

Week-1

- \* These Generative AI models are designed by mimicing Human Activity
- \* So these LLM (Large Language Models) are trained on trained using trillions of words over many weeks & months with large amounts of compute power
- \* LLM's are able to take natural language or human written instructions & perform all the tasks which are assigned
- \* The text passed to LLM is known as prompt
- \* The space or memory avail. to prompt is called the context window
- \* The output of the model is called completion
- \* Act of using the model to generate text is known as inference

LLM tasks & Use cases:-

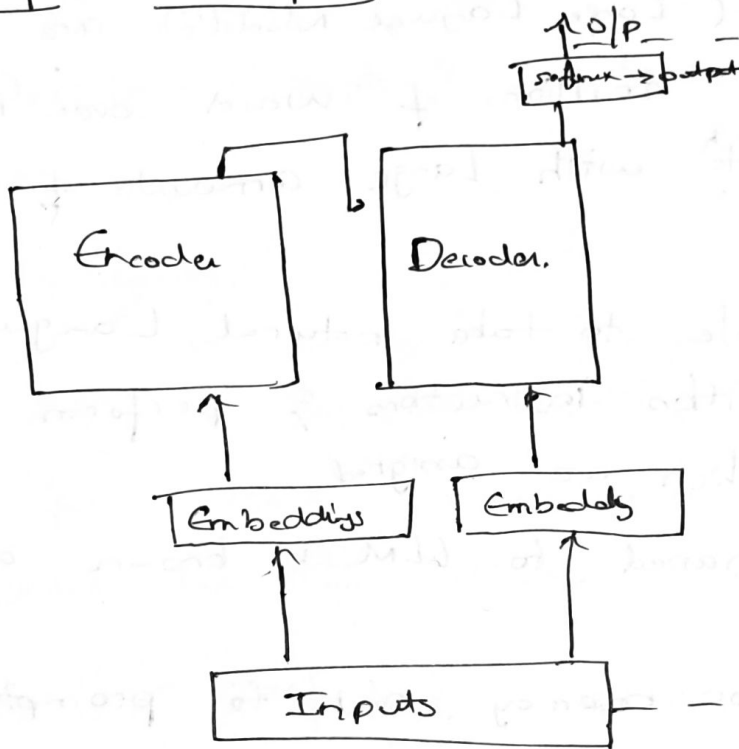
- \* It is used for Information retrieval
- \* Models can be tuned to perform specific tasks
- \* previous RNN's used to generate text
- \* Used to check relation between phrases

words,

transformer is

- It checks the relations between each of the words
- Attention map used to define attention weights.

### Simplified transformer Architecture



Types

- Encoder only
- Decoder only
- Encoder-Decoder.

→ Inputs are the words which are converted into numbers by tokenizer, where each number is a dictionary of all the words model can work with.

→ Tokenizer methods

→ token ID's matches two complete words

→ token ID's used to represent parts of words.

→ Tokenizer used to train the model we must use the same tokenizer when we

## Generate text

→ User inputs are converted into Embeds where the tokenizer is used to modify to accordingly to the desired output

\* Encoder :- It encodes input i.e. prompt with contextual understand & produces one vector per input token

\* Decoder :- Accepts input token & generates new tokens.

## Prompt Engineering:-

\* work to develop & improve the <sup>prom</sup>pt known as prompt engineering

\* In context Learning provides extra information needed for the model to perform better

