

## Retrieval-Augmented Generation

API

RAG vs Fine-tuning

Upfront latency  
Response relevance

Separate search

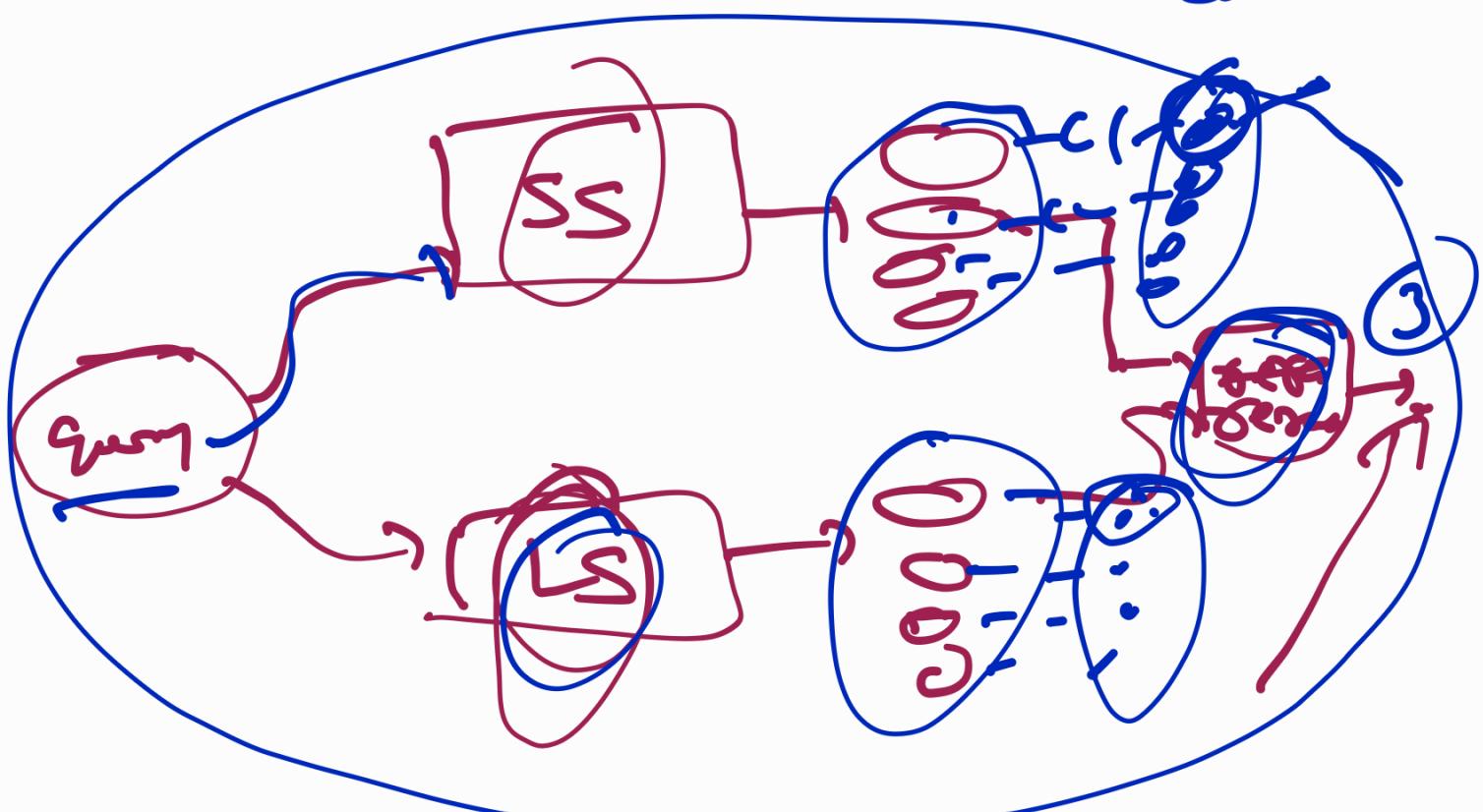
Process  
Recall  
Extract  
MRPC

