

Assignment - 3

Generative AI Assignment

Page No.

Date:

Generative AI - ~~the~~ the life cycle walks you through the

FLAN-T5 individual stages & decisions you have to make when you're developing Generative AI
→ Finetuned applications

language NLP it is a subset of traditional machine learning
text-to-text

transformer Large Language models (LLM) :- BERT, GPT, FLAN-T5.
LLM4, PaLM, BLOOM

LLMA

Large language model → In this course, learn how they are built and trained,
Model Hub API * how you can interact with them via text known as
prompt

- * How to fine tune models for your use case and data
- * How you can deploy them with applications to solve
your business & social tasks

- The text you pass in the LLM is called Prompts
- The space or memory that is available to the
prompt is called the context window

* Generative Algorithms are not new, previous generations of
language models made use of an architecture called
RNNs (were limited by the amount of compute &
memory needed to perform well at generative tasks)

Google & the University of Toronto → Transformer Architecture
had arrived

Transformers:- Scale Efficiently, Parallel process
Attention to input meaning

Output

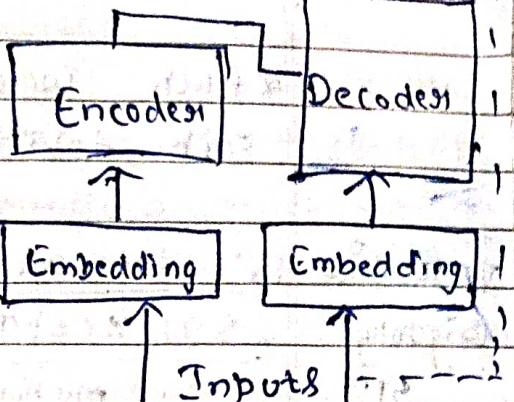
Page No.

Date :

