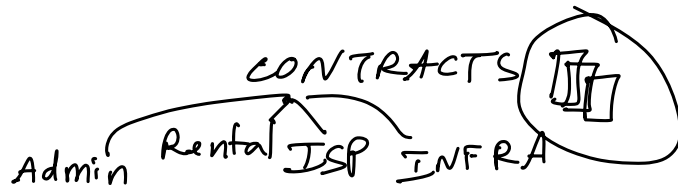


# RAG (Retrieval-Augmented Generation)



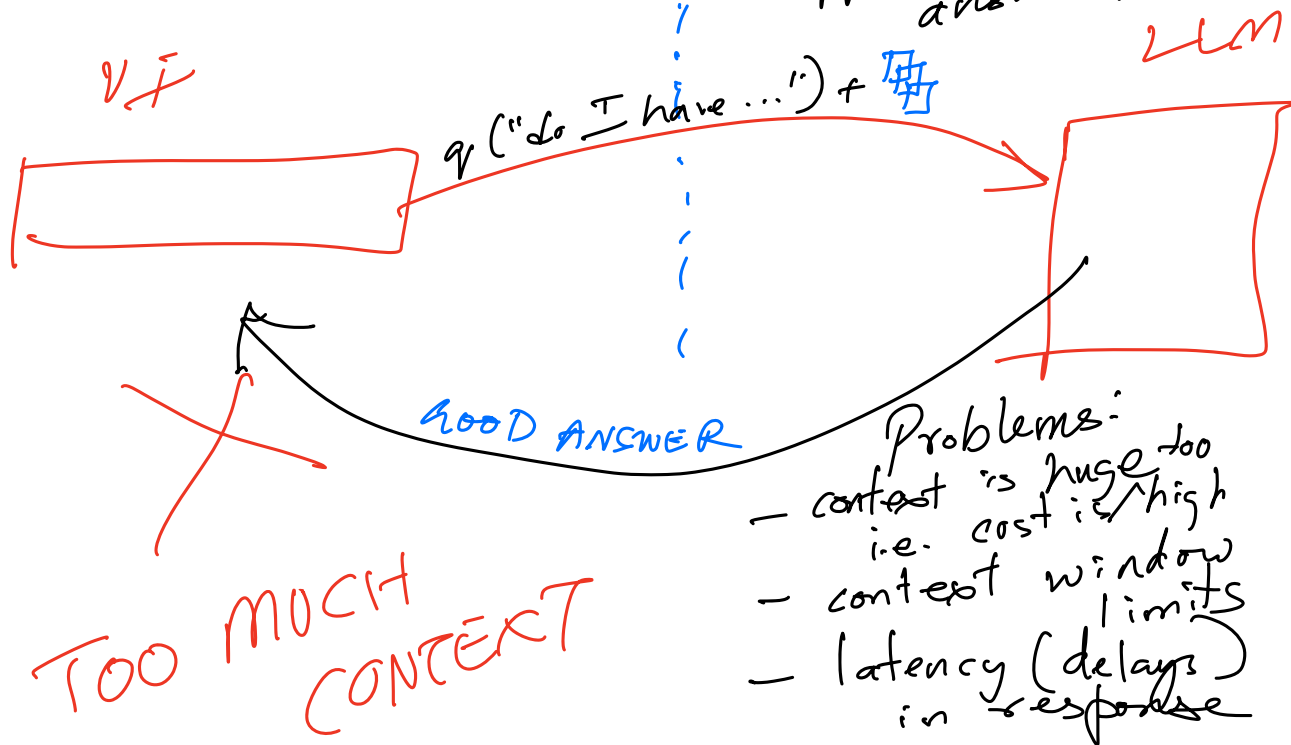
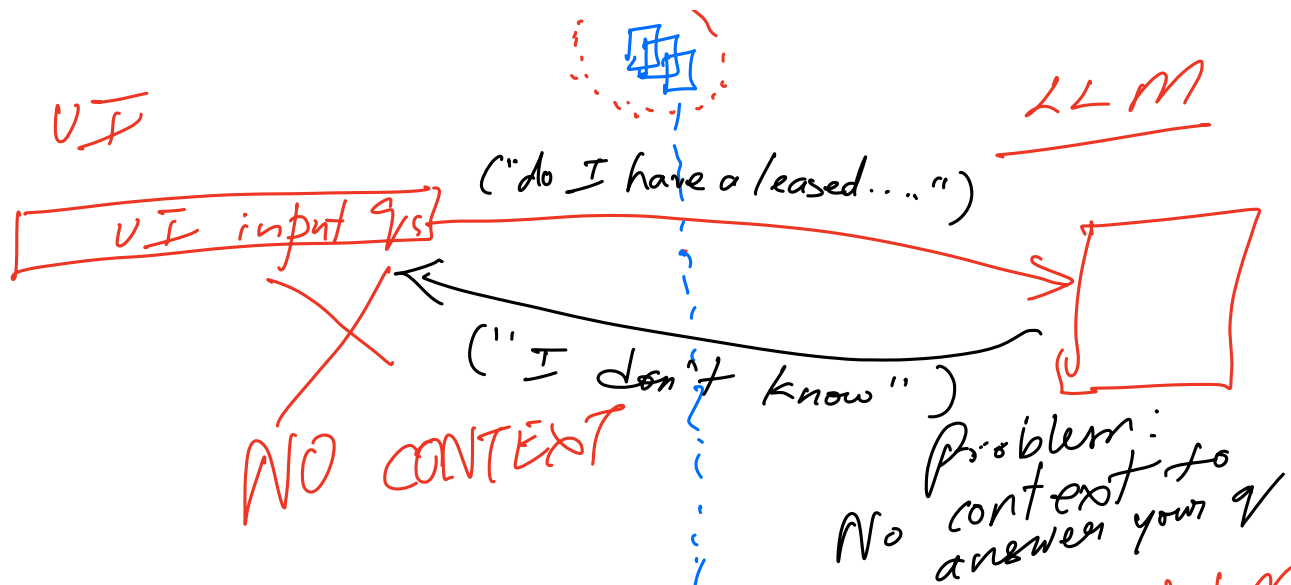
Implement a UI / chat iff that allows for questions to be asked about the contracts.

