

What	date was Iribe opened to public?
4	retrière Inde Wikipedia article/
·	UMD CS website,
4	Concatenate these does to my prompt
	6 doc 1> wikipedia 4/doc 1>
	Ldoc 27 UMD LS page 2/ Acc 23
6	feed x; Z k to Po to geneate y

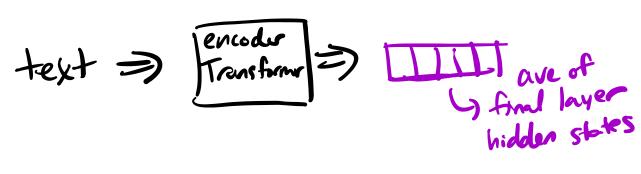
how do we perform retrieval?

5 measure similarity between prompt X
and docs Z

15 how?

15 string motiching 3 cheap, effective
15 BM-25 3 not good at
15 BM-25 3 measuring semantic similarity

La vector database la relies on a pretrained encoder model



b) embed prompt x

embed each doc Z; m datastore

b) identify k nearest neighbors

of x

(x.Z, = 5

products x.Zz = -8

(x.Zz = 0

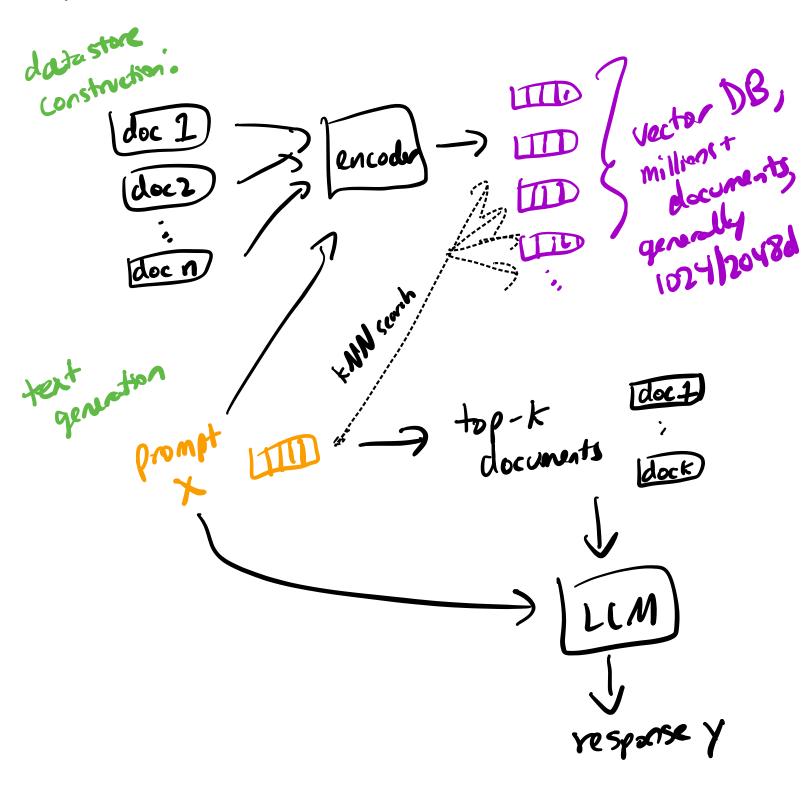
y.Zy = 29

La exact nearest reighbors is intractable

1) generally use approximate kNN LIFAISS

Dencoder is fine-timed on datasets of query/ context pairs

## LARAG WI rector DB



## RAG benefits:

L) specifie facts / domain knowledge may not be memorized by LLM during pretraining

L) fresh knowledge

b) reduce hallucinations, do attribution

L) connect generated text to Source docs

L) not always reliable, LLM can misinterpret retrieved does lignore

RAG cons:

1) retrieved is really hard!

1) poor quality retrieval often leads to worse perf than no retrieval

L) parametrie knowledge
L) knowledge encoded in the
models params

## Long parametric knowledge lo external clock Context

L) information loss when embedding docs

L) chunking, loses context the chunks of retrieved

L) increase inference cost as the increases therefore the chunks of the chunk

Tool use:

- to retrieved is an example of a tool that the LLM can use
- b) in RAF, retrieved happens prior to response generation
- b) what if the model could call diff tools during generation?

4) Toolformer

Agents

L) LLMs +tools + dynamic control flow L) some notion of state / memory