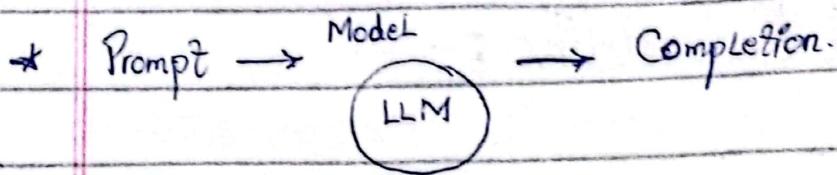


"Generative AI & LLM"



* Text Generation-

- ↳ RNN's were used before as model to generate next word in sequence.
- Problem of RNN was that it could see only one word or some previous word to predict next word, resulting wrong ans.

* Transformer Architectures-

- ↳ Tokenization.
- ↳ Vector Creation.

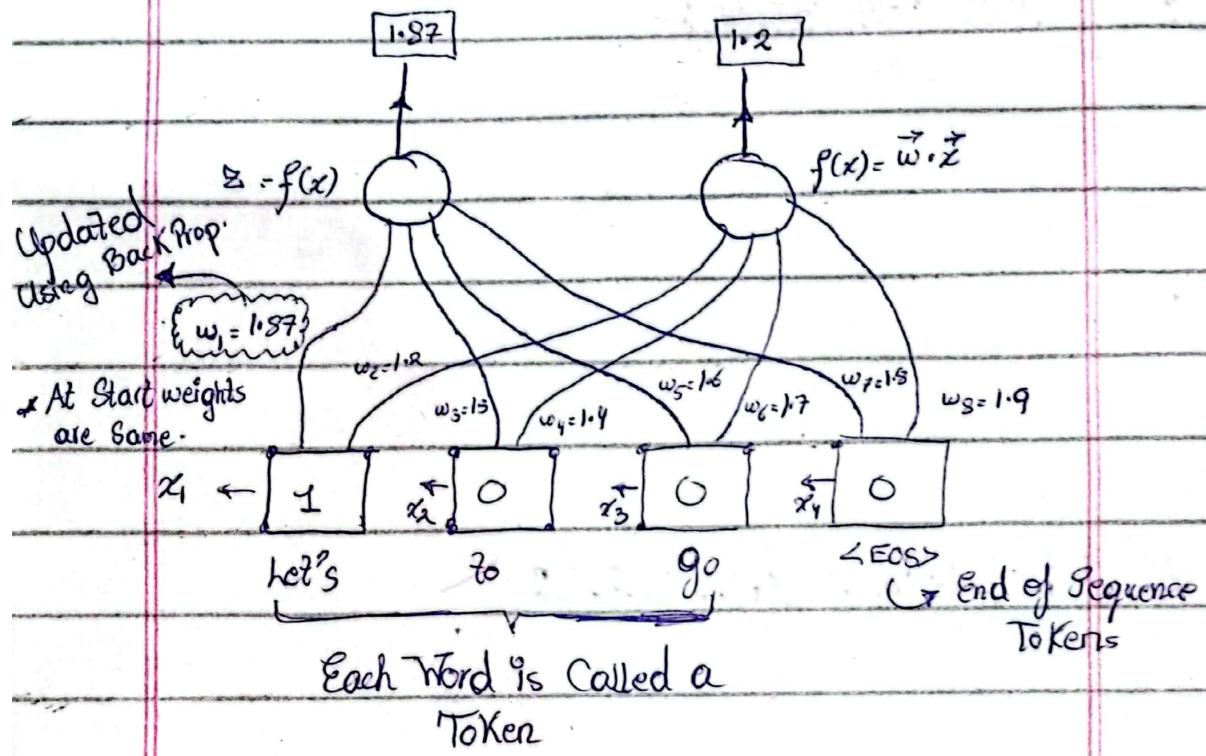
→ Embedding Space of GPT-3: 12,288

↳ Dictionary Size: 50,000 Tokens