Softmax
Output

Transformers architecture :-

↳ Attention is all you need



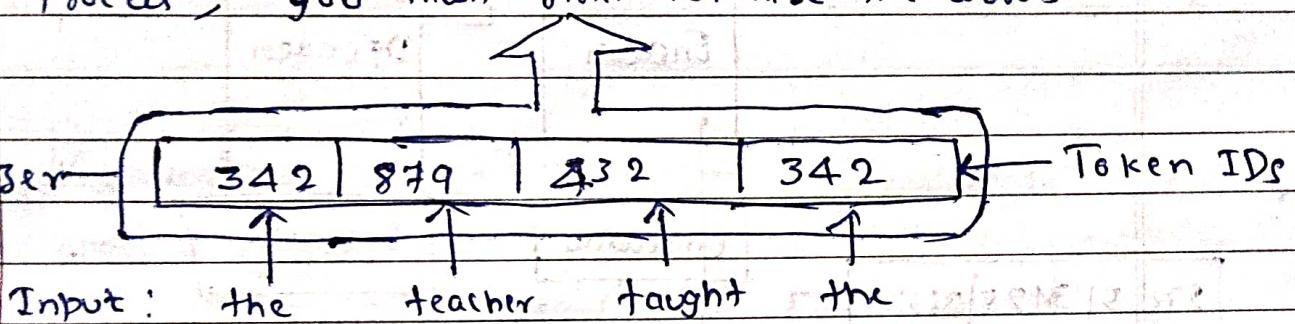
Machine learning models are just

big statistical calculators &

they work with numbers, not words. So

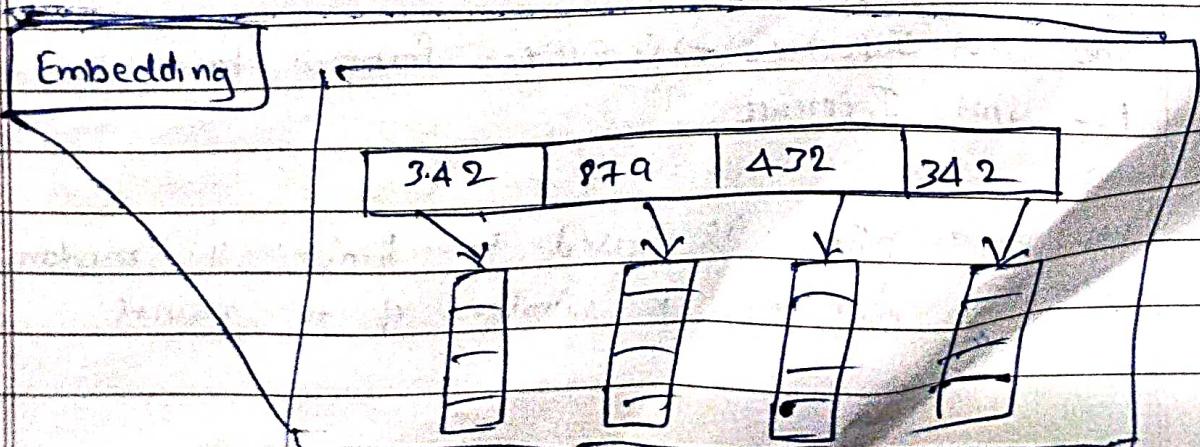
before passing text into the model to

process, you must first tokenize the words

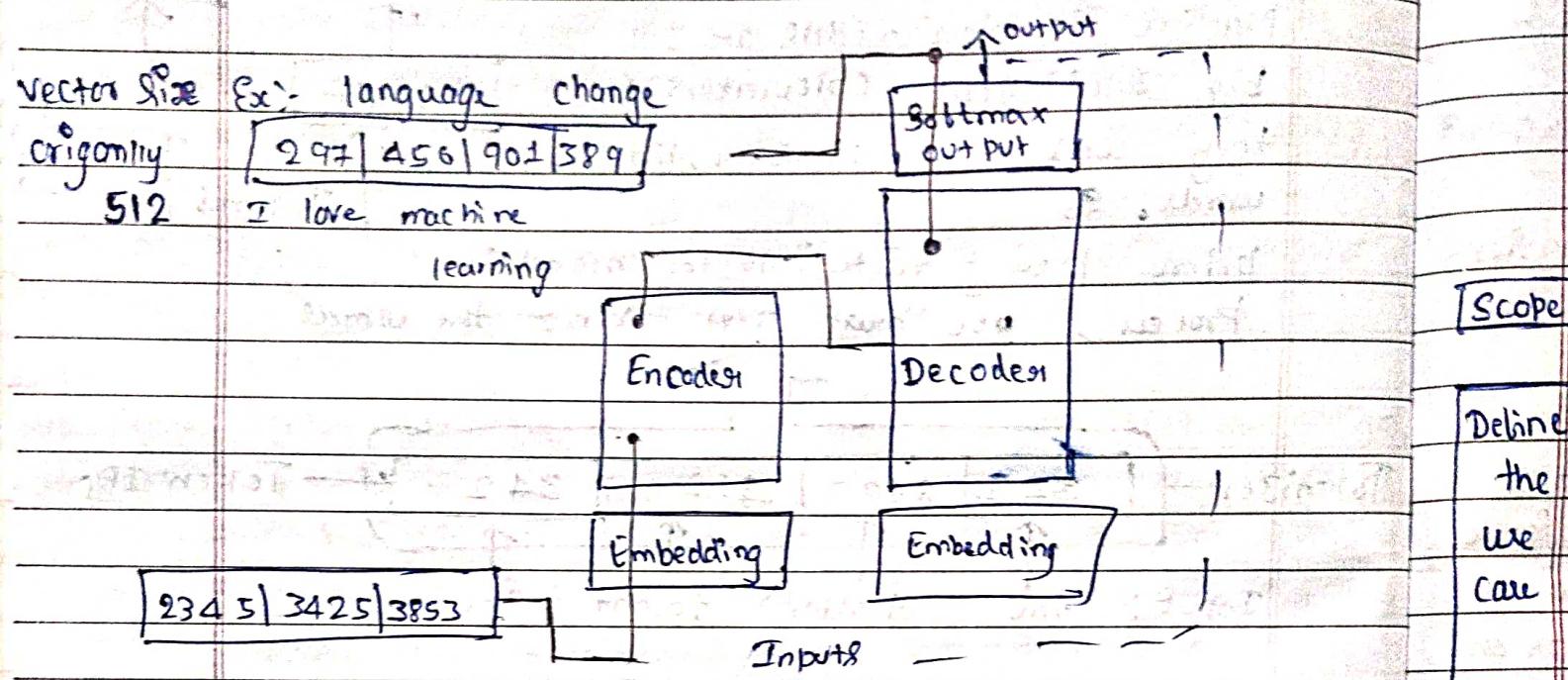


Transformers work:-

- *> Input represented numbers, pass it to the Embedding layer
- *> Embedding layer is a trainable vector space, a high-dimensional space where each token is represented as a vector & occupies unique location within that space. Each tokenID is matched to a multi-dimensional vector, these vectors learn to encode the meaning & context of individual tokens in the input sequence.



- * Each word has been matched to a tokenID, & que
Each token is mapped into a vector



J'aime l'apprentissage

automatique

[BERT] Encoder:- encodes inputs ("Prompts") with contextual understanding and produces one vector per input token. [work as Sequence-to-Sequence models]

[ELAN-TG] Decoder:- Accepts input tokens and generates new tokens.

