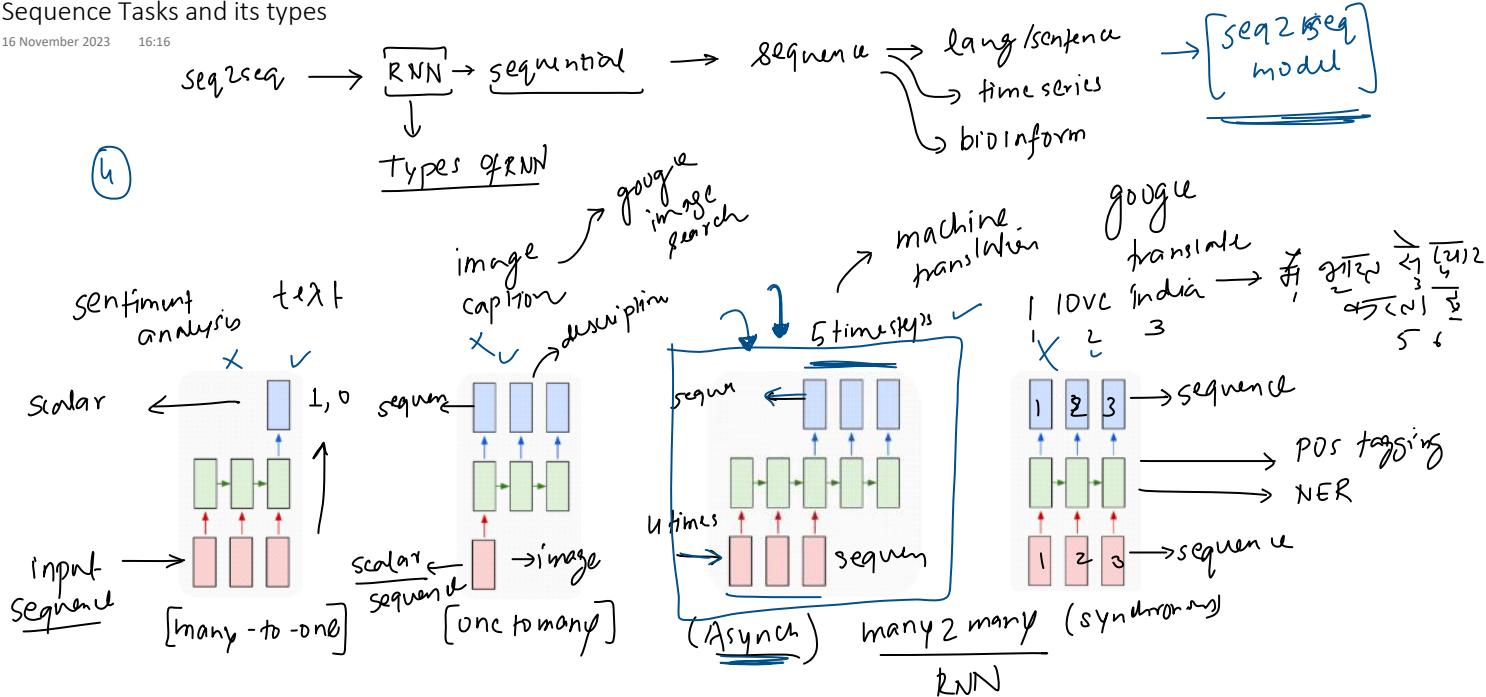


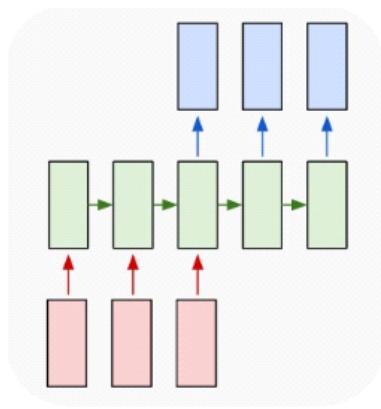
Sequence Tasks and its types

16 November 2023 16:16



Seq2Seq tasks

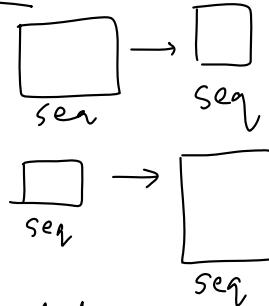
16 November 2023 16:16



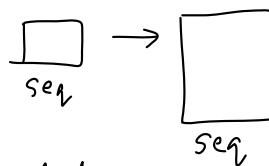
NLP

Seq2seq → machine trans

1) text summariz →



2) Question answer →



knowledge base

3) chatbot → input (text) → output (text)