→ Context size: 2,018 Tokens

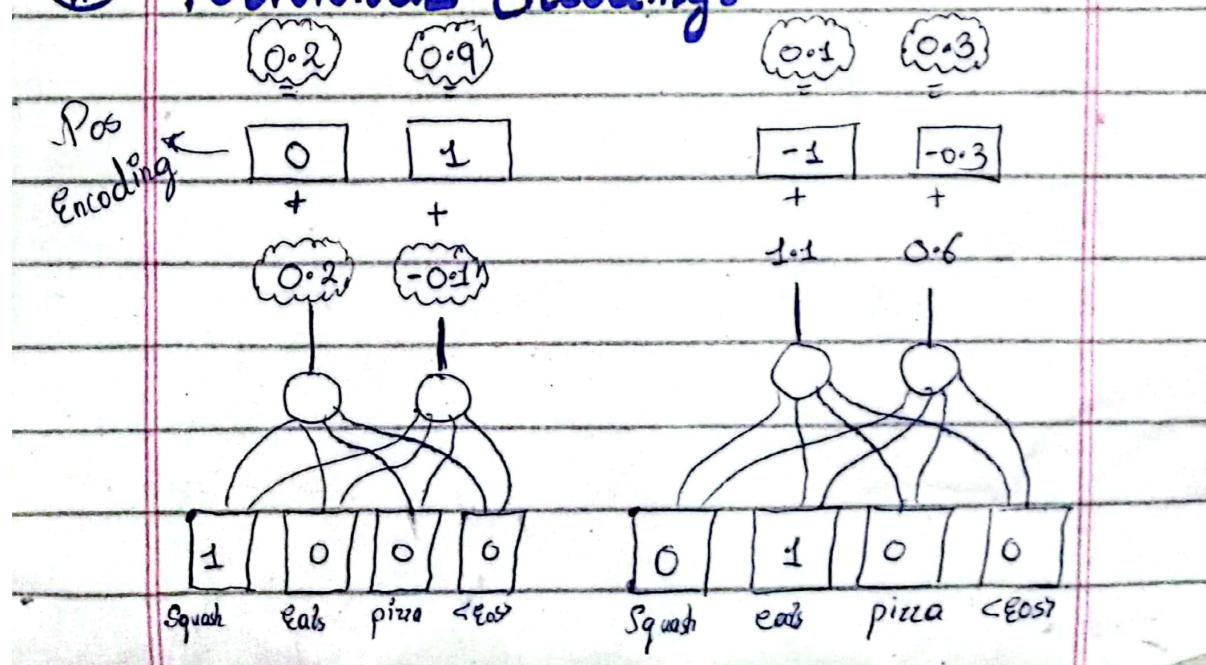
① Word Embeddings -

→ Select words into numbers using its vocabulary created using word2vec:

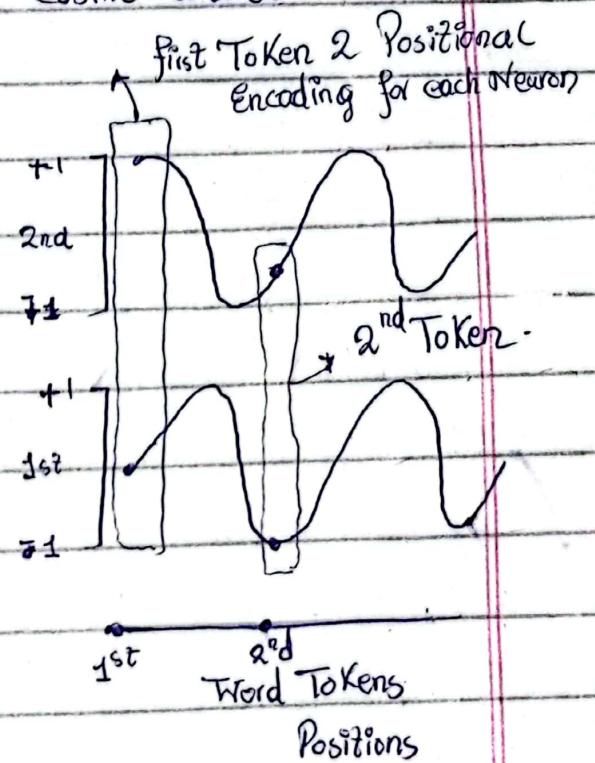


② Positional Encoding -

Positional Encoding -



↳ we get Positional Encoding using alternating Sine and Cosine Waves.



↳ When we Swap word's in Sentences, these word Embedding swap but the positional encoding ^{values} remains the same, so we end up with new pos encodings

