

DEK 1

Language model, probabilistic model that assign prob to sequence of words. $P(\text{this}/\text{context})$.

Base models can be used using finetuning / prompting -

Fine-tuning LLM. (seq-seq model)

- input \rightarrow [SOS] / sentence(i_0) $\xrightarrow{\text{T. Encoder}}$ output sentence(i_N)

No of input tokens = No of output tokens.

- Then we calculate cross entropy loss b/w output/target word & we use backpropagation to adjust.

= Inference.

prompt \rightarrow forcing tokens \rightarrow [SOS] $\xrightarrow{\text{PTM}}$ output token

Embedding

tok \rightarrow tok \rightarrow $\boxed{\quad}$ + $\boxed{\quad}$
 tok ID embedding PE

\rightarrow Embedding (words represented in high dimns by vector)
 to capture more

e.g. for 2-Dmer

\rightarrow the angle b/w similar word is small and opposite is big

$$\text{PE}_2(\text{pos}, 2i) \approx \frac{\text{pos}}{\text{dim}} \\ \text{PE}(1 \rightarrow 2i+1) \approx \cos \left(\frac{\text{pos}}{1000} \frac{2i}{\text{dim}} \right)$$

BERT - overview.

BERT - Base

- 12 encoder layers
- 12 attention heads
- 3072 feedforward size
- 768 d_model

BERT - large

- 24 encoder layers
- 16 attention heads
- 4096 feed ...
- 1024 d_model

- linear head changes according to application
 - PE(⁷⁶⁸) are absolute and learnt along training & 512 pos.
- word piece tokenizer (size ≈ 30,000 tokens)
- max len (512 tokens)

BERT vs GPT/LM

- not for chat, but for fine-tuning
- not like other, trained on left-right context.
- not specifically for text generation.
- = not trained on next token prediction, but masked in & sentence prediction task.

RAG (Retrieval-Augmented)

very useful for Q/A

used in GPT.

can access latest data.

- LLM assigned probab to sequence of words

- Inference, create a prompt to do so.

Fine-tuning Cons.

→ expensive

→ no of parameters might not be enough to capture all the knowledge we want to teach.

→ Fine-tuning is non-additive. E.g. (heavily finetuned might forget some previous knowledge)

few-shot (prompting) to rescue. → LLM can be taught to perform new task using prompt engineering

Q/A with prompt Engineering - {instruct} {context} {Q} {A}

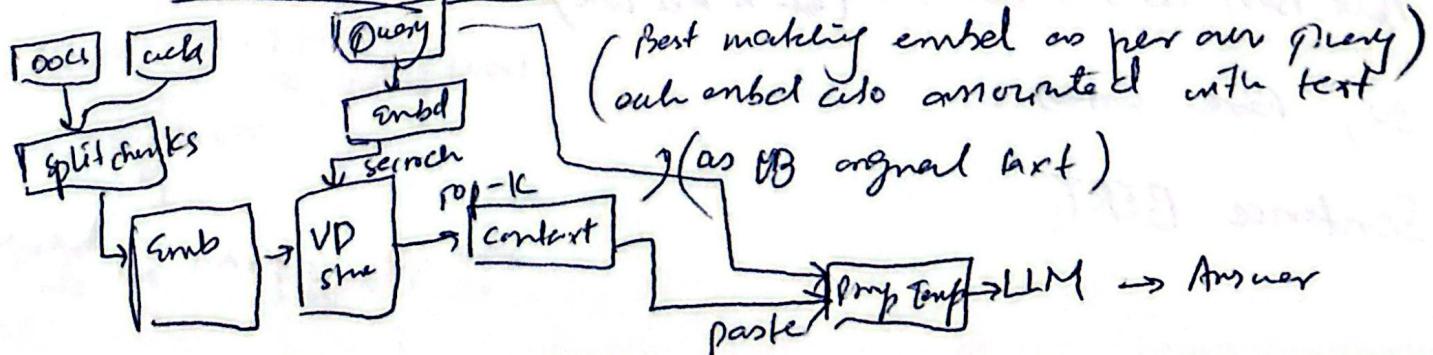
Answer, using access to the context on which it's not trained on

Pns of Finetuning.