4) speech-to-text →

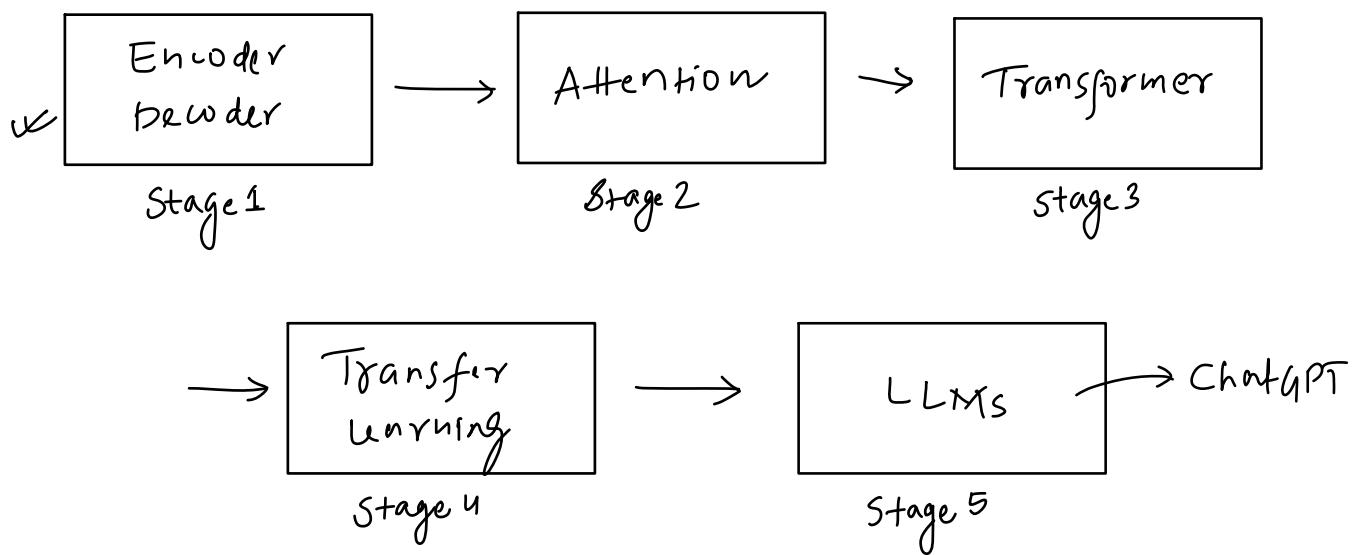
seq seq

→ seq2seq → chatgpt

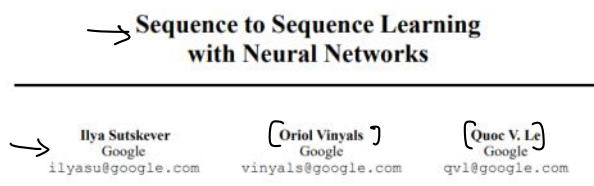
History of Seq2Seq Models

16 November 2023 16:16

ChatGPT



2014 Seminal

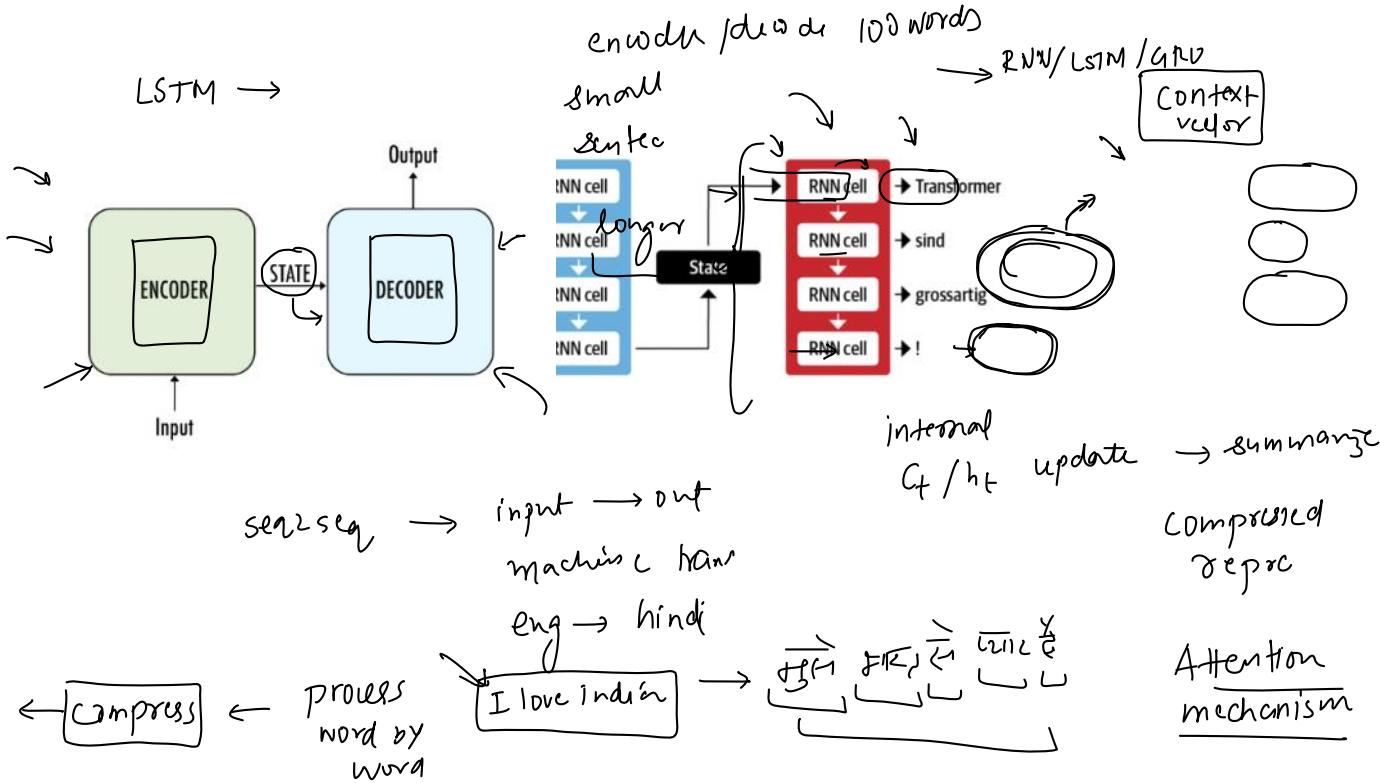
**Abstract**

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.