III

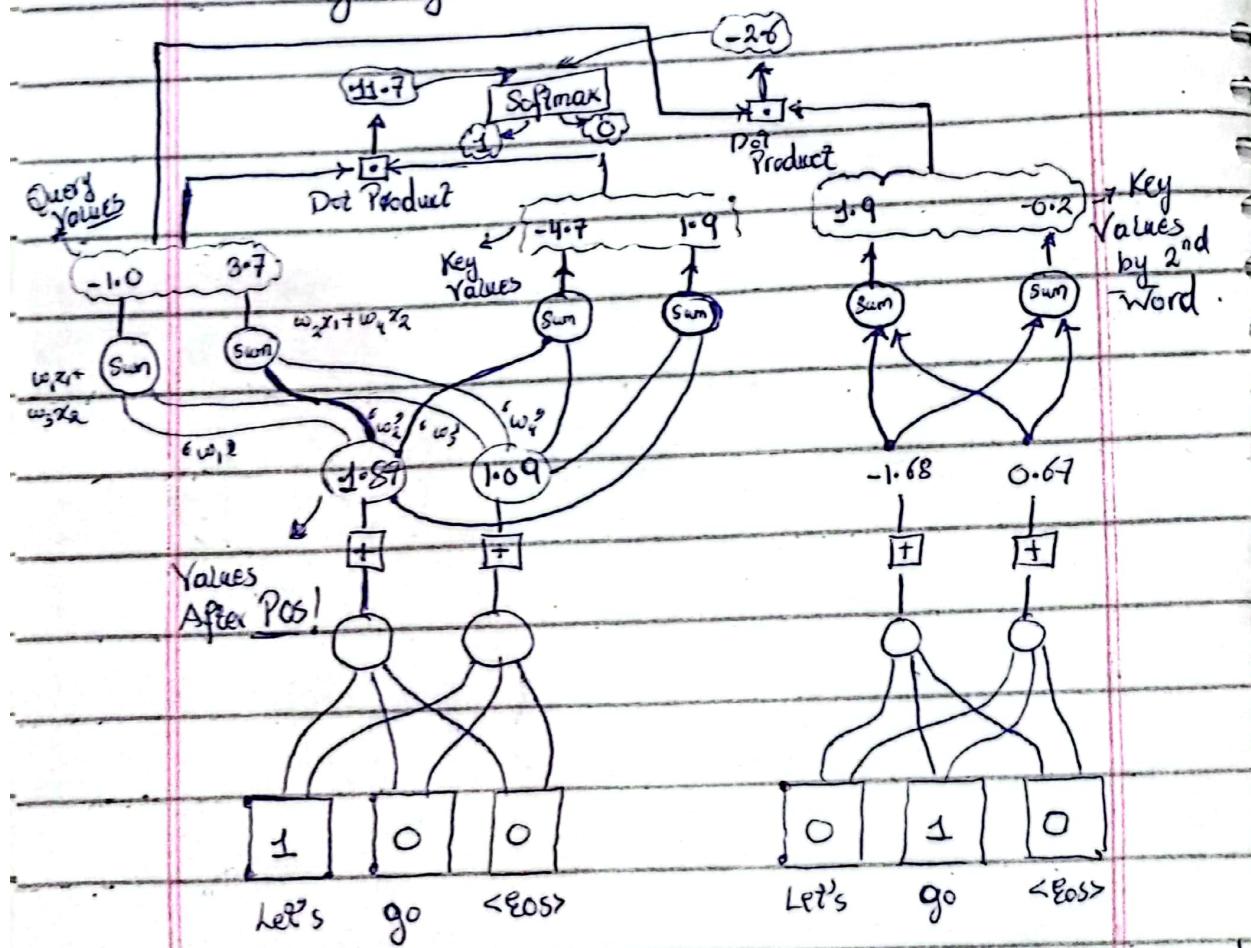
Self-Attention-

↳ Method to Associate words with each other. It gives context that

$$\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 2+2 \\ 1+2 \end{bmatrix} = \begin{bmatrix} 4 \\ 3 \end{bmatrix}$$

Ex 2

how much attention should a word
be giving to other word in sentences



→ You multiply Attention Scores with Every word Value vector

* Multiply Softmax Values by Values matrix to generate Vector in Vector Space

After Value Vectors Generated add

them together like $\begin{bmatrix} 2.5 \\ -2.1 \end{bmatrix} \rightarrow$ This

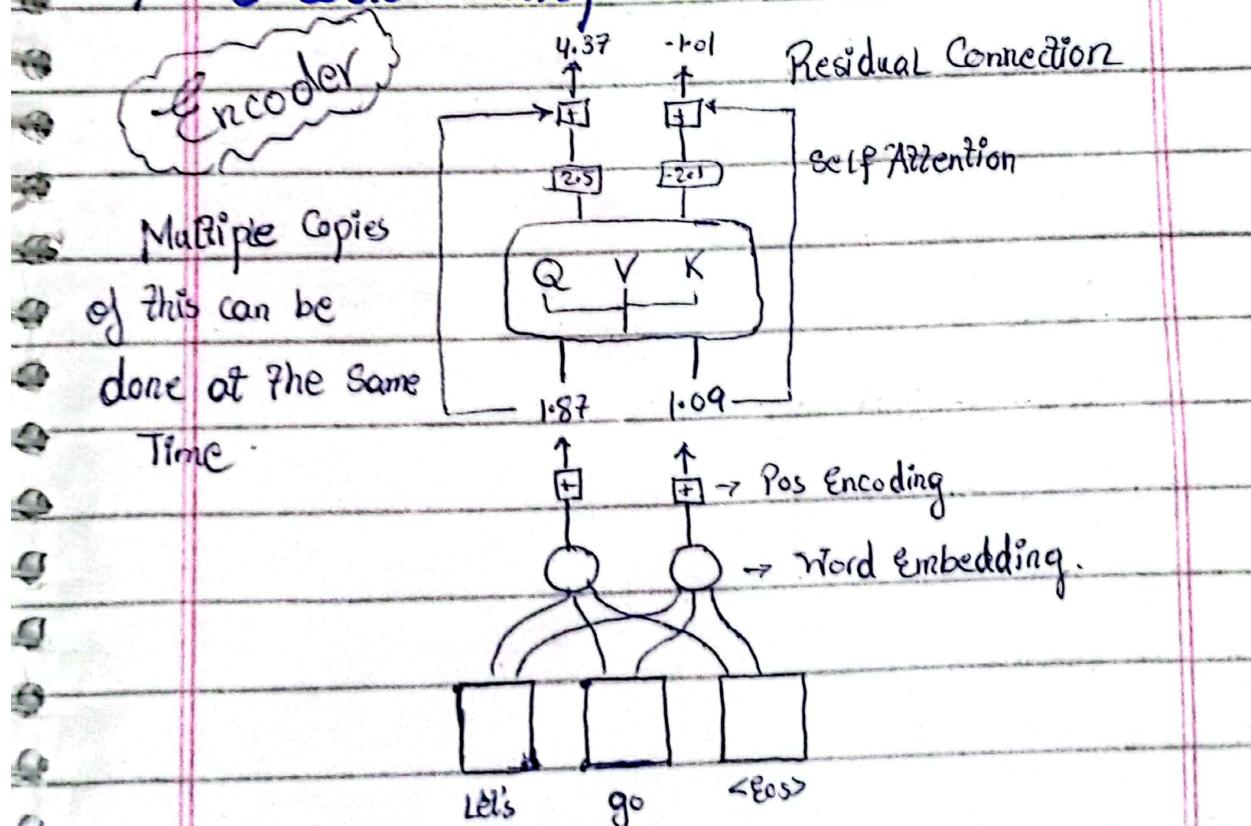
Shows Self-Attention Values of word 'Let's'.

