating the loss.

**Forward Pass:**

- Input:** A sequence of words represented by vectors:  $[0, 0, 1, 0, 0, 0, 0]$ .
- Encoder:** Processes the input through two hidden states. The second hidden state is  $[0, 1, 0, 0, 0, 0, 0]$ .
- Decoder:** Generates a sequence of words represented by vectors. The first generated vector is  $[0, 0, 0, 0, 0, 0, 0]$ .
- Attention:** Compares the encoder's hidden state with the decoder's hidden state to produce weights. The weights for the first word are  $[0.1, 0.2, 0.1, 0.15, 0.3, 1.5, 0.05]$ . The highest weight (1.5) corresponds to the word "end".
- Output:** The final output vector is  $[0, 0, 0, 0, 0, 0, 1]$ , where 1 indicates the word "end".

**Backward Pass (Gradient Flow):**

- Loss:** The loss is calculated as the negative log-likelihood of the generated sequence. The loss value is  $-0.070$ .
- Gradients:** Gradients flow from the loss back through the network. The gradients for the last hidden state are  $[0.1, 0.2, 0.1, 0.15, 0.3, 0.1, 0.4]$ .
- Encoder Gradients:** The gradients for the encoder's hidden state are  $[0, 1, 0, 0, 0, 0, 0]$ .
- Decoder Gradients:** The gradients for the decoder's hidden state are  $[0, 0, 1, 0, 0, 0, 0]$ .

Encoder (encoder) [teacher forcing] ↓ (decoder) [gradient update] 1st step द्वारा output अपना रखा हिं [0 1 0 0 0 0 0] यहि अनुरूप output प्राप्त करने के लिए अपना value तक बढ़ावा देना फिर -स्थानानुरूप convergence जापाणी है !

Diagram illustrating the forward pass of a neural network layer:

- Inputs:**  $y_{true}$  and  $y_{pred}$
- Transformation:**  $y_{true} \rightarrow y_{pred}$  (using weight matrix  $W$  and bias  $b$ )
- Outputs:**  $z$  and  $a$

Annotations in the diagram:

- Yellow arrows point from  $y_{true}$  to  $y_{pred}$  and from  $y_{pred}$  to  $z$ .
- Red annotations highlight the weight matrix  $W$  and bias  $b$ .
- Green annotations show the intermediate values  $z$  and  $a$ , along with their activation function  $\sigma(z)$ .
- A circled "0.4" is shown near the output  $a$ .
- A circled "0.4" is also shown near the output  $a$ .

$$\rightarrow \underline{y_{\text{true}}} [0.1, 0.0, 0.0, 0.0] \text{ সম্প} \rightarrow \underline{L} [0, 0, 1, 0, 0, \dots] \rightarrow \underline{\hat{y}_{\text{pred}}} [0.2, 0.1, 0.3, 0.2, 0.1, 0.1] \text{ এব} \rightarrow \underline{L} [0.1, 0.2, 0.1, 0.15, 0.3, 0.5] \rightarrow \underline{L} [0.1, 0.2, 0.2, 0.1, 0.3, 0.1, 0.4] \rightarrow \underline{avg} \rightarrow 0.3$$

loss:  $-1 \times \log(0.1) = -1$  (1)

$= -1$

$= -0.3$  (0.3)

১০-আমাদের forward propagation করে। এখন, আমরা loss calculation করুন। সুবিধা, আমাদের  
প্রেরণামূলক বৈজ্ঞানিক পদ্ধতি হচ্ছে multi-class classification তাই, loss function হিসেবে categorical-cross-entropy ব্যবহার  
করুন।

$$\text{Categorical-Cross-Entropy} = -\sum_1^T y^{\text{true}} \log(y^{\text{pred}})$$

$$\text{Total Loss} = -1 - 1 - 0.3 = -2.39 \quad \left[ \begin{array}{l} L_{t=1} = -1 \times \log(0.2) \\ L_{t=2} = -1 \times \log(0.1) \\ L_{t=3} = -1 \times \log(0.4) \end{array} \right] \quad \left[ \begin{array}{l} \text{নথ্যতাম্বীয় হয়, যাকুন prediction মানিকা রয়েছে তাহুন loss বাস্তু} \\ \text{training} \end{array} \right]$$

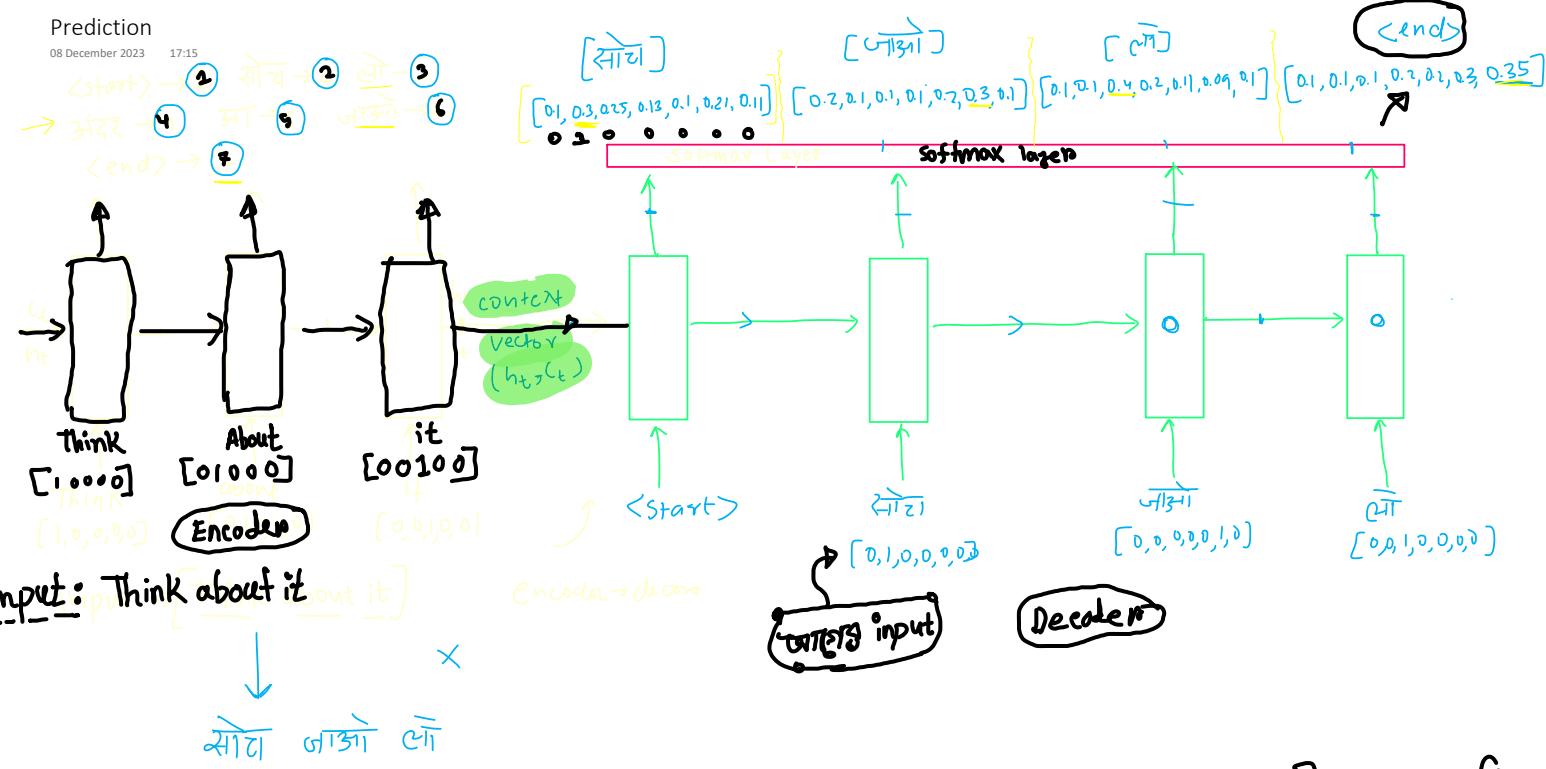
আমাদের, forward propagation করে, loss calculation করে। এখন, back propagation করুন।  
আমরা জানি back propagation কী step এ হচ্ছে। ① gradient calculation ② update parameters.