For e.g. q: when does my lease on the Road No. 10 property end?

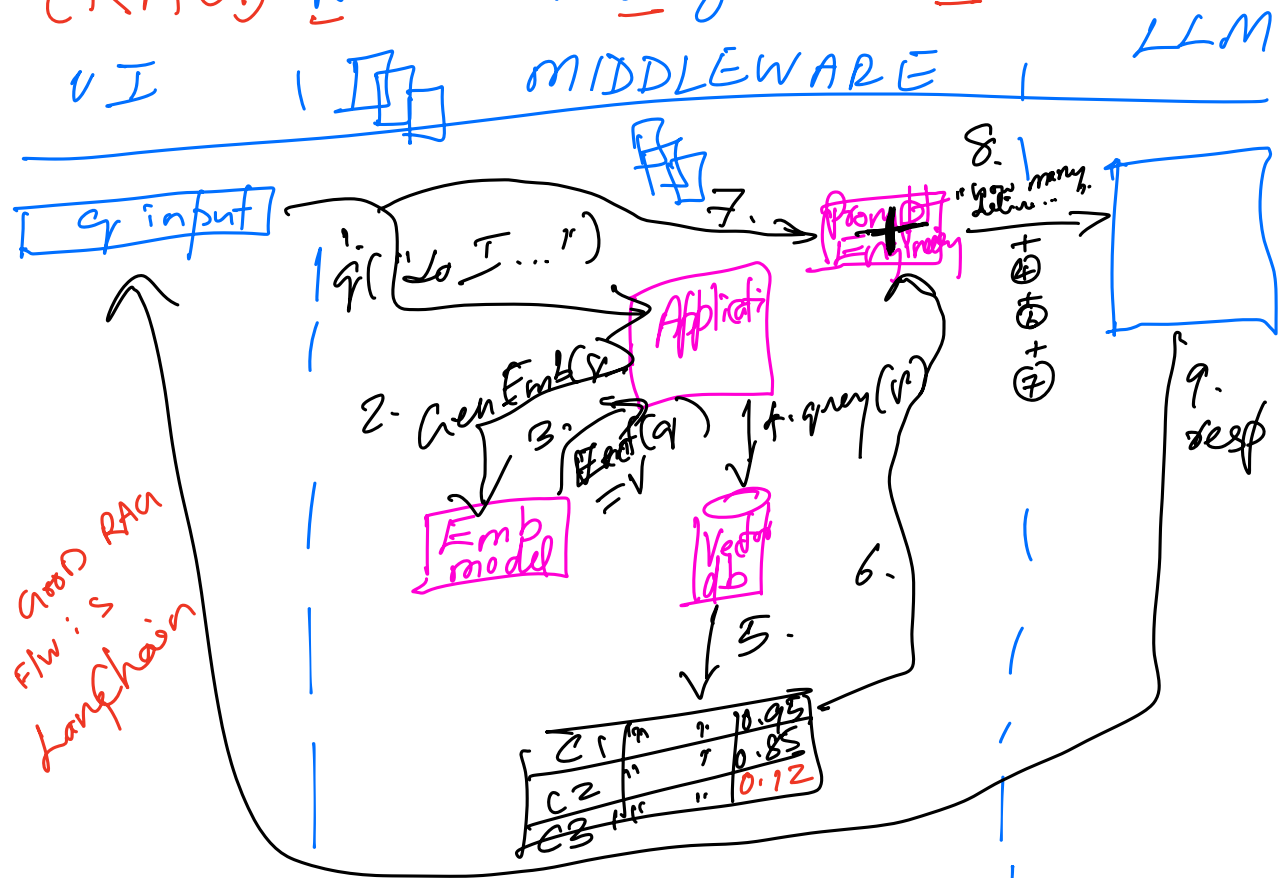
q: do I have a leased line in my ISO ISP service