This can Run Parallel

↑

- * Note the Weights Values used for Queries, Keys and Values remains the same for all the words in the input of Transformer
- * For Multi-head, each uses its own set of weights for Query, Key and Values

iv) Encoder Transformer 8-



① Decoder Transformer -

→ After encoding values are generated they are passed to Encoder-Decoder Attention Part.

* Note it has its own pair of Query,

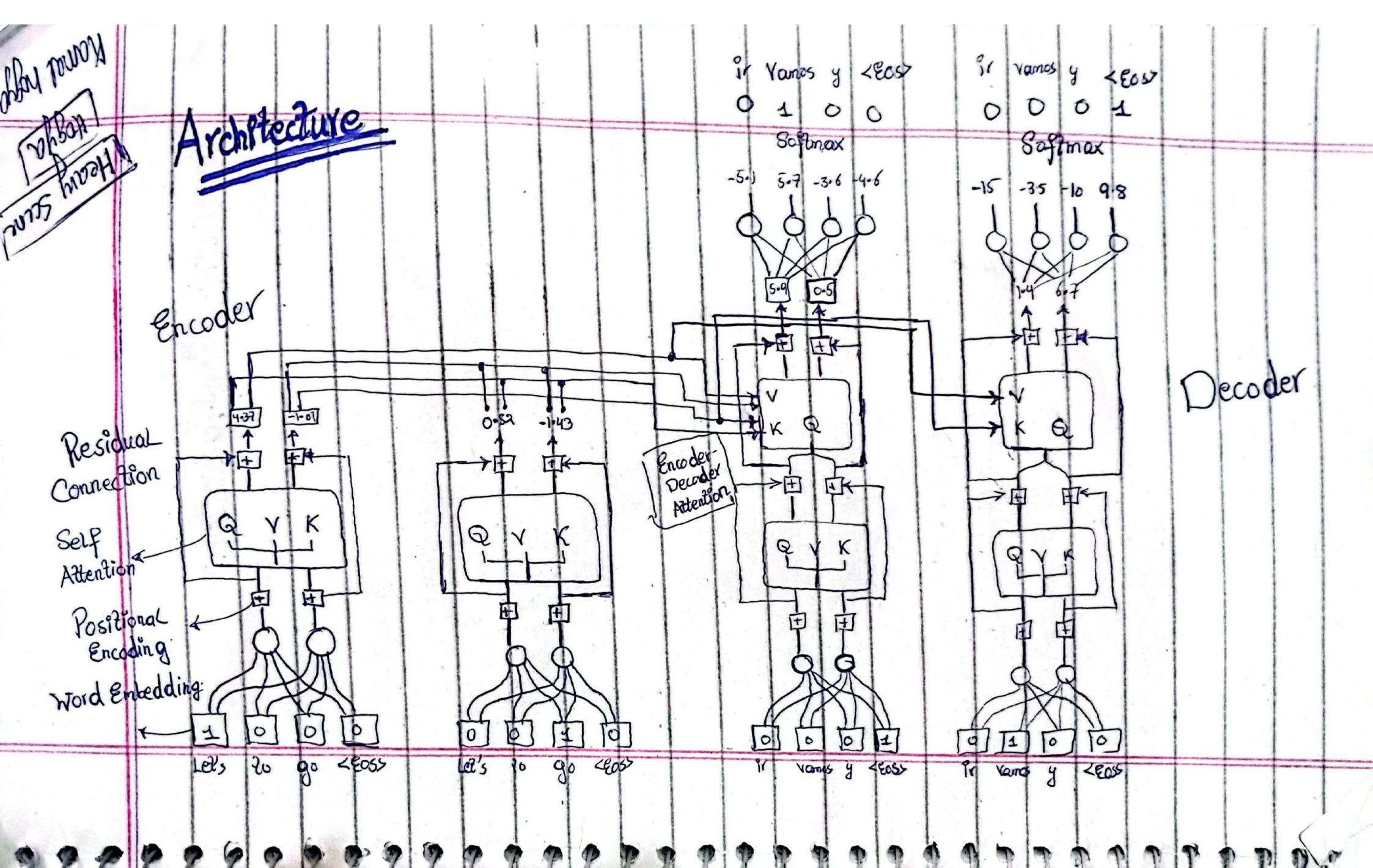
Key and value weights. After Attention

Apply Residual Connection with Output

embeddings, to generate a vector which

is passed to fully Connected layer having

no. of Outputs equal to decoded vocabulary.



* Generative Configurations -

① Max New Tokens

→ Greedy Sampling take highest Probability word [for Short Text Generation].

↳ Random [-weighted] Sampling :- select a token using random-weighted strategy.

② Top K :-

↳ Select Output from Top - K results and apply random-weighted Strategy.

③ Top - P :- Top - K

↳ Combined words probability should be Less than or equal to p .

④ Temperatures -

Broader, flatter Probability Distribution.

↳ high Value → Increase Randomness.

↳ low Value → Decrease Randomness

* Prompt-Engineering 8-

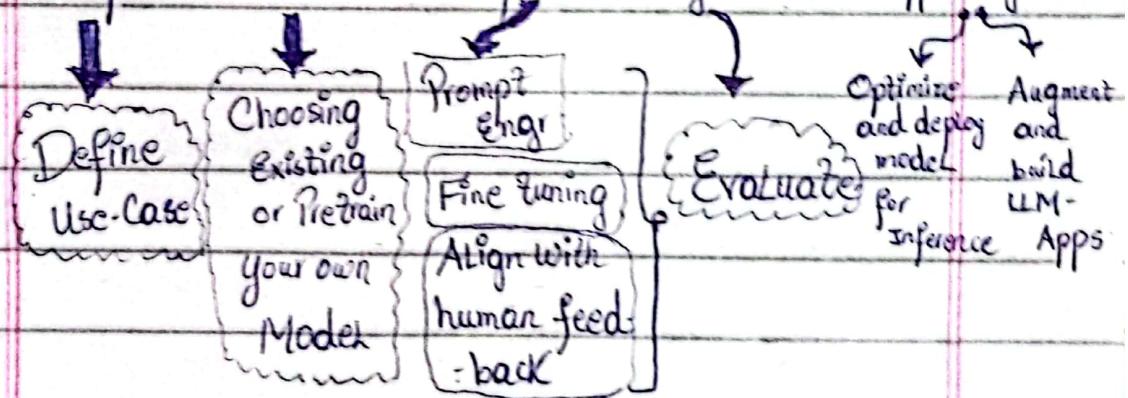