zero shot ~~performance~~ - provide no example

one shot ~~performance~~ - provides one example

few shot Inference - provides few example

## Generative Configuration:-

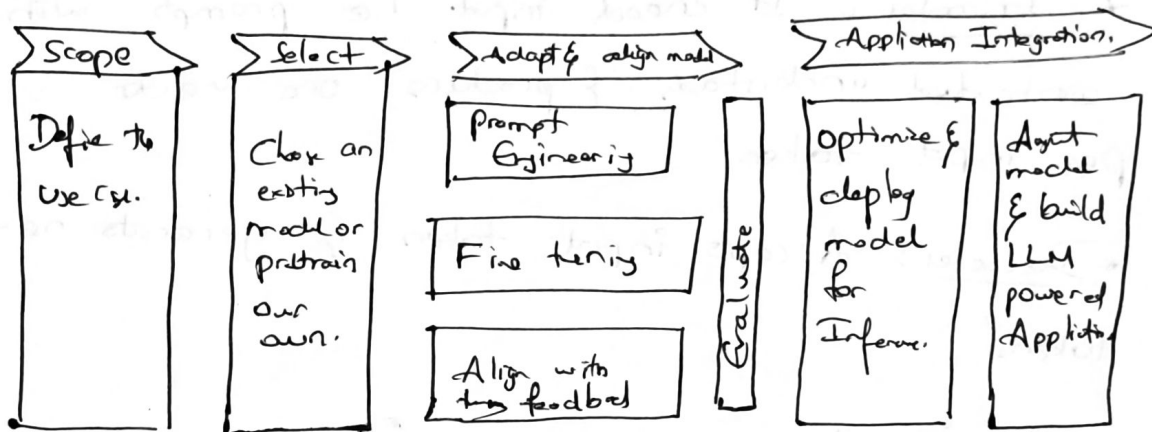
\* Max new tokens

→ Greedy vs random Sampling → Selects a token via a random word with highest probability is selected / weighted strategy.

top k	top p
→ number of tokens to randomly choose from	total probability - but we want the model to choose from.

\* temperature :- Higher the temperature Higher the randomness.

## Generative AI project Lifecycle :-



### \* Pre train LLM's :-

- \* Autoencoders models are pre-trained using masked Language Models
- \* Auto Regressive models are pre-trained using causal Language models
- \* Sequence to sequence models are pre-trained generally using span competition.

### \* Computational challenges :-

→ out of memory

Quantization is used to reduce this problem

→ to reduce required memory to store & train models.

→ Projects original 32-bit floating point numbers into

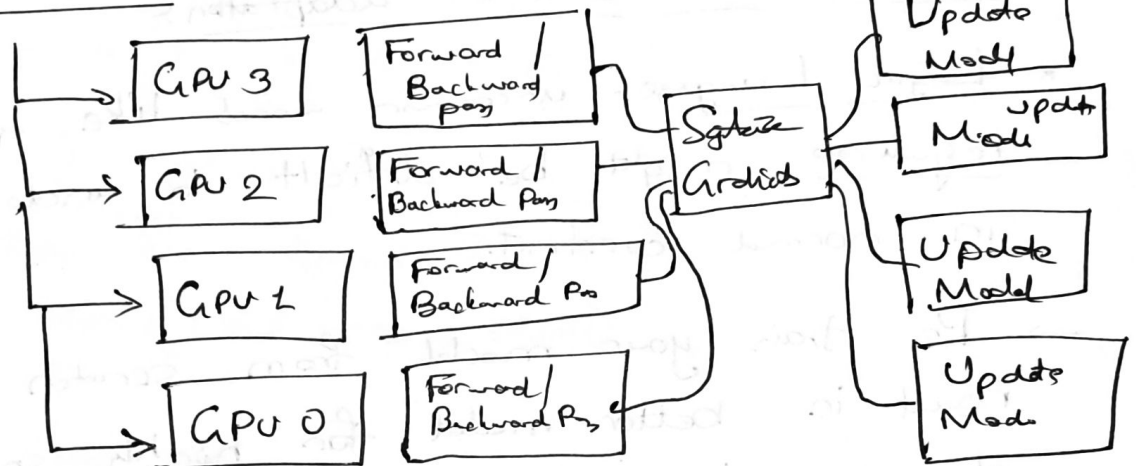
lower precision spaces

## Efficient Multi-GPU Compute strategies

Distributed Data Parallel

→ Zero Redundancy Optimizer

### Data Loader



- \* In fully shared Data Parallel requires to collect this data from all of the GPU's before the forward & Backward Pass.
- \* It helps to reduce overall GPU memory utilization.
- \* Supports offloading to CPU if needed.

### Scaling Laws & Compute-optimal Models:-

- \* Goal is pretraining to maximize model performance which is minimize model loss.
- \* It can be done by scaling choice
  - ↳ Dataset size (number of tokens)
  - ↳ Model size (Number of parameters)

→ Constraint

↳ Compute budgets:-

→ "petaFlop-day" } flops per operation performed at rate of 1 petaFlop per second for one day

\* Pre-Train for Domain adaptation:-

\* Legal Language:- uncommon words like mensae, resolute, might be difficult to understand in normal context.

→ Pre train your model from scratch will result in better model for highly specialised domains like law, medicine or finance.

Ex:- BloomBert APT was specialised in financial tasks.

## Week 2

Instruction Fine tuning:-

\* Here Fine tune LLM with Instruction prompts.

\* This has disadvantage like in context Learn may not work for smaller models

\* Example take up space in the context window.

Fine tuning is a supervised Learn process when you use a dataset of labelled examples to update the weights of LLM.

### Instructions the lens:

→ train model to demonstrate how it should respond to a specific Instruction.