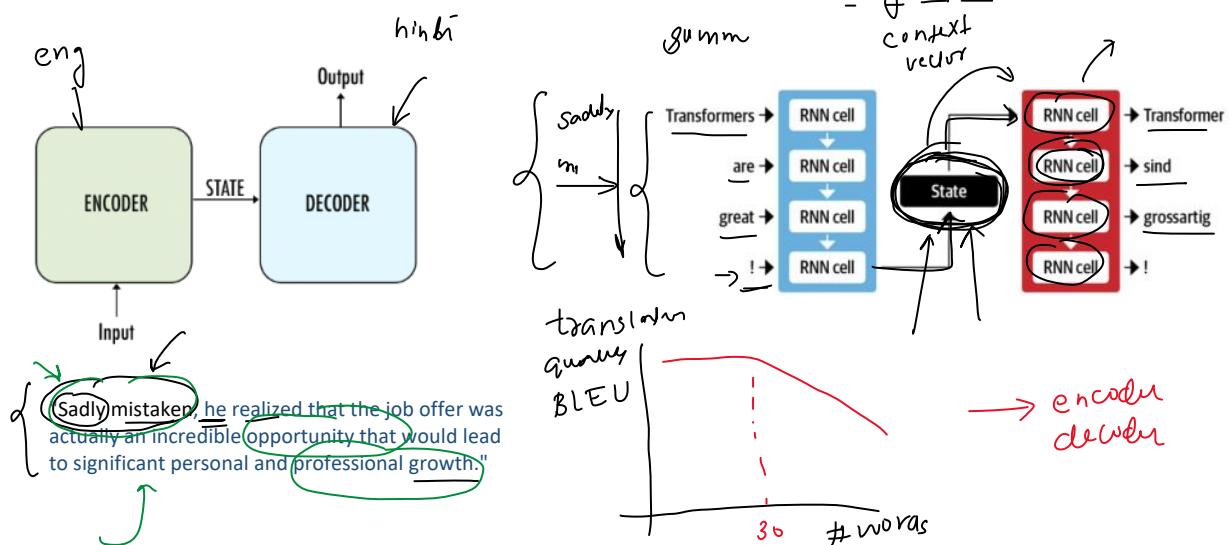
seq2seq
↓
diff
↳ encoder
decoder

[Ilya Sutskever] → cofounder openAI



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

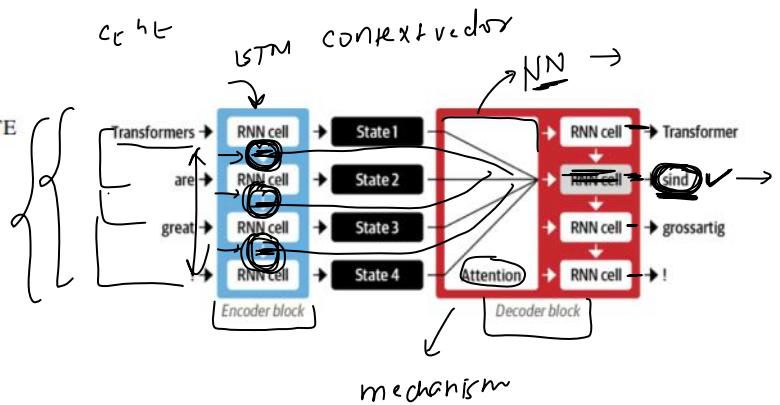
Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently 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.