- In-Context Learning
- ① Zero-Shot Prompting → No context.
 - ② One-Shot Prompting → Single Example
 - ③ few-Shot Prompting → two or more examples in prompt

→ In-Context learning has limitation of Context window size and for too small model low performance

* Generative Configurations-

↳ Generative Life Cycle

Scope → Select → Adapt and Align Model → App Integration -



↳ LLM-Scope and Use-Case -

- ① Essay Writing
- ② Summarization
- ③ Translation
- ④ Information Retrieval
- ⑤ Invoke API and Actions

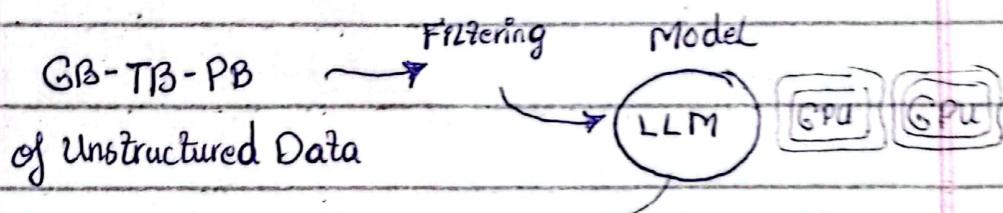
* Pre-training Large Language Models-

↳ Foundational Models

↳ Custom Models (Train your own Model)

* Model hubs { Model usecases & limitation }

↳ Pre-training At high level -



<u>Token</u>	<u>TokenID</u>	<u>Embedding</u>
--------------	----------------	------------------

'the'	37	[-0.05, 2.37, ...]
-------	----	--------------------

'-teacher'	3145	[2.5, 3.2, ...]
------------	------	-----------------

'-teaches'	11749	[3.5, 4.7, ...]
------------	-------	-----------------

'-the'	8	[2.2, 5.2, 2.8, ...]
--------	---	----------------------

① Encoder-Only Models - [Auto-encoding Models]

↳ Masking Language Models

↳ Each word has full context of the words before and after it.