### 1: Gradient calculation:

We calculate gradient of the loss with each of the trainable parameters.

এখন, LSTM, Dense layers, softmax এ যোগ্যতা পৰামিত্ৰ আছে গুৱাখনা দিলে gradient  
calculate কৰিব। Gradient calculate-কৰে loss function এ তাৰ পৰামিত্ৰ কোটকো contribute  
কৰব। যাকুন পৰামিত্ৰ কোটো দিনৰ কৰিব এবে কোটো loss function এ গুৱাখনা কৰুন।  
প্ৰাণ্যত কোটো parameter update কৰি optimizer দিলে।

# Training Complete এবং Prediction করুণা

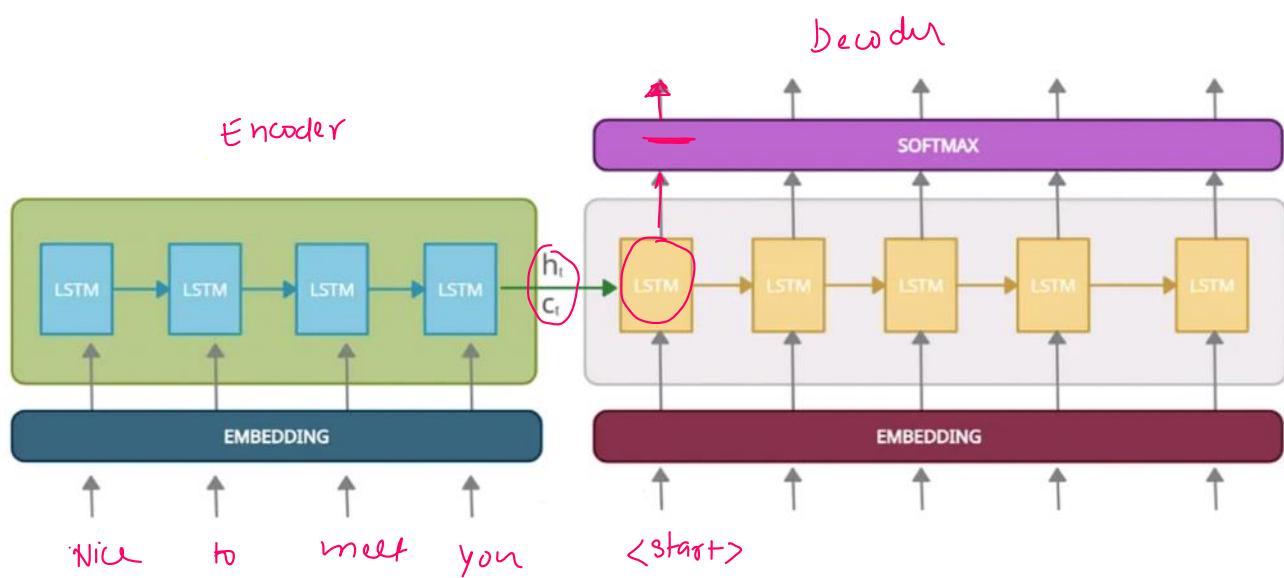
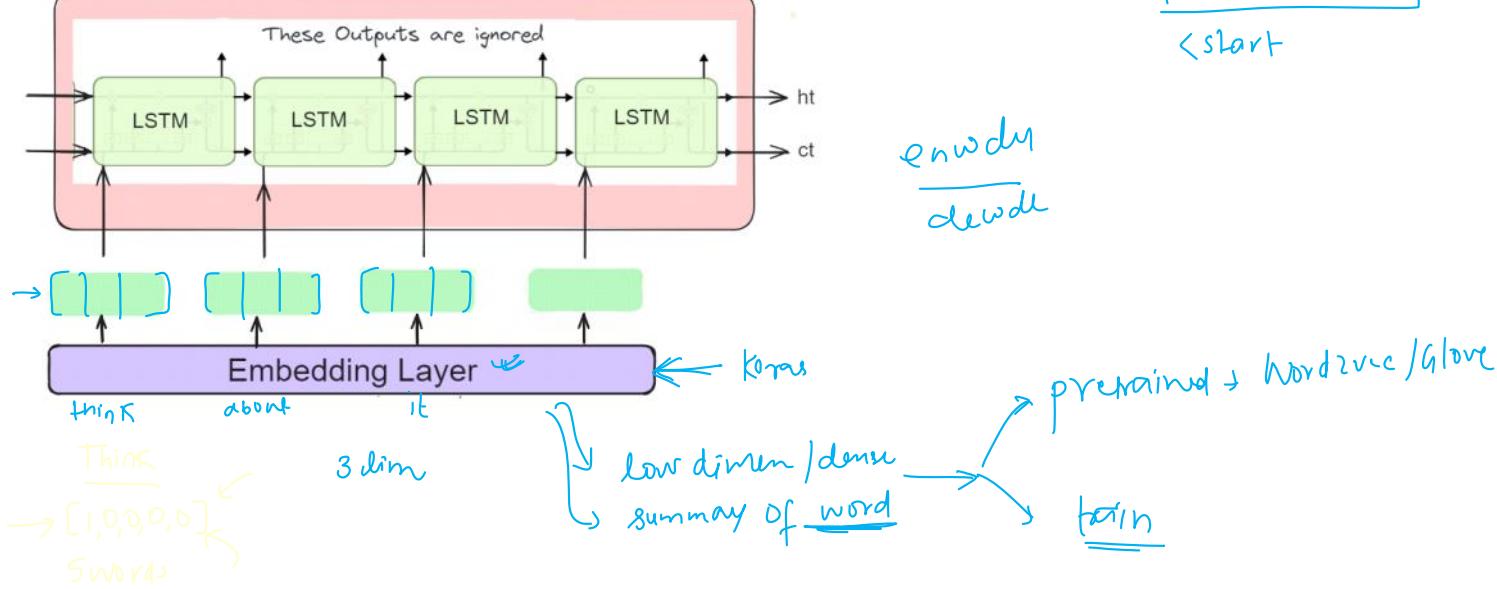


prediction-এর জন্মে Decoder-এ আমাদের কাছে ক্ষেত্র প্রদর্শিত হয়। এখন প্রাপ্তি এবং প্রক্রিয়া  
 নেক্সট  
 ক্ষেত্র ক্ষেত্র