q: what is the penalty for ending our catering contract?

q: summarize my maintenance contract of my sector-7 lease



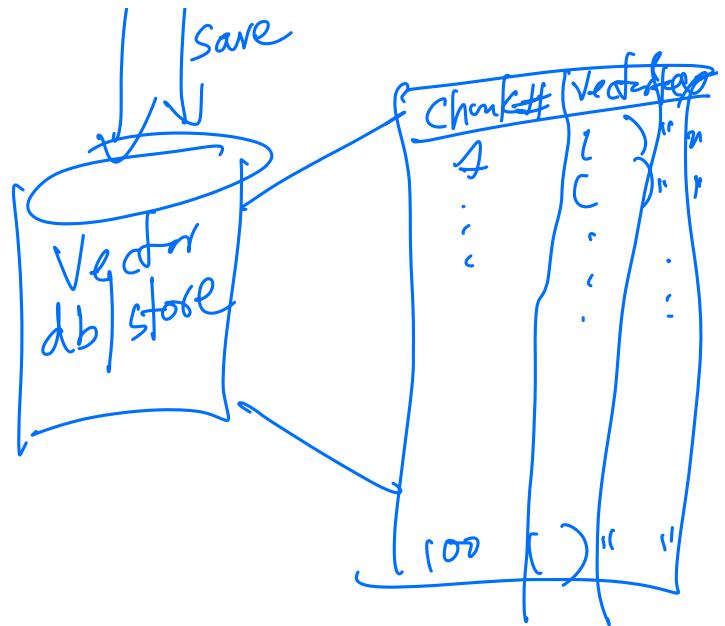
# (RAG) Retrieval-Augmented Generation



ONE-TIME PREP

(100 docs x 10,000 chars)

chunk it up into 1000 characters / chunk  
 (10 chunks x 1000 chars = 10000 chars)  
 vectorize (gen embedding)  
 Emb model



## CONTROLS WITH YOU

COSINE  
VS  
EUCLIDEAN

- PROMPT ENGINEERING
- CHUNK SIZE
- RELEVANCE SCORE FILTERING
- VECTOR DB SEARCH METHOD

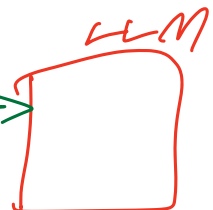
Prompting  
Retrieval

Q: are any of my contracts with a Maharashtra-based company?