Use Cases :-

- (i) Sentiment Analysis
- (ii) Named Entity Recognition
- (iii) Word Classification

Example Models :- (i) BERT (ii) ROBERTA

(ii) Decoder-Only Models - [Auto-Regressive Models]

- ↳ For Text-Generation
- ↳ Each Token Sees Token Before it

(iii) Encoder-Decoder Models: { Seq-to-Seq }

- ↳ Translation
- ↳ Text Summarization
- ↳ Question Answering.



Computational Challenges :-

↳ GPU RAM for 1B Parameters.

∴ 1 Param = 4 bytes (32-bit float)

1B Param = $4 \times 1B = 4 \text{ GB}$ { GPU Memory }

Not Only

This

Adam Optimizer :- 8 bytes Per Param

Gradients: +4 bytes per Param

Activations: +8 bytes per Param

∴ Total for a Single = 4 bytes + 20 extra bytes
Parameter

This makes it 24 GB - RAM Required

∴ To Reduce This Use Quantization Process

↳ Reduction of 32-bit floating Point to 16-bit FP or 8-bit Integer, reduces precision but are memory efficient

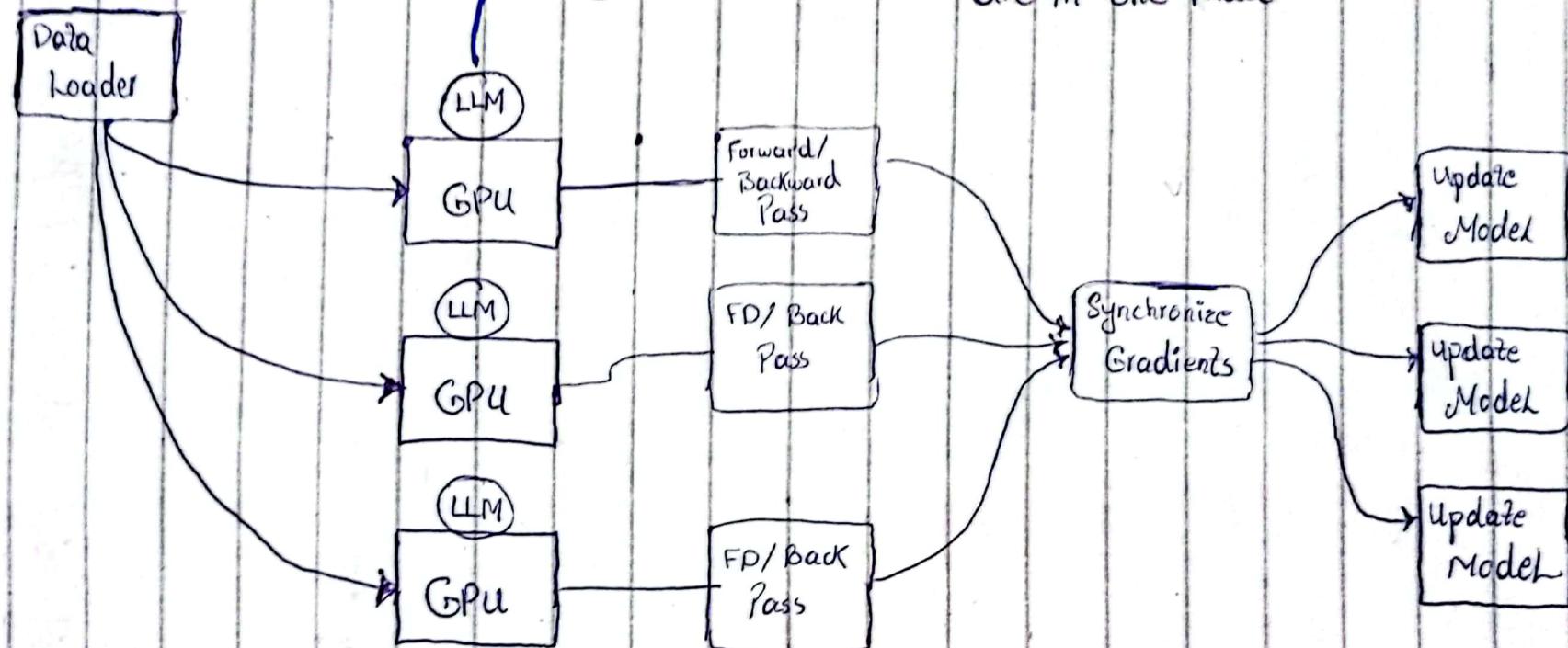
↳ Another Alternative is to Use BFLOAT-16. Which Optimizes the 16 bit - Precision.