→ higher quality result as compared to prompt

→ smaller context size (input) during inference (do need of context)

Q/A with RAG:



Embedding

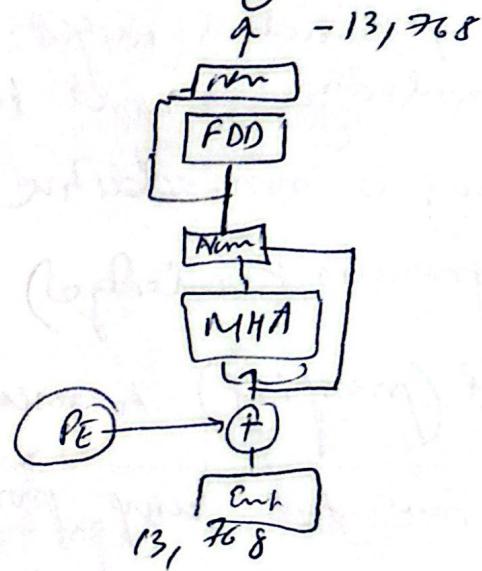
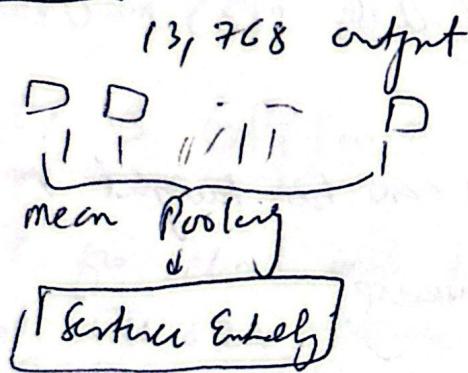
→ similarly by cosine.

words that are synonyms tend to occur in some extend (distributional hypothesis). To capture meaning we need to have access to context.

Self-attention access 2 ways word emb and position encoding and modify embeddings to capture context as well

→ In BERT pre-training input embeddings are also modified.

Sentence Embedding with BERT



Sentence Embed comparison

By cosine similarity → cosine of angle b/w 2-vects.

$$\text{Sc}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

$$\frac{\sum_{i=1}^n A_i \cdot B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

2 similar
sentence
should produce
same embedding

But BERT doesn't know this (how to tell them)

so, here comes

Sentence BERT

Art-Practicing.

Masked Language Model task (cloze)

2 Randomly select words in a sentence and mark
it and let model predict word in left/right context.
left/right context in part

→ we don't mask next value after ($\Phi \times R^t$), so each token attend token to its left and right.

- pretraining procedure selects (15%) of tokens from a sentence
 - 80% by smart
 - 10% by zebra
 - 10% by avg (saver and)

Next-Predictor Task (NPT)

Input Emb \rightarrow + Segment Emb \rightarrow + PE

↑ Represent this token is Sent A Nlu or B

□ □ - □ □ □

out [st] []

\downarrow
linear layer (2 features)
 \downarrow softmax
 \downarrow next { ! max

\rightarrow C.E loss \uparrow backgnd noise

[CLS] token in BERT: capture info from all other tokens
because not masked in (Φ^{k_i}).

BERT-Finetuning

(prompting technique)

one prompt, class of noun, or

→ Text-classification

[CLS]

→ $\Phi(A)$:

2 issues

→ which part of input in context, which is question?
Ans: segment analysis $\rightarrow Q \Rightarrow A, C \Rightarrow B$