# Improvement of the encoder & decoder

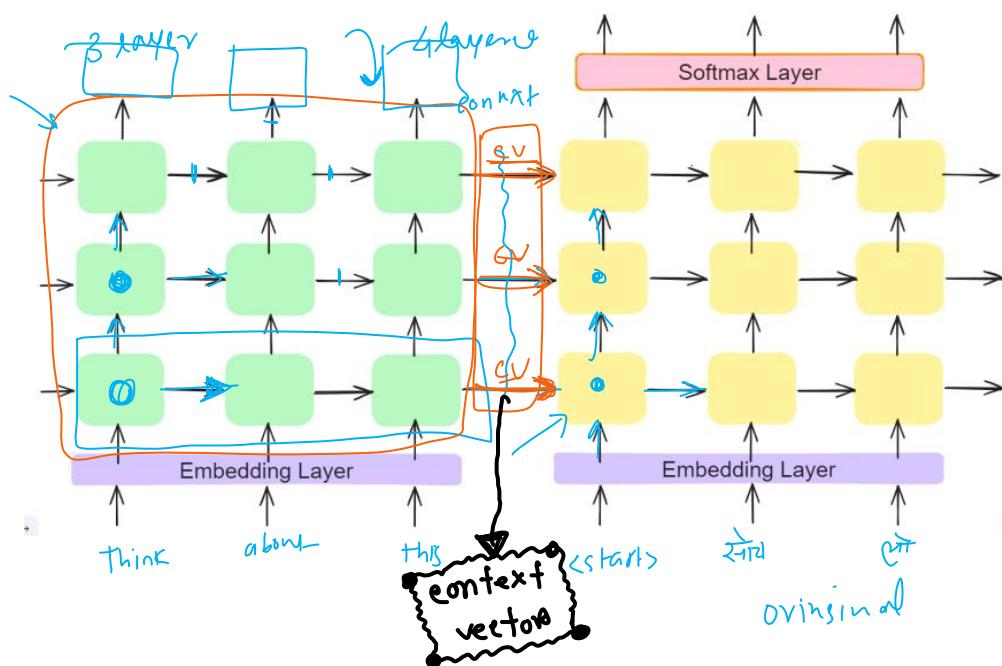
## Improvement 1 - Embeddings

08 December 2023 17:16



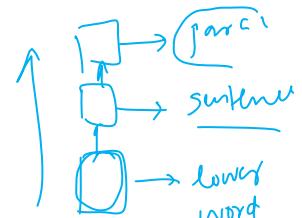
→ আমরা improvement কিভাবে embedding -এ কৈবল্য করবেন। Embedding technique  
আমরা আপনা দুঃখী। এখানে কুকুল আমরা embedding layer (Integer encoded data ফিরে)  
একটা small dimension এর output পাই। সুলভে মেরো আমরা 3 dimensional ফিরে তুলবলুম।  
আমরা আমাকু একটা word এর summary পাই। অৱশ্য, embedding low dimension & dense হব।  
এখানে আমরা pre-trained embedding (Word2vec/Glove) -এ গুরুত্ব দেওয়া পাই। এইটা নি (গুরুত্ব  
দেওয়া) আমরা training dataset -এ কৈবল্য পাই। Now, see the final diagram of embedding  
both for in encoder and decoder!

long para single layer



3) deep learnin NN learning

- 1) long term dependency
- 2) layered represent

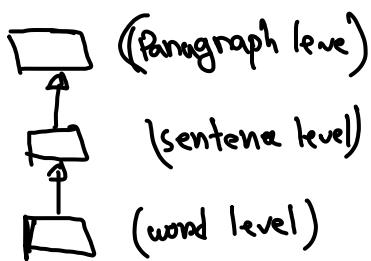


the phone battery is bad  
but the memory is great

Multi-layer LSTM এ্যথাৰ কষ্টলৈ encoder-decoder Architecture-এ improvement আসে। Input দিবলৈ embedding layer এতু মধ্য দিলৈ LSTM-এ শাম্ভাৰ তাৰপাত্ৰ, যেই LSTM অনলাইনে LSTM এতু দিকে update কৰা ও উপরোক্ত Layer-এতু LSTM কুলোকে হিপুল প্ৰোড়ি কৰা।

Multi-layer LSTM এ্যথাৰ কষ্ট স্থিতি:

- i) Long term dependency গুণো handle কৰত আৰি। Long term dependency বলতু কোন input ঘৰুৱা word সংজ্ঞা-সন্তোষৰেখি আ-বৰুৱা paragraph।
- ii) Layer representation সুবিধা আৰি।

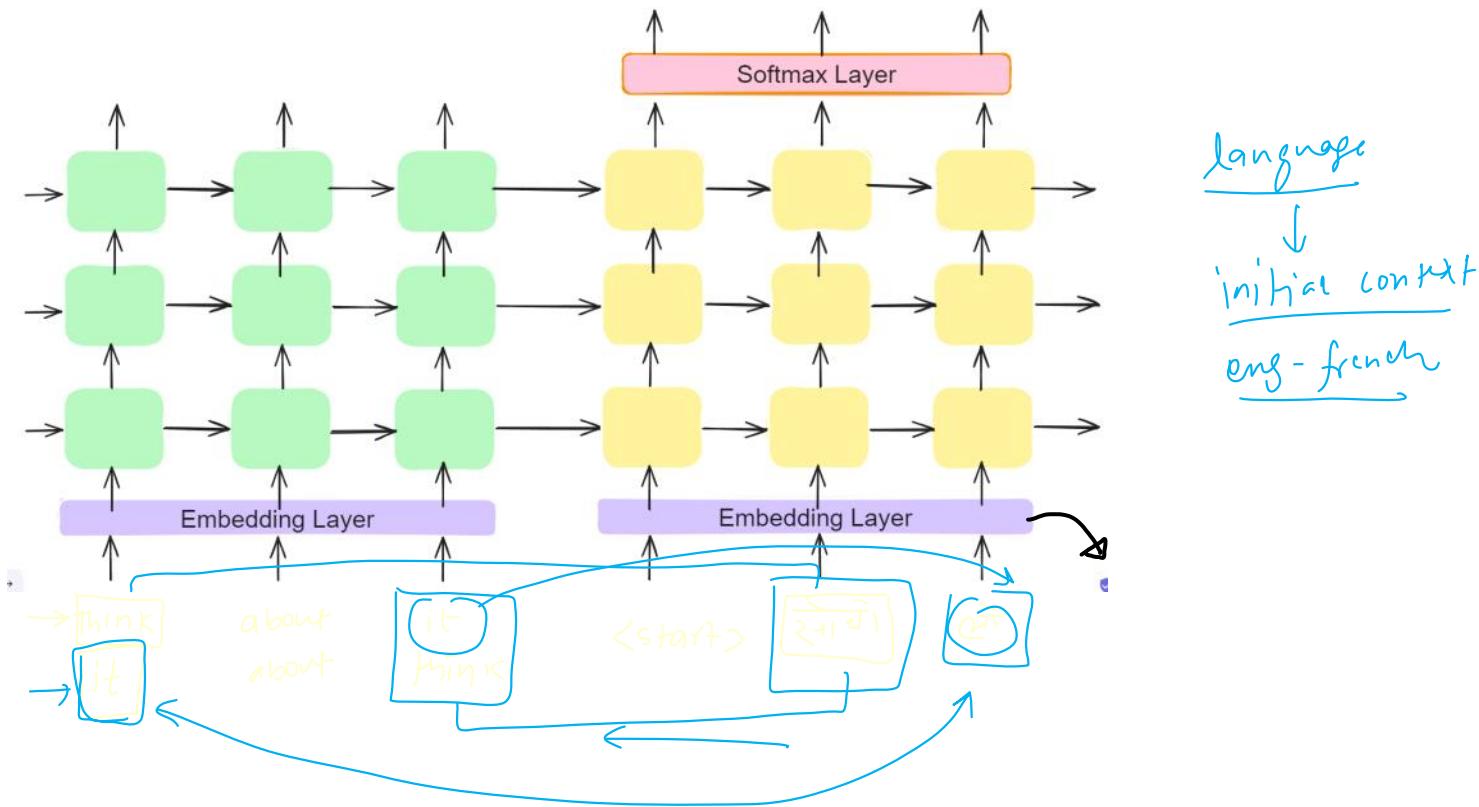


NN কুলোজ layer দেখি হৈলৈ semantic meaning  
কোনোৰাৰে capture কৰত আৰুৱে।

- iii) এখন, আমৰা Multi-layer LSTM এ্যথাৰ কৰচি এতু কোনো আমদানি parameter কোনো  
অনেক সুট্টি পাইছু। parameters এতু সহজে সুট্টি লাগিয়া আন্তৰ complex pattern কোনো  
কুলোজ পাইবে। আবগিন আবগিন overfitting কৰিব না।

## Improvement 3 - Reversing the Input

08 December 2023 17:16

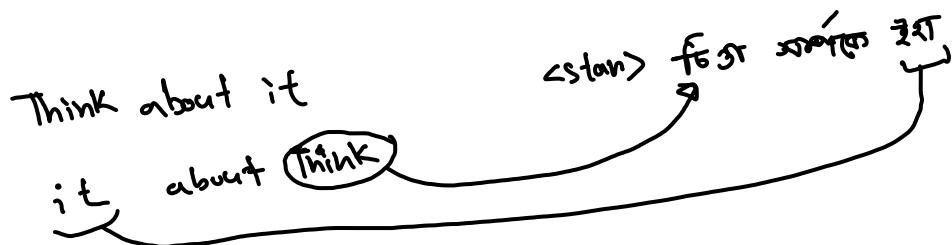


Input କେ reverse କରିବ ।

Original input: Think about it.

reverse: It about think.