RAG m/w

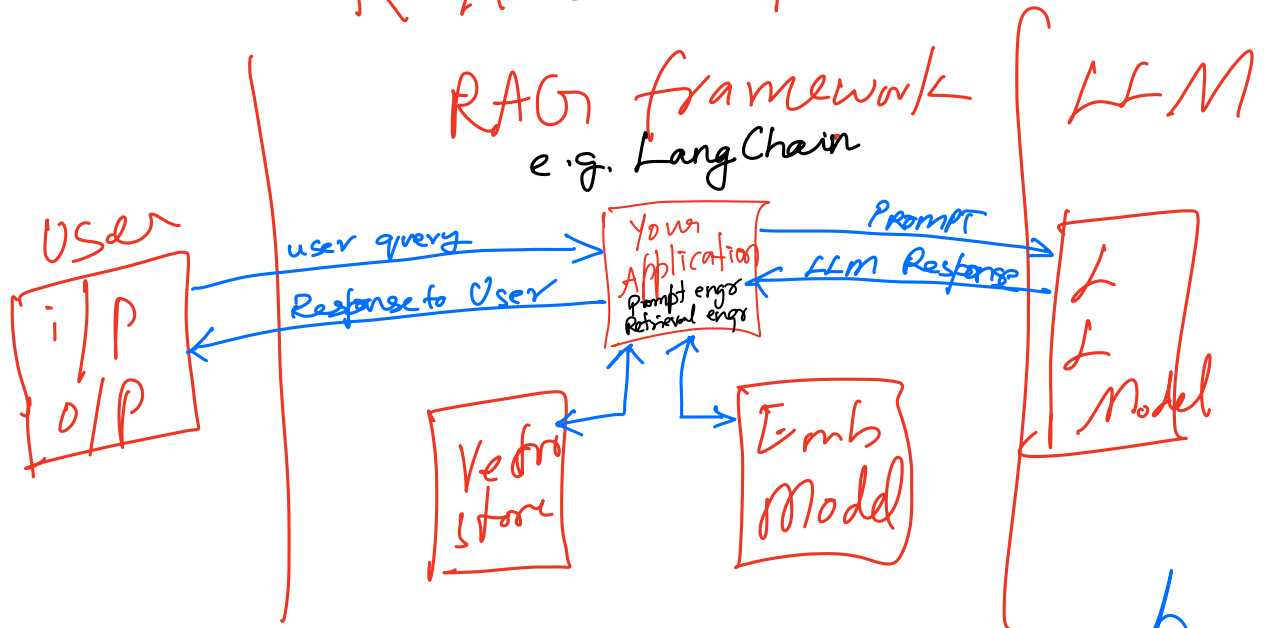
UI

"are any of my Maharashtra..."  
"use the below as context"  
"82, P. H. Road, ... New Delhi"  
"75, ... Mumbai"  
"Contract - ABC Holdings  
37, C. H. Road, Pune, ..."



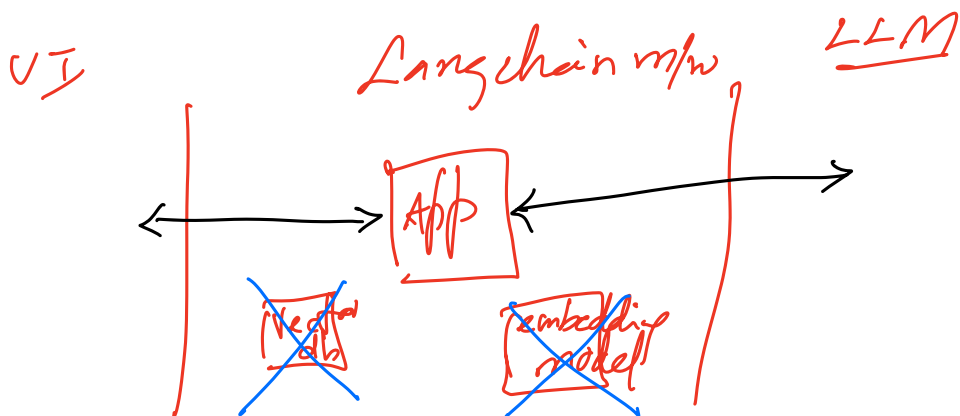


RAG m/ware

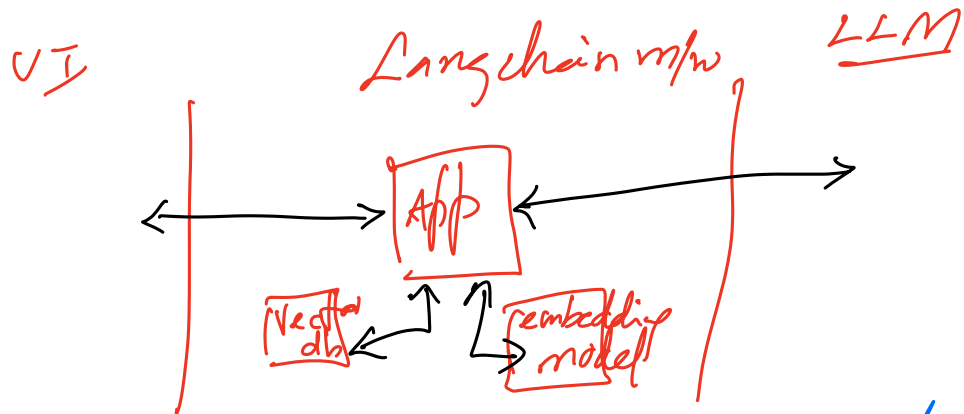


demos coming up

- 1) Embeddings demo (nothing to do with Langchain)
- 2) Langchain demo (direct with LLM)



### 3) Langchain RAG end-to-end demo



### 4) Lang Chain Advanced