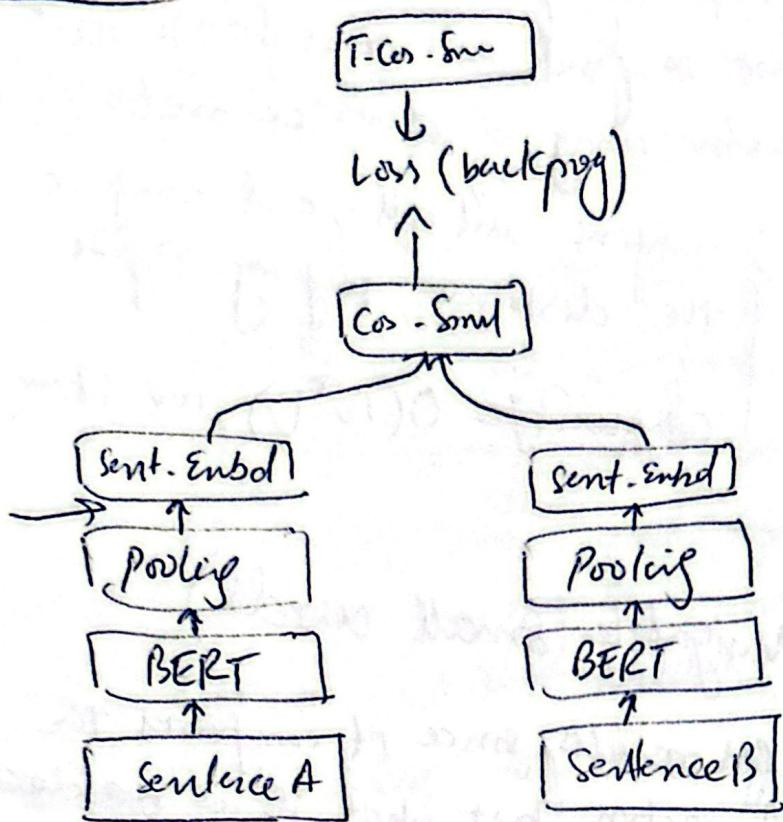
→ Start & end of answer

[CLS] \rightarrow (start token or end token) \rightarrow we know [zero is answer]

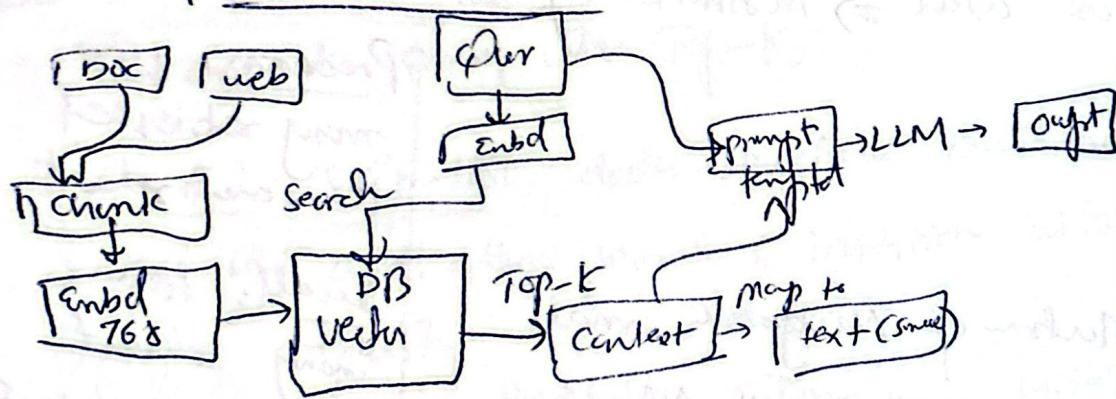
Sentence in BERT

Architecture

lower layer
if you want
to reduce size



Q/A with RAG:



- Strategies
- Fneture
 - RAG
 - Feature + RAG

gradient platform-

Siamese Network

- same architecture
- same parameters
- same weights

Implementation

Just one network,
other is branch

Search -
cosine similarity

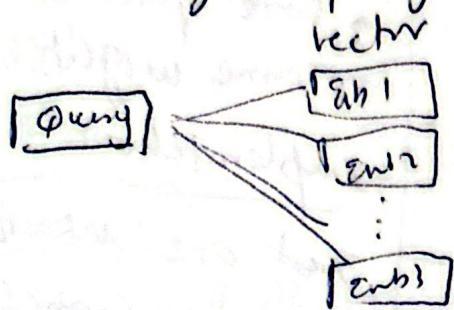
have high
values -

Euclidean
dist
< given -

vector DB

→ similarity algorithm → (KNN variant)

→ vector database stores vectors of fixed dimension such that we can query the database to find the embeddings that are closest to given query vector using a distance metric (cosine similarity).