↳ Quantization Aware - Training learns the scaling factors during Training.

* Efficient GPU - Strategies:-

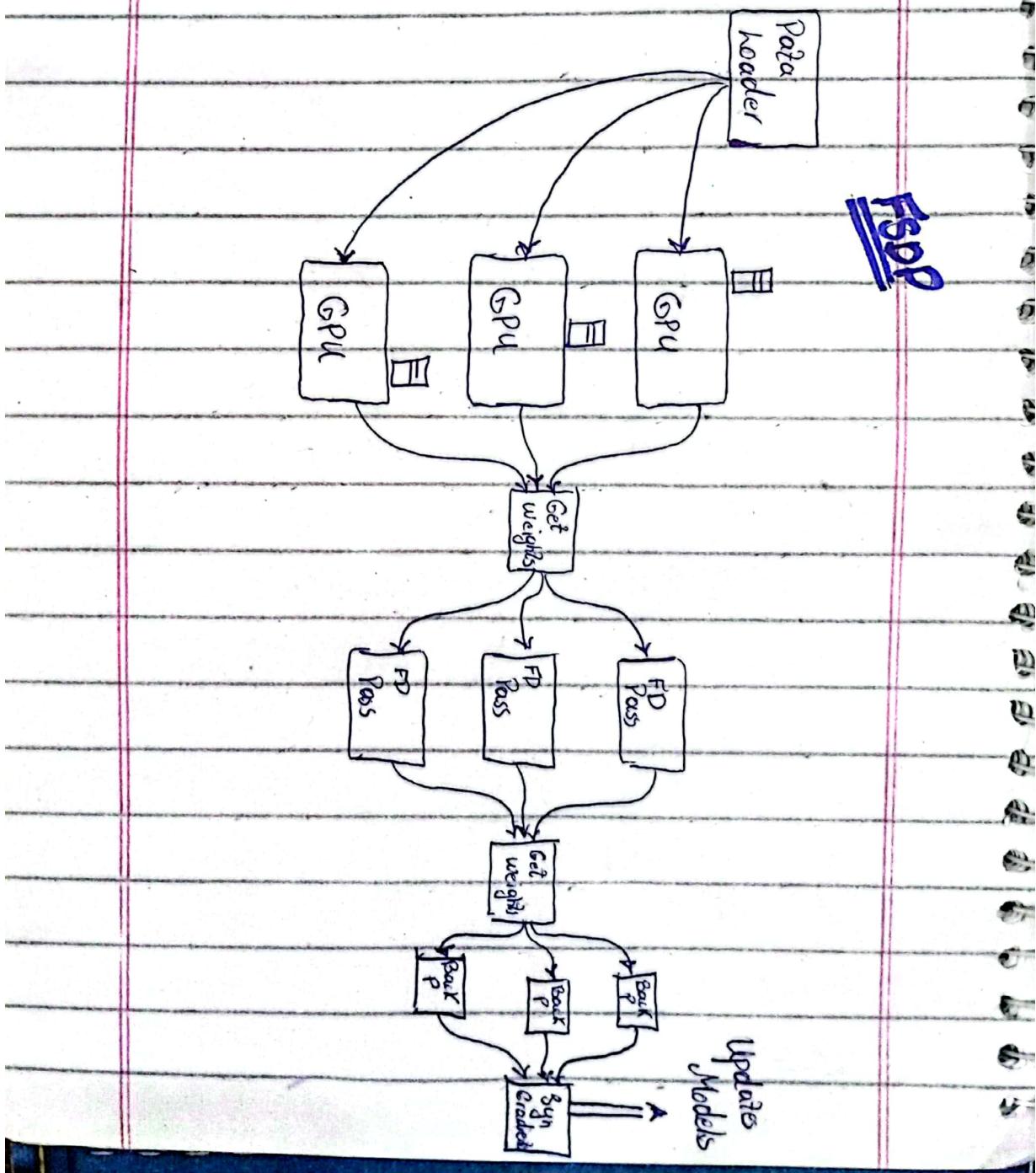
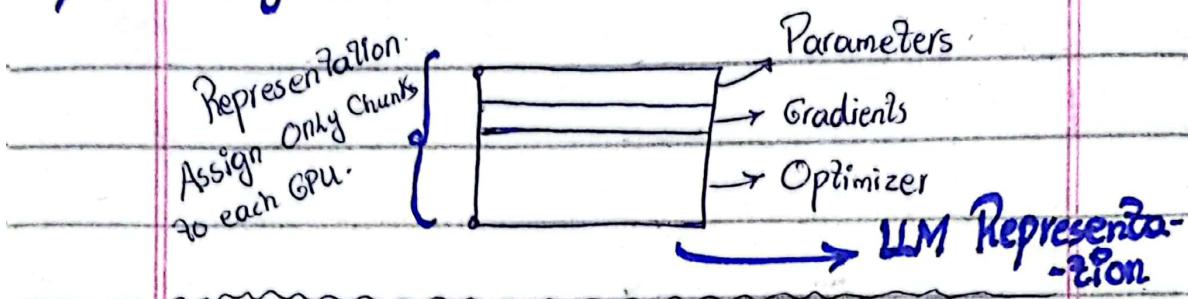
① Distributed Data Parallelism -