Prompt [ ]      Completion [ ]

These are feeded to the LLM (Large Language Model)

### Steps for fine tuning

- 1> Prepare Your Dataset
- 2> train - test - Split
- 3> Pass it to the Model
- 4> LLM Completion.
- 5> Calculate loss using Cross-Entropy

### Fine tuning on single task:

#### \* Catastrophic forgetting:

\* Fine tuning can significantly increase the performance of a model on a specific task but can lead to reduction

in ability on other tasks.

### How to avoid Catastrophic Forgetting:-

- \* First check whether catastrophic forgetting actually impacts your use case.
- \* Fine tune on multiple task at the same time which would require more resources & compute time.
- \* Consider Parameter Efficient Fine-tuning (PEFT)
  - It preserves the weights of the optimal LLM and trains only a small number of task specific adapt layers & parameters.

### Multi-task instruction fine tuning:-

- \* Instruction fine-tuning with FLAN (Fine tuned language)
- \* FLAN models refer to a specific set of Instruction fine tuning (FLAN-TS, PALM).
- \* FLAN-TS :- It is a great, general purpose Instruct Model which uses 473 datasets across 146 task categories, includes SQuAD, A dialogue Dataset.



## Model Evaluation Metrics:

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Prediction}}$$

While the metric can be considered for small models, but while LLM's can be much more challenges.

→ Example → "Mike really loves drinks tea" & "Mike like sippp tea" are of same context, So how does measure the similarity.

\* ROUGE (recall oriented under stry for gsh evaluation, It is primarily employed to assess the quality of automatically generated summaries by comparing them to human-generated reference summaries.

\* BLEU score: (Bilingual evaluation understudy) is an algorithm designed to evaluate the quality of machine-translated text again by comparing it to human-generated translations.

## Terminology:-

Unigram → Single words

Bigram → Two words

n gram → Is a group of n-words.

$$\text{Range-1 Precision} = \frac{\text{unigram matches}}{\text{unigram in output}}$$

$$\text{Range 1 Recall} = \frac{\text{unigram matches}}{\text{unigram in reference}}$$

$$\text{Range-1 F1} = 2 * \frac{\text{Precision} * \text{recall}}{\text{Precision} + \text{recall}}$$

[ Harmonic mean ]

\* Here in unigrams order is not considered which would create conflict in some cases to overcome this we use bigrams.

\* Another method is taken the longest common subsequence present in both generated output and the reference string

LCS (gen, ref)

\* BLEU score is useful in evaluating the quality of machine-translated text

\* It is calculated using the average precision over multiple n-gram sizes & then averaged.

BLEU metric = Avg (precision across range of n-gram sizes)

Benchmark:

→ GLUE

→ SuperGLUE

→ HELM

→ MMLU

→ BIG-bench

Benchmarks for massive models

→ MMLU (Massive multitask language understanding)

→ Big Bench

→ Big Bench Hard

→ Lite

→ HELM (Holistic Evaluation of Language Models)

Parameter Efficient Fine Tuning (PEFT)

Full time tun of Large Language Models is challenging to handle weights, optimizer states, gradients, Forward Activation, temporary memory weight which would then increase the storage size to petabytes.

→ PEFT only update only a small subset of parameters

→ Some PEFT techniques represent most of the model weights & focus on fine tuning a subset of model parameters

PEFT tradeoffs include

- Parameter Efficiency
- Training Speed
- Inference Cost
- Model Performance
- Memory Efficiency

## PEFT Methods

- \* Selective :- Selects the subset of initial LLM parameters to fine-tune.
- \* Reparameterization :- of a model weights using a low-rank representation
- \* Additive :- Add trainable layers or parameters to model

LORA :- Low Rank Adaptation of LLMs

1. Freeze most of the original LLM weights
2. Inject 2 rank decomposition matrices
3. Train the weights of the smaller matrices.

## Steps to update model for Inferences

- 1) Matrix Multiplies the low rank matrices
- 2) Add to original weights.

\* Applying LORA to just the self-Attention layers of model is often enough to fine tune for a task

## Prompt tuning with soft prompts

- \* With the LORA the goal was to find an efficient way to update the weights of the model without having to train every single parameter again
- \* In prompt tuning we add additional tokens & leave it up to the supervised Learn process to determine their optimal values.
- \* The set of trainable tokens is called a soft prompt & it gets prepended to embedding vectors that represent your input text.
- \* Prompt tuning is a very parameter efficient strategy because only a few parameters are being trained
- \* Prompt tuning can be effective as full fine tuning for larger models.