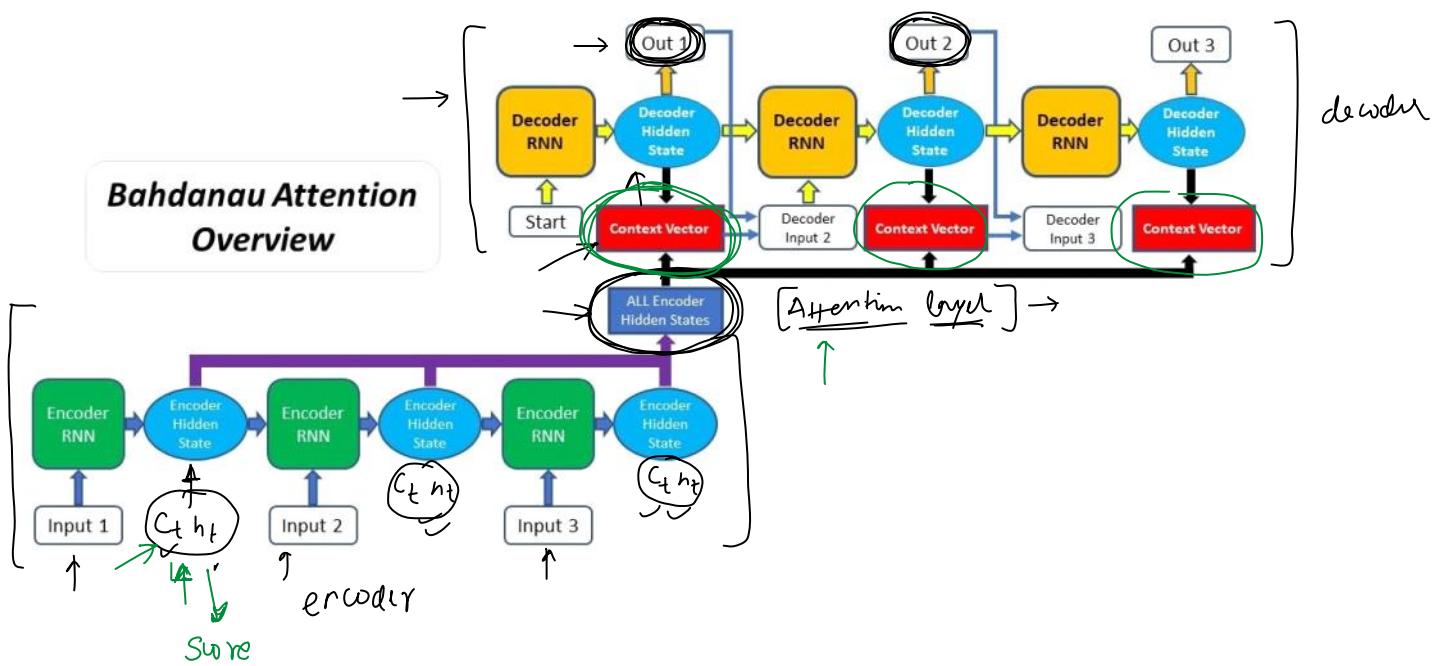
1 INTRODUCTION

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

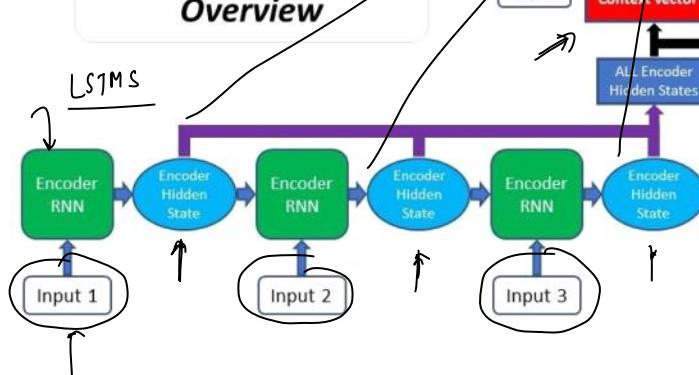
2015-2017

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* [†]	Lukasz Kaiser*
Google Research	University of Toronto	Google Brain
llion@google.com	aidan@cs.toronto.edu	lukaszkaiser@google.com

Ilia Polosukhin*[‡]
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.