Zero Shot Inference

One Shot Inference → first question followed by a second question
 Few Shot Inference

Top - K Sampling is used to limit the random sampling from the K number of top result

zero shot → absolutely no labeled data available
one shot → each new class has one labeled ex
few shot → limited no. of labeled examples for each new class

Page No.

Date:

Ques: Temp: It is used to affect the randomness of the output of the softmax layer. A lower temp results in reduced variability while a higher temp results in ↑ randomness of the output.

Generative AI Project lifecycle

Scope	Select	Adapt & align model	Application integration
Define the we care	Choose an existing model or pretrain Your own	Prompt Engineering Fine-tuning Align with human feedback	Optimize & deploy model for inference Augment model & build LLM-powered applications

python libraries

- `torchdata` → It helps with the data loading and some other aspects for PyTorch specific to datasets
- `Transformers` → This is a library from Hugging Face, (open source tool for LLMs)
- `Datasets` → It can load in many of the common public datasets that people use to either train models, fine tune models or just experiment

Tokenizers are all come from the Hugging Face library

AutoEncoding Models

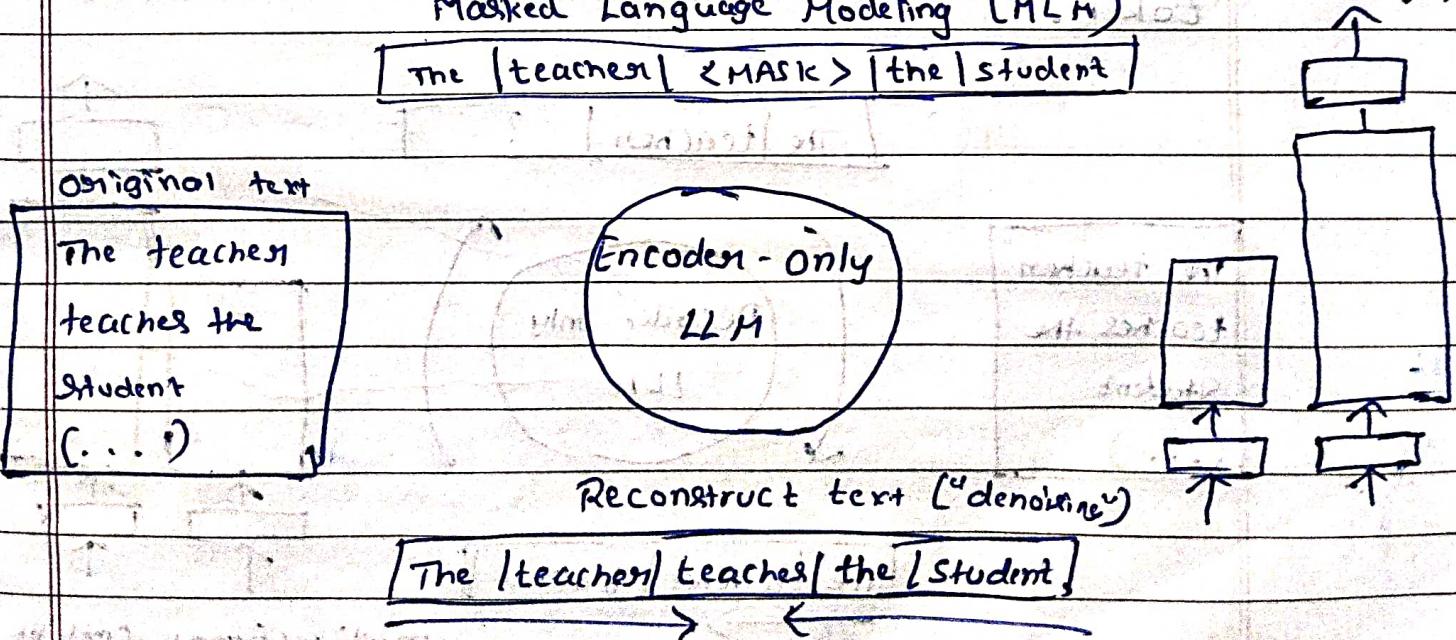
The Encoder-only Variants of Transformers are also called autoEncoding models.

They are pre-trained using Masked Language Modeling (MLM). In MLM, tokens in the input sequence are randomly masked and the training objective is to predict the masked tokens in order to reconstruct the original input sequence.

This is also called a denoising objective since the masking of the tokens can be thought of as adding noise to the input sequence and then predicting the masked tokens can be thought of as removing that noise from the input sequence.