Decoder କେବଳ first reverse କରିବ



Think & first କୁଣ୍ଡି କାହାକାଟି ତାକୁ ପଞ୍ଚମ କରିବ । As a

a result, gradient କେ propagate କରସି ଉପରେ କମ କାଗଜ । କିନ୍ତୁ, it କୁଣ୍ଡି ପଦ୍ଧତି - କୁଣ୍ଡି ଏବଂ - କୁଣ୍ଡି । ଯାନ୍ କୁଣ୍ଡି ରସ ଲୋ ଅଛି Language ଏବଂ ଏଣ୍ ବାକୀ ବାକୀ । ସମ୍ଭାବନା କେବଳ language ଏହା Eng. sentence ରୁ meaning ଓ ଏହା language ଏ ଏକାକି word କୁଣ୍ଡି କୁଣ୍ଡି କରିବାକୁ ବାକୀ ବାକୀ ଦିଇ ।

&lt;start&gt; &lt;end&gt;

**Application to Translation:** The model focused on translating English to French, demonstrating the effectiveness of sequence-to-sequence learning in neural machine translation.

**Special End-of-Sentence Symbol:** Each sentence in the dataset was terminated with a unique end-of-sentence symbol ("<EOS>"), enabling the model to recognize the end of a sequence.

**Dataset:** The model was trained on a subset of 12 million sentences, comprising 348 million French words and 304 million English words, taken from a publicly available dataset.

**Vocabulary Limitation:** To manage computational complexity, fixed vocabularies for both languages were used, with 160,000 most frequent words for English and 80,000 for French. Words not in these vocabularies were replaced with a special "UNK" token.

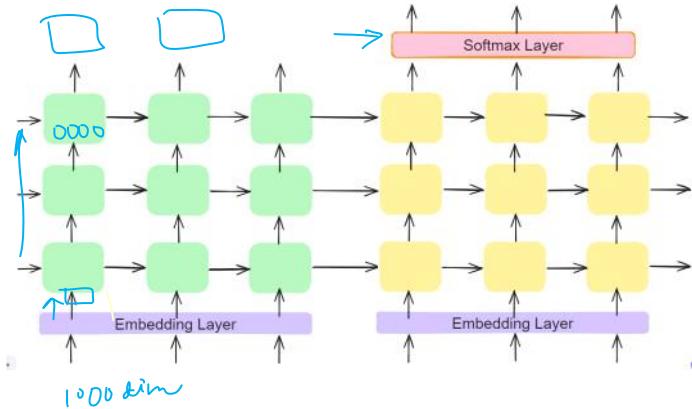
**Reversing Input Sequences:** The input sentences (English) were reversed before feeding them into the model, which was found to significantly improve the model's learning efficiency, especially for longer sentences.

**Word Embeddings:** The model used a 1000-dimensional word embedding layer to represent input words, providing dense, meaningful representations of each word.

**Architecture Details:** Both the input (encoder) and output (decoder) models had 4 layers, with each layer containing 1,000 units, showcasing a deep LSTM-based architecture.

**Output Layer and Training:** The output layer employed a Softmax function to generate the probability distribution over the target vocabulary. The model was trained end-to-end with these settings.

**Performance - BLEU Score:** The model achieved a BLEU score of 34.81, surpassing the baseline Statistical Machine Translation (SMT) system's score of 33.30 on the same dataset, marking a significant advancement in neural machine translation.



# Attention Mechanism

Lecture: 60

## The Why

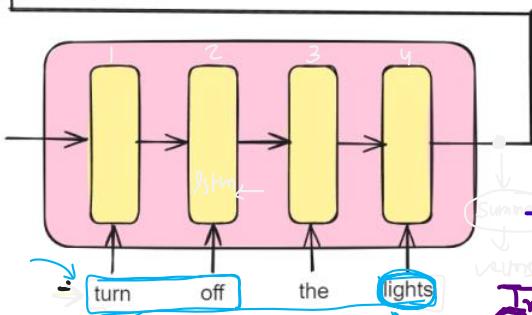
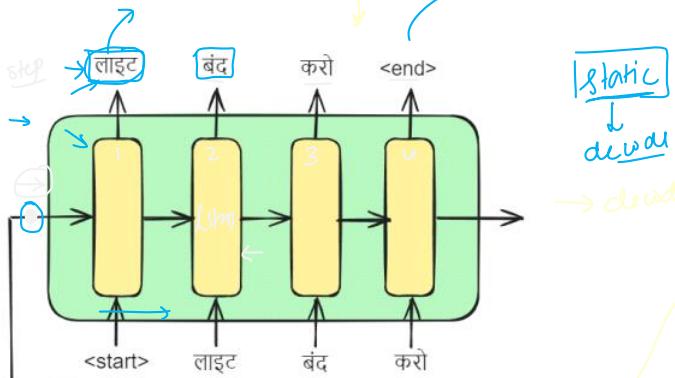
20 December 2023 13:35

improvement of encoder-decoders

dynamic

50 words

Ehudor



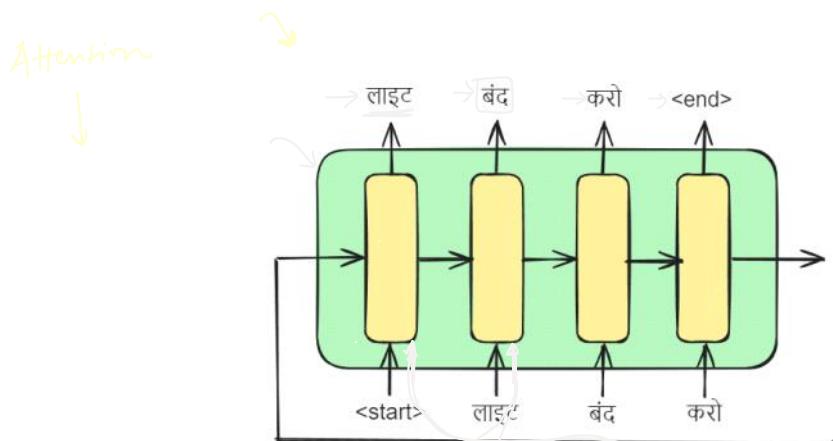
Once upon a time in a small Indian village, a mischievous monkey stole a turban from a sleeping barber, wore it to a wedding, danced with the bewildered guests, accidentally got crowned the 'Banana King' by the local kids, and ended up leading a vibrant, impromptu parade of laughing villagers, cows, and street dogs, all while balancing a stack of mangoes on its head, creating a hilariously unforgettable spectacle and an amusing legend that the village still chuckles about every monsoon season.