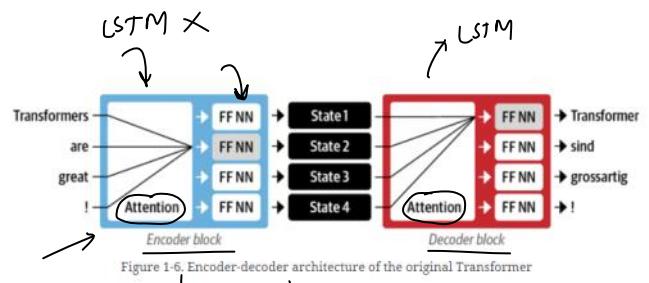


Figure 1-6: Encoder-decoder architecture of the original Transformer

LSTM / RNN cell
Attention
Self-attention
stage

arch

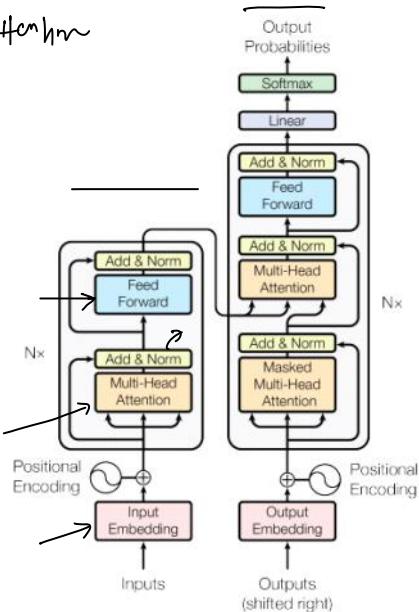