Masked Language Modeling (MLM)

The [teacher] <MASK> [the] student

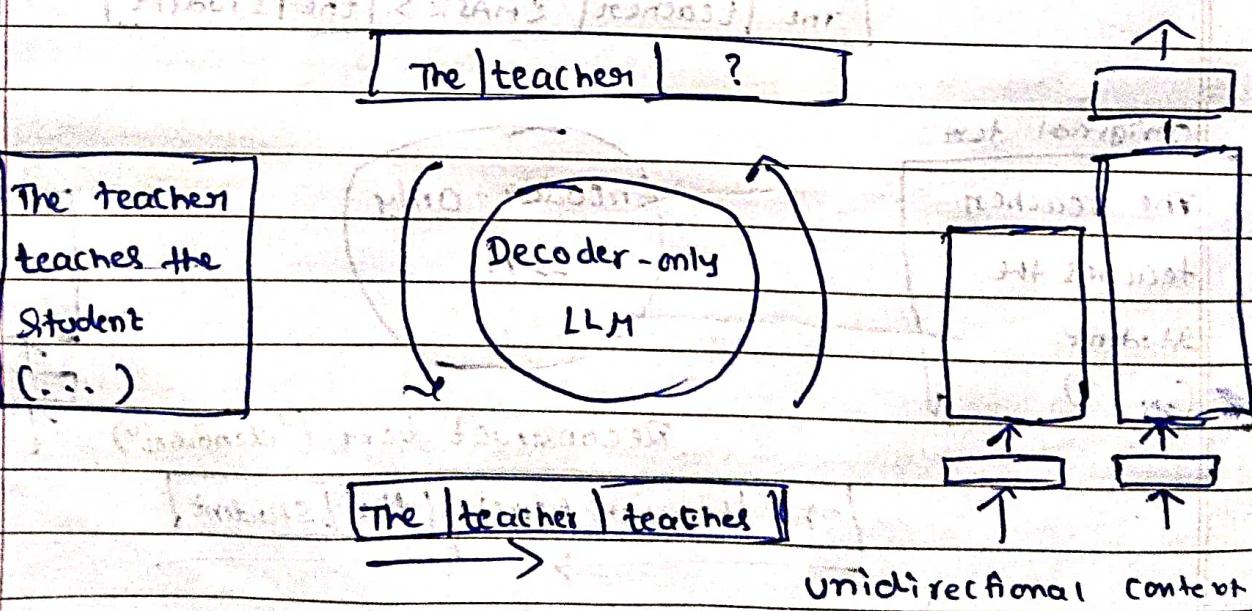


Example: BERT, ROBERTA

Decoder-only Models (Autoregressive Models)

The decoder-only variants of Transformers are also called autoregressive models.

They are pre-trained using Causal Language Modeling (CLM). In CLM, the training objective is to predict the next token based on the previous sequence of tokens. The tokens of the input sequence are masked and the model can only see the input tokens leading up to the token being predicted at the moment. The model has no knowledge of the tokens that come after this token. The model then iterates over the input sequence one-by-one to predict the next token. Thus, in contrast to autoencoding models, the model builds a unidirectional context for each token.