ये यहाँ से a रातु चाकू 1-जश यकू काढ़ु जन्म  
language (उ) translation कम्बल पायु (र) ? वा पायु  
एउ गुड़ paragraph - गत - दुधा कम्बल वा, जापा (र)  
encoder-decoder द्वारा अवलम्ब होइया काढ़ु।)

In decoder side, अप्सा timestep एउटे उग्नि translation उड़ा जन्म  
light word कोरे यहुस्ते फूल। आप्सा, word यकू एउं लाज turn off एउटे तुरीय यहुस्ते फूल। फिर आप्सा  
अउज्ज्ञान घुर्हे sentence ट्राई decoder उड़ा बाढ़ु - पाठीछा। एउटे खजाने static representation उड़ा  
अप्सा decoder उड़ा translation बाढ़ु अप्सान्न रस्ते। अथवा, यहन आमाज्ज्ञाय लाइट एउटे English को हिन्दीग्य  
राखे अप्सा यहन घपि घुर्हे sentence उड़ा एकटो पाढ़ु (lights) एउटे attention-फिल्टर आधारमा उपलब्ध भालो राखे।  
ग्राम्स, static nature ना इसे dynamic इन्हों आवा राखे।

Information is valuable  
to individual companies in  
determine what information part  
**Information Security**  
strategy is knowledge based  
part as general however intellectual and knowledge-based assets

Once upon a time in a small Indian village, a mischievous monkey stole a turban from a sleeping barber, wore it to a wedding, danced with the bewildered guests, accidentally got crowned the 'Banana King' by the local kids, and ended up leading a vibrant, impromptu parade of laughing villagers, cows, and street dogs, all while balancing a stack of mangoes on its head, creating a hilariously unforgettable spectacle and an amusing legend that the village still chuckles about every monsoon season.



प्रृथम solution कि ? आमदौर human रुपर तोन sentence translate  
कर्ने ज्यादा translation वाले ग्रामलय दिके शाखाएँ आमदौर sentence  
अवृत्त अन्तर्भुक्त हिस्से राज्य translation कायाते वायाते शहरे।  
ऐसे में एक एक वार्ता की राज्य translation कर्ने ज्यादा शहरे।  
जैसे, light एक translation पर आमदौर decoder-के एन्ड्रोइड्स  
एक्सेस इविधु माजे आड़कता । अर्थात्, light की translation पर आमदौर decoder-के एन्ड्रोइड्स  
encoder पर एक timestep के शुरूआती स्थिति । Decoder के शुरूआती राज्य encoder पर आड़कता या छाप  
set of timestep currently decoder पर आमदौर शुरूआती स्थिति । ऐसे mechanism के attention mechanism  
-एल्टे ।

## What is attention Mechanism?

The What

21 December 2023 06:04

अथवा इसकी जड़ है उदाहरण, नेटवर्क encoder के hidden state. यह decoder परिस्थिति LSTM के output होने से current timestep का input रिसेप्ट करता है। 1, 2, 3, 4 इन्हें timestep. Normal, encoder-decoder Architecture का decoder अपना जगह कहि timestep 2 का,  $(y_1 \& s_1)$  LSTM input फिलहाल रखता है।

Attention mechanism का द्वारा timestep 2 का input एक weighted sum है।

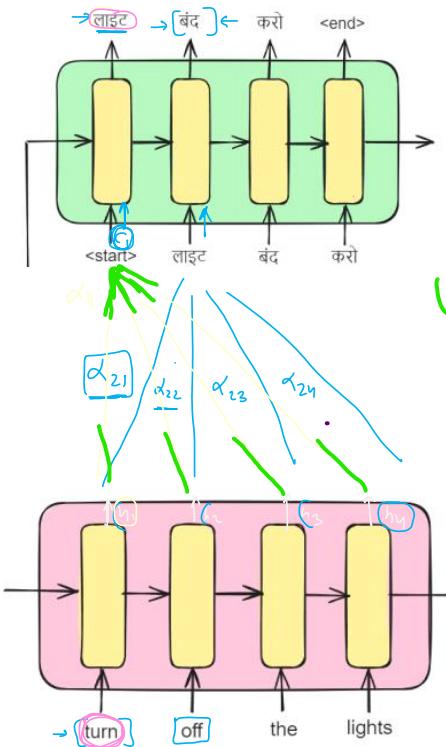
$\sum_i \alpha_{i,j} (y_i, s_j, c_i)$   $c_i$  current timestep  $i$  encoder's hidden state at timestep  $j$  is important

$\alpha_{i,j} = \frac{h_i \cdot h_j}{\|h_i\| \|h_j\|}$   $\alpha_{i,j} \in [0, 1]$   $y_i, s_j, c_i$  को  $h_i, h_j$  के scalar, vector?

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

जैसा कि  $y_i, s_j, c_i$  को  $h_i, h_j$  के scalar, vector?

यहाँ  $c_i$  current translation का  $i$  तक encoder timestep important तबसे जानका  $h_i$  को pass करता है। LSTM के रख तक  $h_i$  रख वेक्टर- $c_i$  dimension  $h_i$  के जमान रहता है। यहाँ  $h_i$  का अधिकारी है और जानका रखता है weighted sum का रखा जाता है।



decoder  
 $i=1$

$\alpha_{11}$

$\alpha_{12}$

$\alpha_{13}$

$\alpha_{14}$

$$(C_1) C_1 = \alpha_{11} h_1 + \alpha_{12} h_2 + \alpha_{13} h_3 + \alpha_{14} h_4$$

↓ vectors

जैसा कि context vector तो क्या है?

$d_{21}$

encoder

$c_i \rightarrow d_{2i}$

$$(C_2) C_2 = \alpha_{21} h_1 + \alpha_{22} h_2 + \alpha_{23} h_3 + \alpha_{24} h_4$$

encoder  
 $j$

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

$c_2$

$$1 \times j = 16 \rightarrow$$

decoder का timestep  $i$  encoder " " का रखता है।

$$\text{total alphas} = i \times j = 4 \times 4 = 16$$

फिर, परे  $\alpha$  की किसका calculation कहा रखूँ?

Let calculate,  $\alpha_{21}$  का आवश्यक alignment score or similarity score है।  $\alpha_{21} (i=2, j=1)$  का अर्थ,  $\alpha_{21}$  जुहानु, decoder का timestep 2 का, encoder का timestep 1 का असरित करने का अर्थ? फिर  $\alpha_{21}$  का मान  $s_1$  से independent रहता है। आवश्यक अपने पार्टी के असरित translation रखते होंगे एवं किसी भी कठोर,  $h_1$  के कानूनी असरित translation का रखते होंगे, यहाँ,