→ compare with all, and compare & sort the distance, keeping Top-k.

$$\text{complexity} = O(N^2 D), \text{ too slow}$$

HNSW (Hierarchical Navigable Small World)

- always produce accurate results, since it compares the query with all stored vectors but what if we reduced the no of comparison but still obtain accurate results
- The metric we used ⇒ similarity search is called

Recall.

Precision & Recall:

HNSW is an evolution of Navigable small world algorithm for approximating NN Neighbors which is based on the concept of six degrees of separation

Precision: How many retrieved items are relevant

Recall: How many relevant items are retrieved

Six-degree of Perception.

Navigable-Small world → it builds a graph that connects close vectors with each other but keeping no of connections small (every vector may be connected upto 6 vectors)

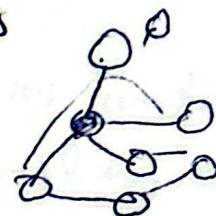
→ Randomly choose a point in graph & compare similarity score with query & we compare query similarity with the friends of entry point and move which has better similarity score. and (repeat). Stop at local best (best score) among its connection (neighbors). Order them best to least-best matches & keep best.

→ Repeat by doing different entry points

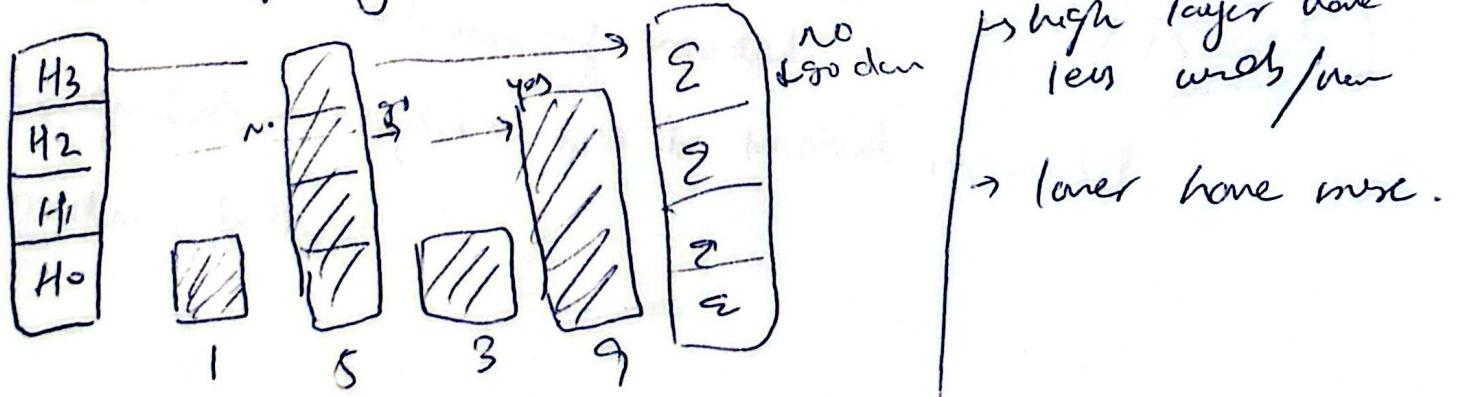
→ Keep Top-K

→ Inserting a new vector.

→ connect query with Top-K



Idea #2: skip-list. data-structure that maintain a sorted list and allow search & insertion with an average of $O(\log N)$ complexity.



HNSW: skip list + Navigah --. dense, sparse

let

P

- Find random entry point at top-level
- calculate Cosine-Similarity with it & his friends.
- local best → go down
- we repeat with other random starting
- we sort & keep Topk.