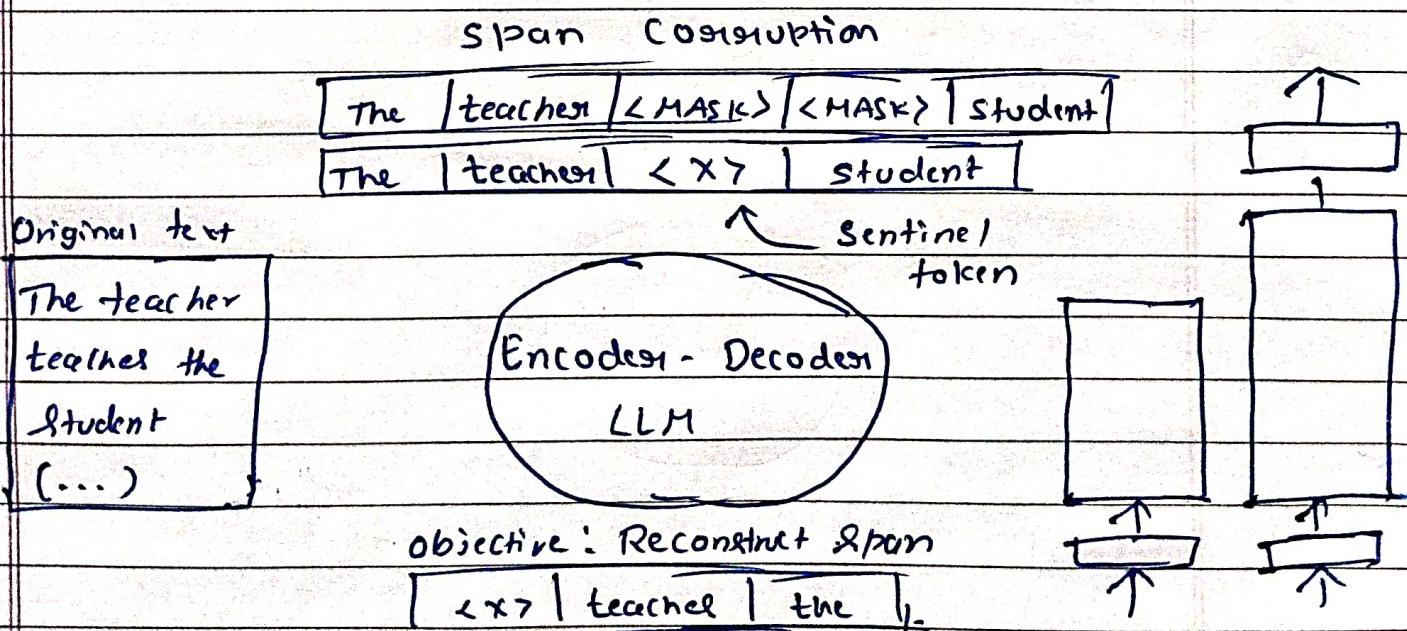


Ex: GPT, BLOOM

Encoder - Decoder Models (Sequence - to - Sequence Models)

The Encoder-decoder variants of Transformers are also called Sequence-to-Sequence models.

The exact details of pre-training objective vary from model to model. For ex, FLAN-T5 is trained using Span corruption. In span corruption, a part of the input sequence is masked and replaced by a sentinel token. These sentinel tokens are special ~~or reserved~~ tokens added to the vocabulary that do not correspond to any actual word from the dataset. The decoder then has to reconstruct the sentence autoregressively, the output is the sentinel token followed by the predicted token.



Ex:- FLAN-T5, BART

Week-2

Fine tuning LMs with instruction

- 1) Instruction Fine-tuning :- when you have your base model, the thing that's initially Pretrained it'll encoded a lot of really good information, [it doesn't necessarily know how to be able to respond to our prompts]
* it helps it to be able to change its behavior to be more helpful for us.

- 2) Parameter Efficient fine-tuning :- So we have a lot of customers that do want to be able to tune for very specific tasks, very specific domains, this be a great way to achieve similar performance results on a lot of tasks.

Fine tuning is a supervised learning process where you use a data set of labeled examples to update the weights of the LLM. The labeled examples are prompt completion pairs.

- Instruction fine-tuning is particularly good at improving a model's performance on a variety of tasks.

During fine tuning, you select prompts from your training data set and pass them to the LLM, which then generates completions. Next you compare the LLM completion with the response specified in the training data.

Catastrophic Forgetting

Page No.

Date :

* one way to mitigate catastrophic forgetting is by using regularization techniques to limit the amount of change that can be made to the weights of the model during training. This can help to preserve the information learned during earlier training phases & prevent overfitting to the new data.

* Single task fine-tune

* Multiple task fine-tune

L4) Evaluation - Metrics

Rouge :- used for text summarization,

- compares a summary to one or more reference summaries

BLEU Score :- used for text translation

- compared to human-generated translations

PEFT (Parameter Efficient Fine-Tuning)

It update updates only a small subset of parameters so it's more robust against this catastrophic forgetting effect

* Parameter Efficient Fine-tuning (PEFT) methods specifically attempt to address some of the challenges of performing full fine-tuning

↳ Storage requirements :- with PEFT, we can change just a small amount of parameters when fine-tuning, so during inference you can combine the original model with the new parameters.

Computational Constraints: Because most parameters are frozen, we typically only need to train 15%-20% of original LLM weights, making the training process less expensive.

* Catastrophic forgetting or with PEFT, most parameters of the LLM are unchanged & that helps making it less prone to catastrophic forgetting.

* BLEU: Evaluation metric below focuses on precision in matching generated output to the reference text & is used for text translation.

* Multi-task finetuning requires separate models for each task being performed.

* LORA works: It decomposes weights into two smaller rank matrices & trains those instead of the full model weights.

Scaling model training steps

* The first step in scaling is model training is to distribute large datasets across multiple GPUs and process these batches of data in parallel.

FML

Submitted by : Karthik C

Reg. No : 240158016

Assignment -3 : Generative-A.I