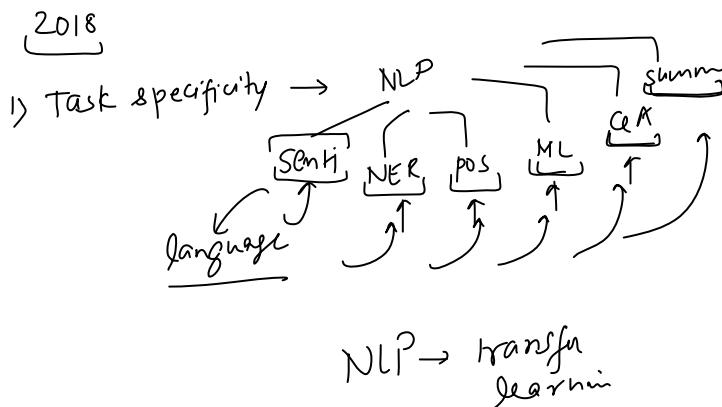
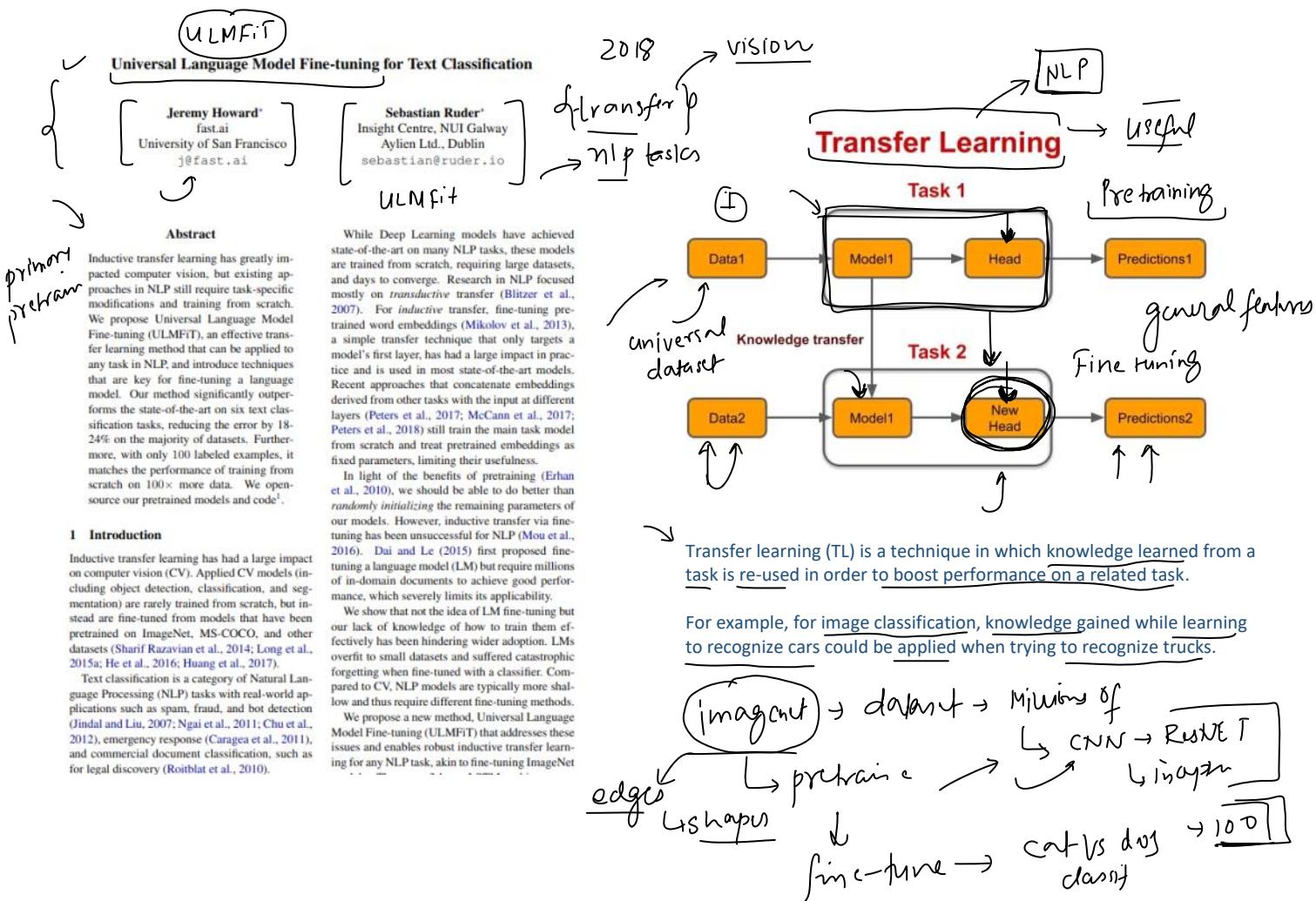
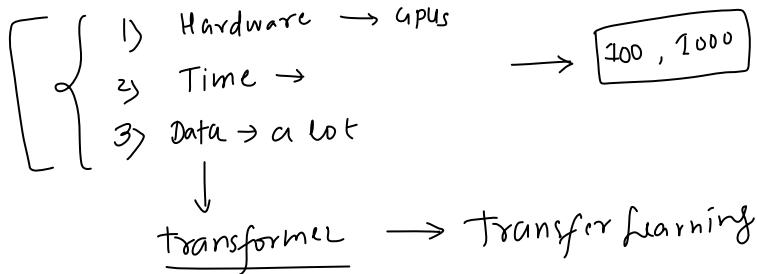
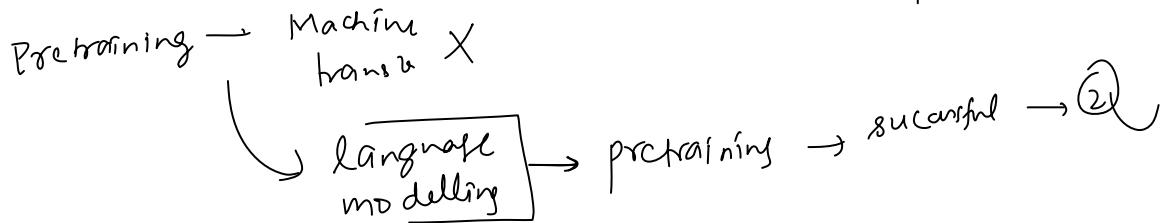
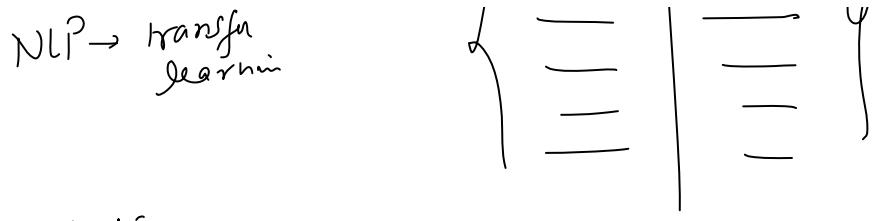