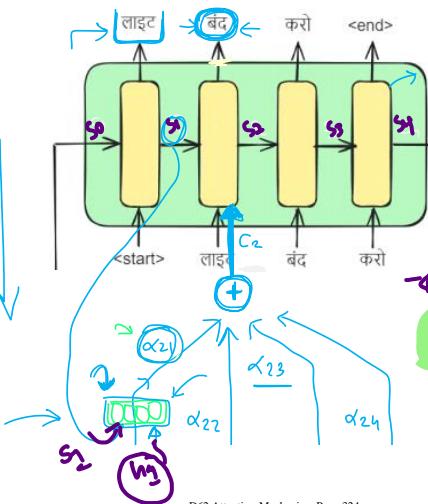
$$\alpha_{21} = f(h_1, s_1)$$

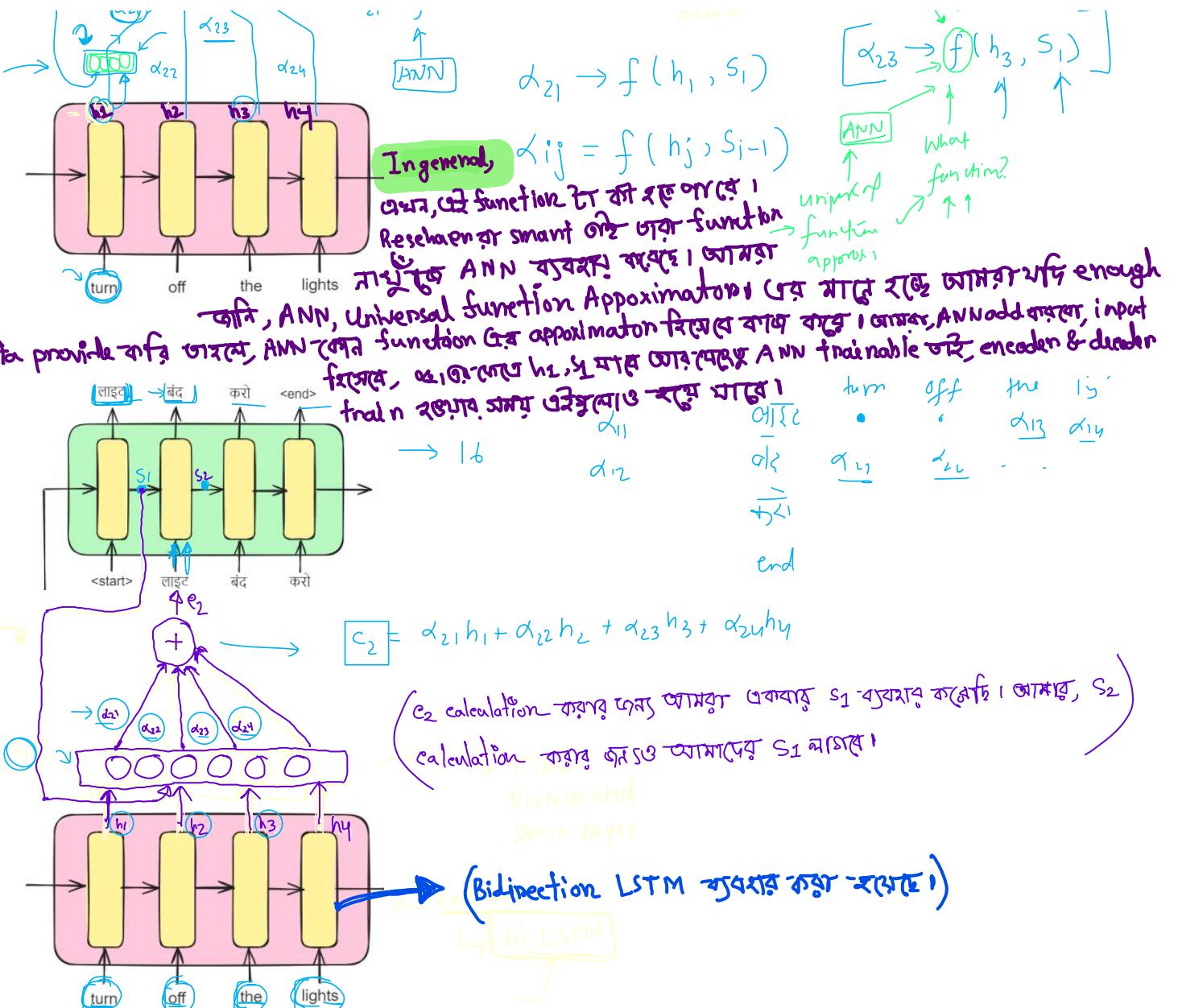
ANN

$$d_{21} \rightarrow f(h_1, s_1)$$

$$\rightarrow \begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \end{matrix}$$

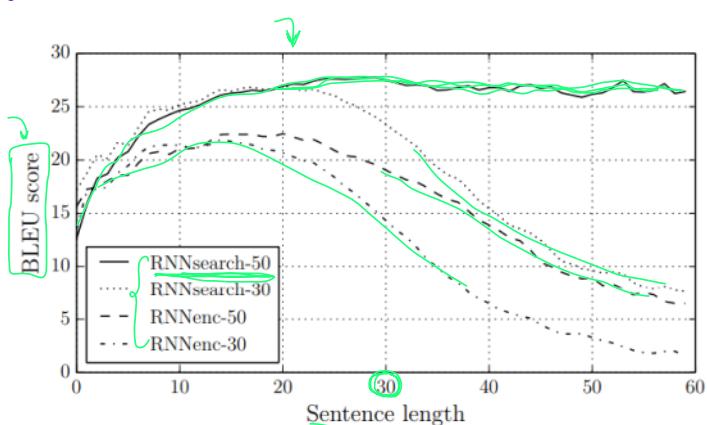
$$\left[ \alpha_{23} \rightarrow f(h_3, s_1) \right]$$





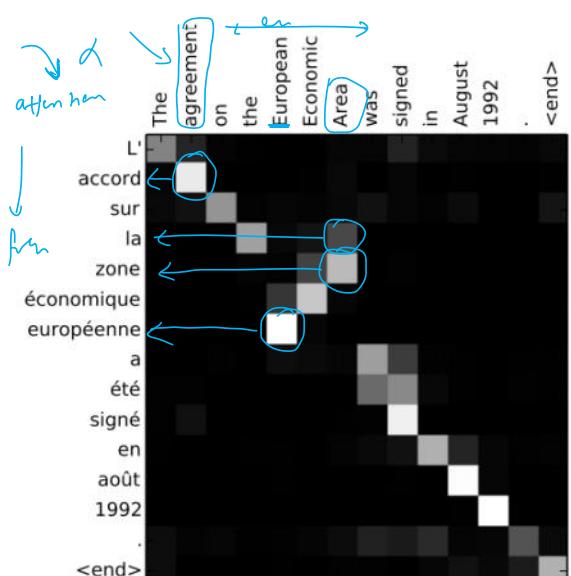
Bleu score 1 तरीके में यामाता translation के quality measure  
-करता है।

eng-fran



RNN-search-50 (Attention Mechanism)

This graph from official research paper.



(a)

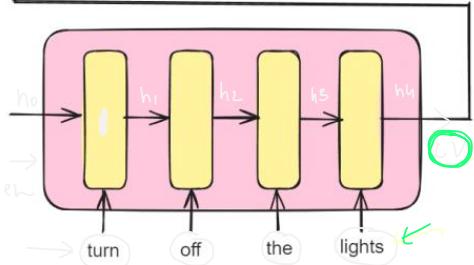
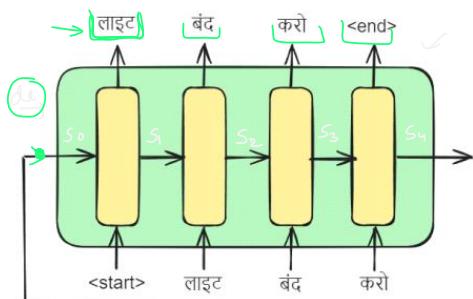
$\alpha_{21}$  value यामाता - इसका heatmap graph यामाता पाएं।

turn	$\alpha_{11}$	of	$\alpha_{12}$	the	$\alpha_{13}$	light	$\alpha_{14}$
off	$\alpha_{21}$		$\alpha_{22}$	$\alpha_{23}$	$\alpha_{24}$		
the			$\alpha_{32}$	$\alpha_{33}$	$\alpha_{34}$		
lights				$\alpha_{42}$	$\alpha_{43}$	$\alpha_{44}$	
end	$\alpha_{41}$						