DDP 8-



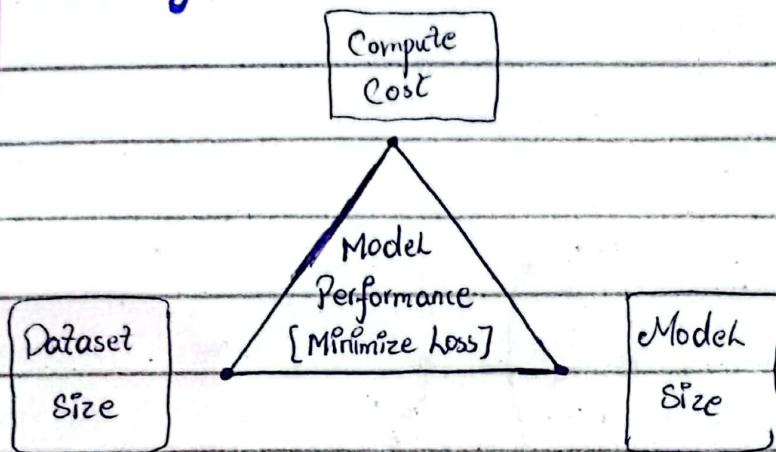


Fully Sharded Data Parallelism [FSDP] :-



1 PetaFlop/s-day → 1 PetaFlop Operations for whole day.
1 Quadrillion floating Point Operations Per Second.

Scaling Laws 8-



Increasing Dataset Size on a fixed

Smaller model can also increase performance and reduce loss.

↳ Finding the Optimal No. of Parameters and Dataset size for max performance.

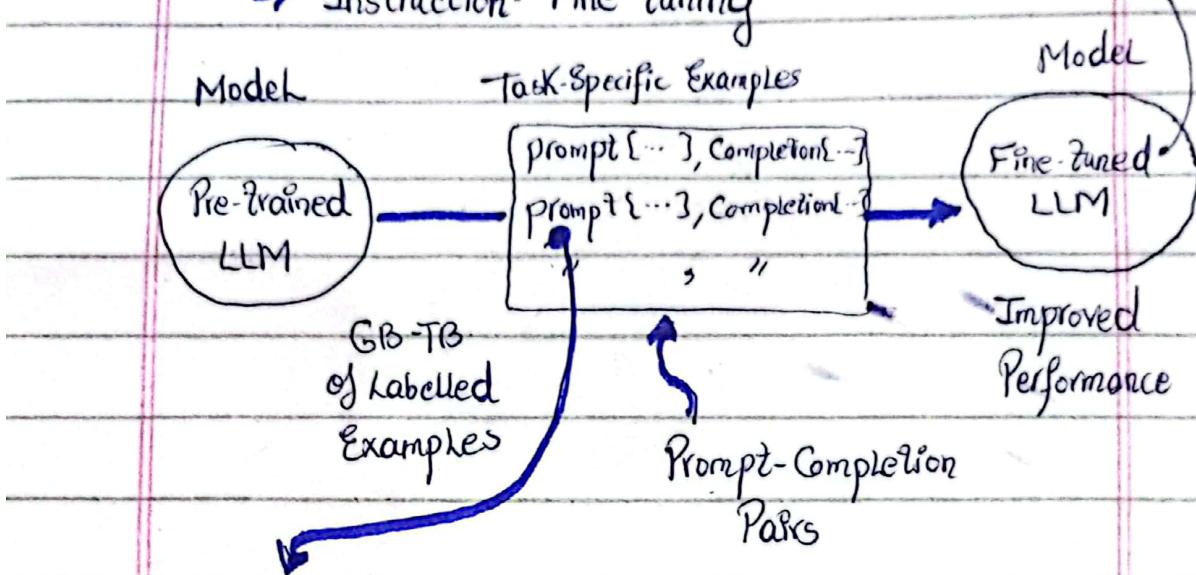
↳ GPT-3 Under Achieve during Training because it has 175B Parameters yet dataset size was only 500B Tokens.

↳ Llama uses 6.8B Parameters but dataset size is 20x times, so 1.4 Trillon Tokens. Utilizing All Parameters during Training.

* Fine-Tuning 8-

↳ Instruction-Fine Tuning

Also Known as
Instruction-Model is
case of Instruction fine-
tune



Each Prompt Consist of Instruction associated
with the text :-

- Like:-
- i) Do Sentiment Analysis / Classify the Review
- ii) Predict the Rating
- iii) Generate a text
- iv) Provide small short summary of sentence.

↳ Split dataset in three Portions :- train, validation and test set.

↳ Full Fine-Tuning deals with using all parameters, gradients and Optimizers while training.