Figure 1: The Transformer - model architecture.

Stage 4 - Transfer Learning

20 November 2023 15:39





NLP task → NLP/PL model next word pred
 I live in India. and the capital is New Delhi

1) Rich feature learning
Language modeling as a Pretraining task
 The hotel was exceptionally clean, yet the service was bad ↓
pathetic

→ know trans
 ↓
 text classif / ques. | textsum) NLP / PLM

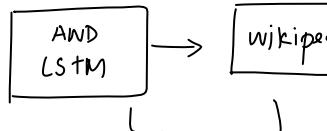
mt (kew → supervised
 eng | hin labeled
 → unsupervised task

2) Huge avail of data
 pdf → dataset
 labelling

fine tuning

[ULMFIE]

X transformer



Unsupervised
 pretrain
 Language
 modeling

classifier

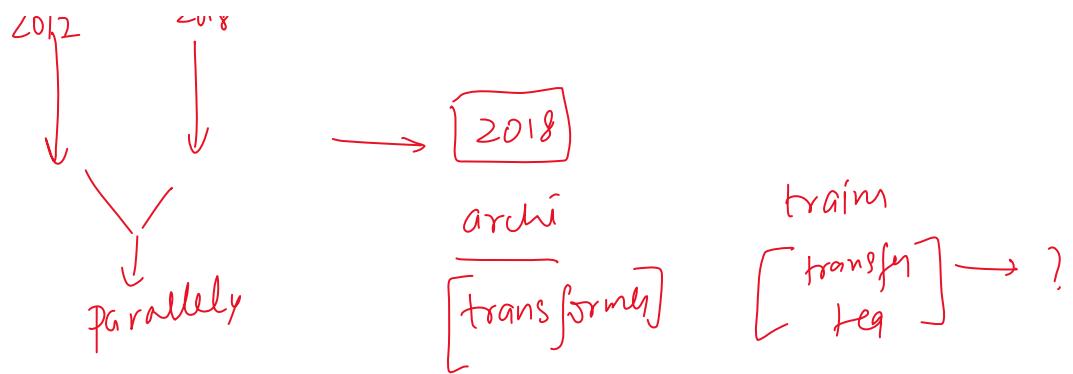
imdb
 yelp
 new dataset

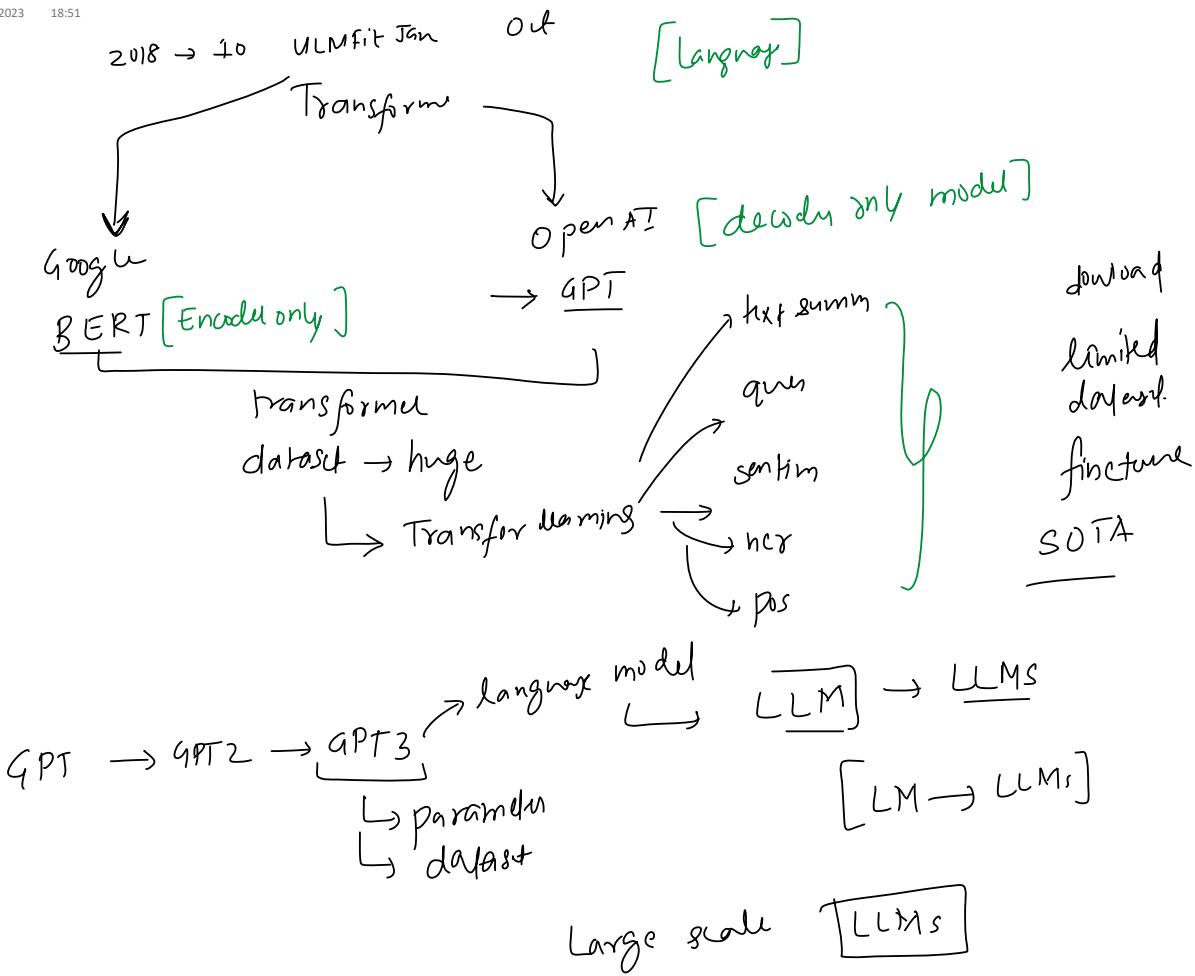
Scratch → 10000 rows
 100 row → better →

model
 ↓
 test

State of the art

2012 2018



Qualities of LLMs

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