# Lecture -70

Recap

16 January 2024 16:10



turn off the lights → लाइट बंद करो

## Encoder-Decoder

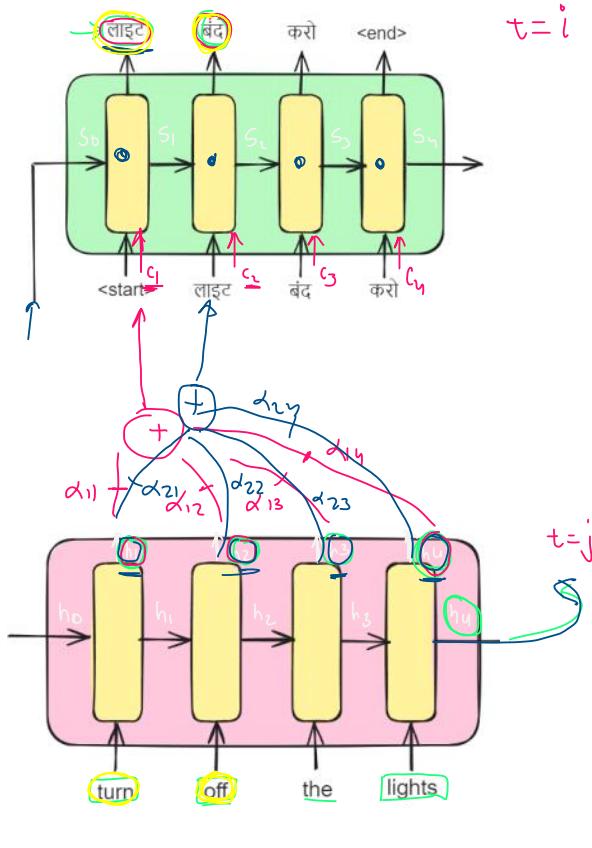
sentence > 30 words  
paragraph  
document

bilstm  
stacked lstm

translation

bottleneck → Attention mechanism

## Attention Mechanism



$c_1 \ c_2 \ c_3 \ c_4$

$$4 \times 4 = 16$$

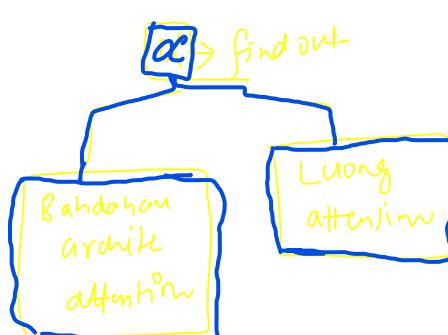
weighted sum

$$c_i^* = \sum_{j=1}^4 \alpha_{ij} h_j$$

$\alpha$  → alignment score

$$c_1 = \alpha_{11} h_1 + \alpha_{12} h_2 + \alpha_{13} h_3 + \alpha_{14} h_4$$

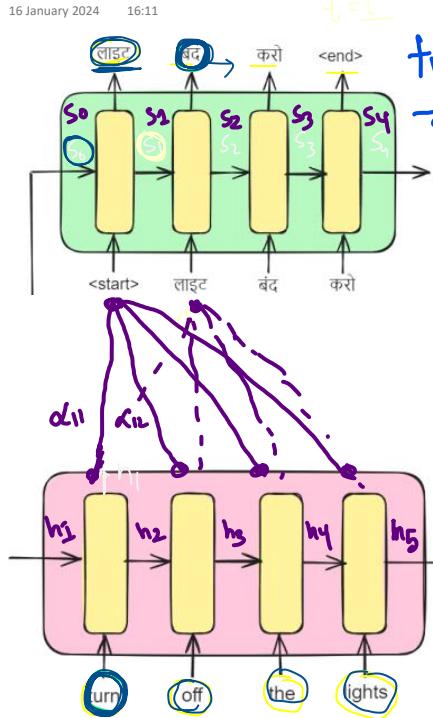
$$c_2 = \alpha_{21} h_1 + \alpha_{22} h_2 + \alpha_{23} h_3 + \alpha_{24} h_4$$



alpha find कराने वाला पथ  
पथ है !

## Bahdanau Attention

16 January 2024 16:11



$\alpha$  (attention score)

$\alpha$  (alignment score) = Decoder  $s_t$ -तों परिवर्तन timestep  $t$  के encoder  $s_i$  hidden state गुणा Weighted करने का राय ! (Based on previous hidden state + encoder timestep h.s.)

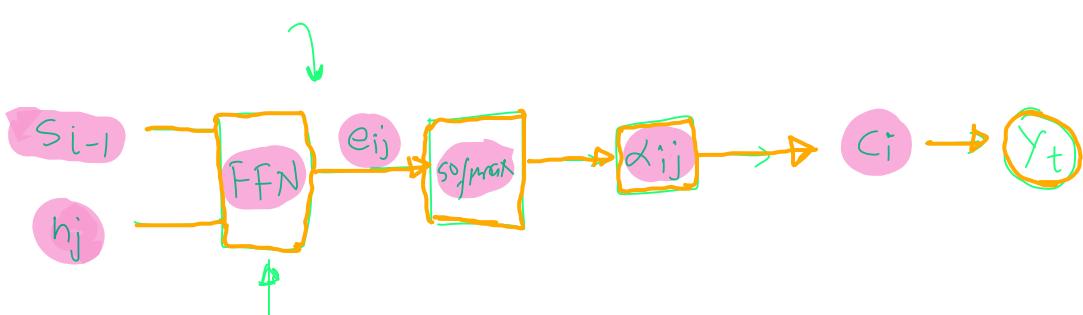
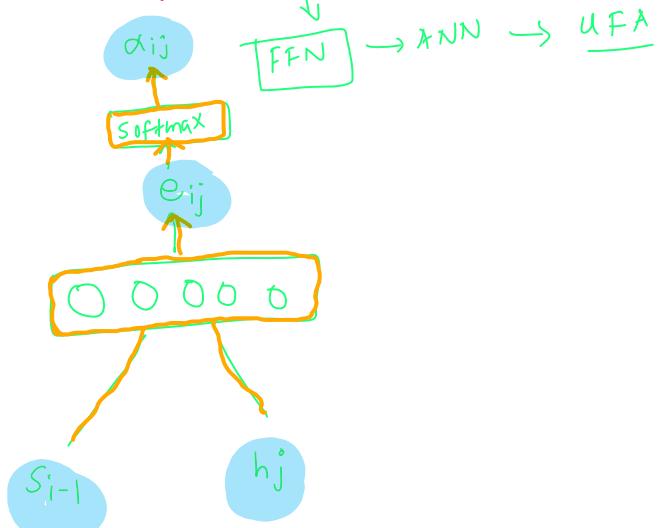
$$\alpha_{ij} = \sum_{\text{alignment}} \alpha_{ij} h_j [i \text{ decoder } s_t \text{ timestep}] \rightarrow \text{prev hidden state}$$

$$\alpha_{11} = f(h_1 + s_0) \quad \alpha_{12} = f(h_2 + s_0)$$

$$\alpha_{21} = (h_1 + s_1) \rightarrow \alpha_{22} = (h_2 + s_1)$$

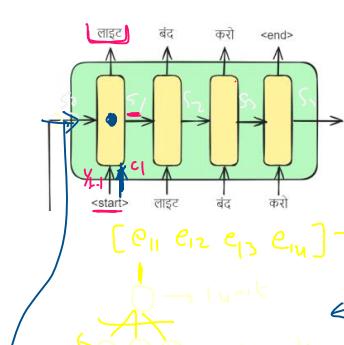
$$\alpha_{ij} = f(h_i, s_{i-1}) \quad \alpha_{21} = f(h_2, s_1)$$

यहाँ  $\alpha$  function कि रूप ? Bahdanau एवं FFN (Feed Forward Network) ANN द्वारा बहुत ही universal Function approximator द्वारा द्वयूरूप !



Flow diagram of Bahdanau

$$S_0 \ Y_{t+1} \ c_i \rightarrow \text{softmax} \rightarrow Y_t \ (\text{लाइट}) [S_1]$$



$i=1 \quad t=2$

$$s_0 = [e \ f \ g \ h]$$

$$c_1 = \sum \alpha_{ij} h_j$$

$$[ \alpha_{11} \ \alpha_{12} \ \alpha_{13} \ \alpha_{14} ] \rightarrow \text{softmax} \rightarrow [ \alpha_{11} \ \alpha_{12} \ \alpha_{13} \ \alpha_{14} ]$$

$$c_1 = \frac{\alpha_{11} h_1}{4} + \frac{\alpha_{12} h_2}{4} + \frac{\alpha_{13} h_3}{4} + \frac{\alpha_{14} h_4}{4}$$

enables synthesis  
by  $h_1, h_2, h_3, h_4$

$S_{i-1} \ h_j$

Diagram illustrating a linear combination of hidden states:

$$C_1 = \alpha_{11} h_1 + \alpha_{12} h_2 + \alpha_{13} h_3 + \alpha_{14} h_4$$

where  $\alpha = [e \ f \ g \ h]$

batch operation:  $u \times 8 \rightarrow 8 \times 3 \rightarrow u \times 3$  (using  $\tan(u \times 3)$ )

Diagram illustrating a neural network layer:

Input: turn, off, the, lights

Output:  $h_0 = [a \ b \ c \ d]$  (4dim)

$(4 \times 3) \rightarrow 3 \times 1$

$\left[ e^{e^{11}} + e^{e^{12}} + e^{e^{13}} + e^{e^{14}} \right] \rightarrow 4 \text{ numbers}$

Diagram illustrating scaling of hidden states:

4 rows / 8 cols → 8 cols

$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{11} \ h_{12} \ h_{13} \ h_{14}$
$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{21} \ h_{22} \ h_{23} \ h_{24}$
$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{31} \ h_{32} \ h_{33} \ h_{34}$
$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{41} \ h_{42} \ h_{43} \ h_{44}$

Diagram illustrating a time-distributed fully connected network:

True label:  $i=2$

$C_2 = \alpha_{21} h_1 + \alpha_{22} h_2 + \alpha_{23} h_3 + \alpha_{24} h_4$

$\alpha = [\alpha_{21} \ \alpha_{22} \ \alpha_{23} \ \alpha_{24}] \rightarrow [e_{21} \ e_{22} \ e_{23} \ e_{24}]$

alignement model:  $e_{ij} = \sum \alpha_{ij} h_j$

additive attention:  $\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^n \exp(e_{ik})}$

time distributed FNN:

$s_{11} \ s_{12} \ s_{13} \ s_{14}$	$h_{11} \ h_{12} \ h_{13} \ h_{14}$
$s_{21} \ s_{22} \ s_{23} \ s_{24}$	$h_{21} \ h_{22} \ h_{23} \ h_{24}$
$s_{31} \ s_{32} \ s_{33} \ s_{34}$	$h_{31} \ h_{32} \ h_{33} \ h_{34}$
$s_{41} \ s_{42} \ s_{43} \ s_{44}$	$h_{41} \ h_{42} \ h_{43} \ h_{44}$

Diagram illustrating the attention mechanism:

Input: turn, off, the, lights

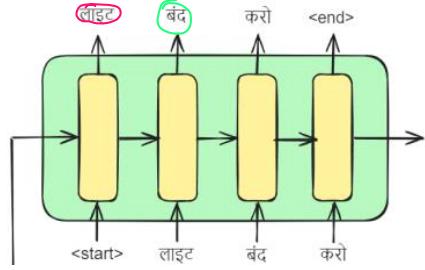
Hidden states:  $h_1, h_2, h_3, h_4$

Attention weights:  $e_{ij} = V \tanh(W[s_{ij} : h_j] + b)$

# Luong Attention, Bahdanau એન્ટેવન્ટેશન એન્ડ્રોયુન્ટ એન્ડ્રોયુન્ટ improve કરીએ ।

## Luong Attention

17 January 2024 00:09



Luong

parameters → slow

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

$$f(s_{i-1}, h_j) \times$$

$$\alpha_{ij} = f(s_i, h_j)$$

updated info current ① diff

$$s_i = [a \ b \ c \ d]$$

$$h_j = [e \ f \ g \ h]$$

dynamic

to adjust

$$a_i j$$

$$\text{softmax} \rightarrow e_{ij}$$

$$[ae + bf + cg + dh]$$

scale → attention

$$c_i = \sum \alpha_{ij} h_j \rightarrow \left[ V \tan(w[s_{i-1}, h_j] + b) \right]^T$$

(T, transpose)

dot product

$$\begin{bmatrix} s_i^T \\ h_j \end{bmatrix}$$

$$\begin{bmatrix} a & b & c & d \end{bmatrix} \begin{bmatrix} e & f & g & h \end{bmatrix}$$

$$[ae + bf + cg + dh]$$

$$e_{ij}$$

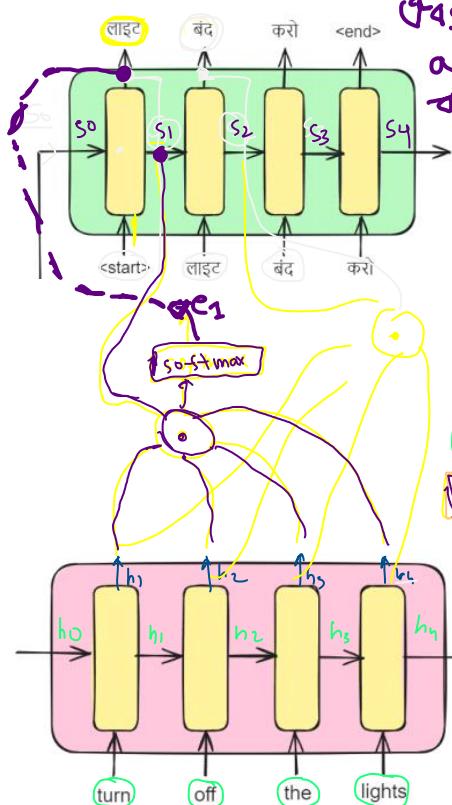
$$\text{softmax} \rightarrow e_{ij}$$

④ Bahdanau એ ડેકોડર એન્ડ્રોયુન્ટ Previous hidden state ( $s_{i-1}$ ) એન્ટેવન્ટેશન

⑤ Luong " " " current " " ( $s_i$ ) " "

⑥ " " " FFN એન્ડ્રોયુંહું dot product એન્ટેવન્ટેશન

અન્ને current hidden state એન્ટેવન્ટેશન એન્ટેવન્ટેશન up to date info આપ્દી એન્ટેવન્ટેશન dynamic રૂએ નાને એન્ટેવન્ટેશન એન્ટેવન્ટેશન Luong એન્ટેવન્ટેશન accuracy એન્ટેવન્ટેશન bahdanau એન્ટેવન્ટેશન એન્ટેવન્ટેશન | Dot product five એન્ટેવન્ટેશન એન્ટેવન્ટેશન એન્ટેવન્ટેશન એન્ટેવન્ટેશન એન્ટેવન્ટેશન



$$[e_{21} \ e_{22} \ e_{23} \ e_{24}] \text{ softmax} \rightarrow \alpha_{21} \ \alpha_{22} \ \alpha_{23} \ \alpha_{24}$$

$$c_2 = \sum \alpha_{2j} h_j$$

$$S_2 h_1 \ S_2 h_2 \ S_2 h_3 \ S_2 h_4$$

$$e_{11} \ e_{12} \ e_{13} \ e_{14}$$

$$\alpha_{11} \ \alpha_{12} \ \alpha_{13} \ \alpha_{14} \rightarrow c_1$$

એન્ટેવન્ટેશન  $s_0$  એન્ટેવન્ટેશન ( $h_i$ ) dot product કરુંનો | એન્ટેવન્ટેશન  $e_{ij}$  એન્ટેવન્ટેશન softmax layer એન્ટેવન્ટેશન  $c_i$  calculate કરુણા|

Bahdanau એ s\_iનું previous hidden state એન્ટેવન્ટેશન current hidden state

જોણે, આમણે  $c_1$  એ  $s_0$  એન્ટેવન્ટેશન એન્ટેવન્ટેશન એન્ટેવન્ટેશન | એન્ટેવન્ટેશન,  $c_1$  add કરુણો  $s_0$  એન્ટેવન્ટેશન એન્ટેવન્ટેશન એન્ટેવન્ટેશન, softmax લાગુણો then 1st word એન્ટેવન્ટેશન print કરુણો |