2) Hardware → Cluster of GPU → GPT3 → Supercomputer → 100s NVIDIA GPU

3) Training → days to wccs

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

4) energy consumption
 ↗ GPT3 → ...

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

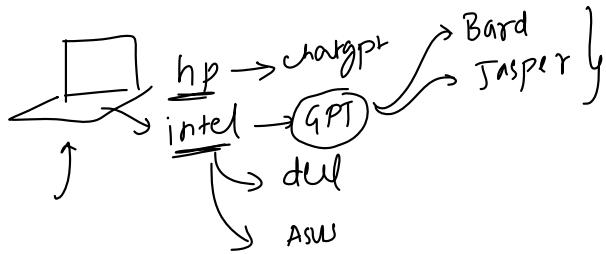


GPT3 → ChatGPT

[diff]

(GPT) → model
[ChatGPT] → application
Chat-bot

NOV



GPT3 → [ChatGPT]

1) RLHF → Reinforcement learning from human feedback

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

===== y labeled
===== y
===== y
===== y

2) Incorporate safety and ethical guideline

- + minimize bias

3) improvement in contextual point

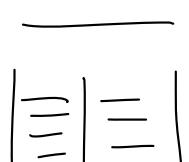
context → maintain context

dialogue
convrs

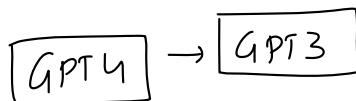
4) Dialogue specific training

- + 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 $\boxed{\dots}$