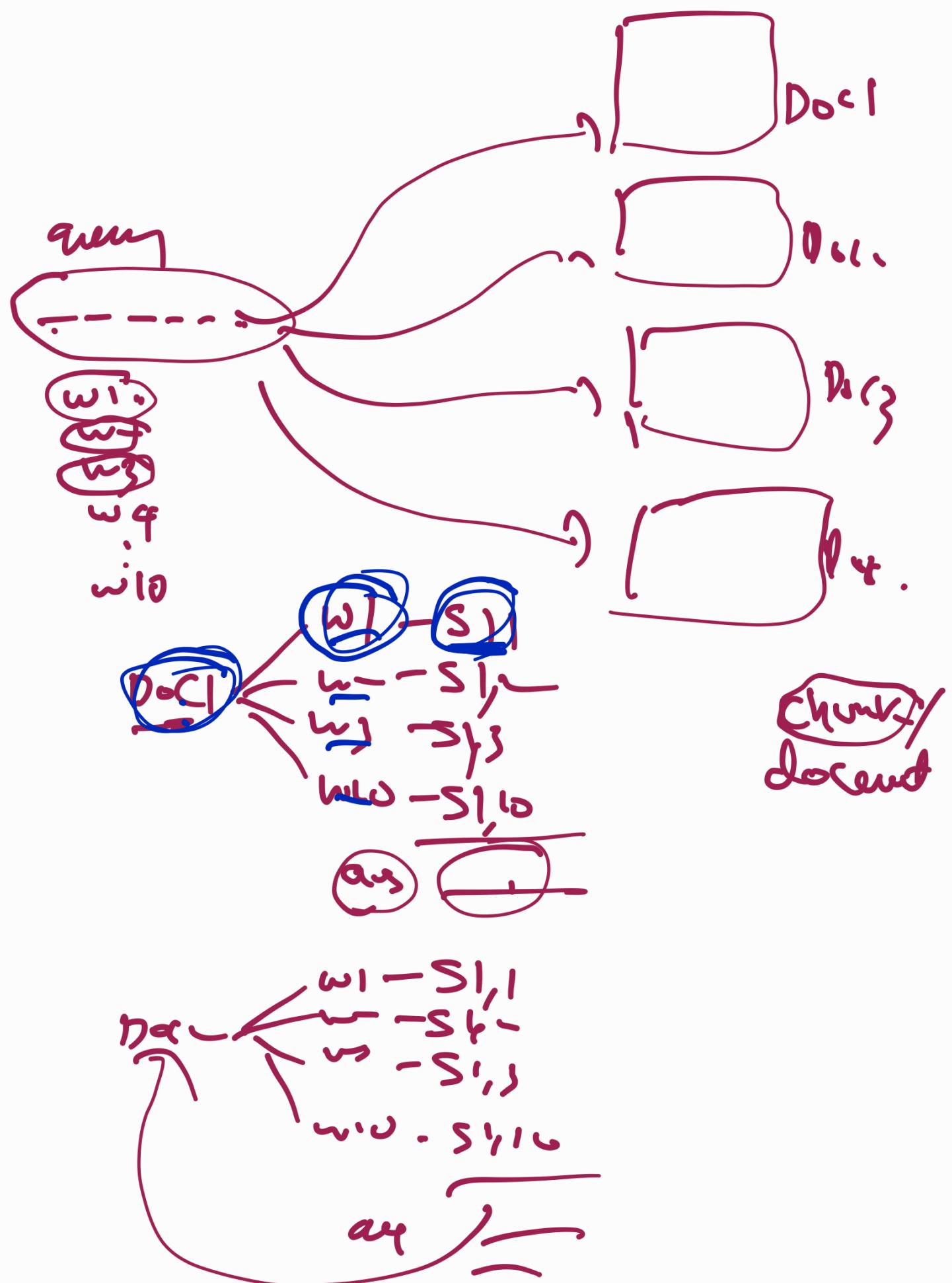
BM25  
Lexical search  
Keyword search

Text-to-Sel

Searchable

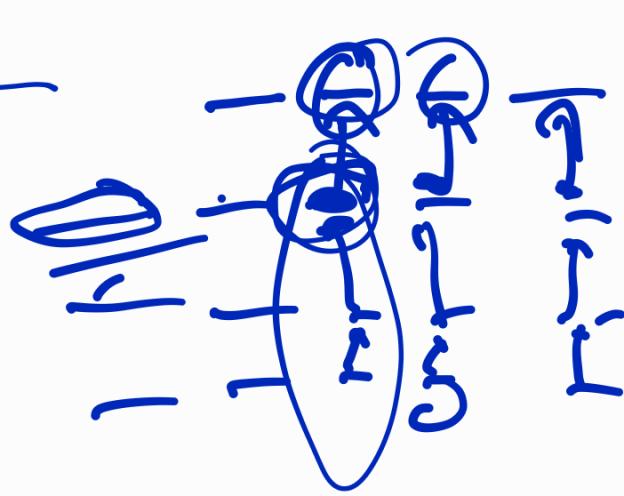
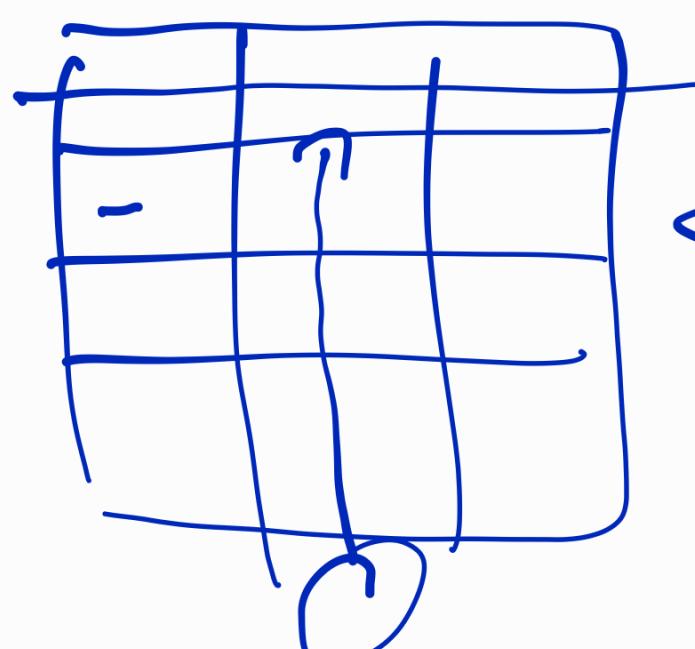
Latency

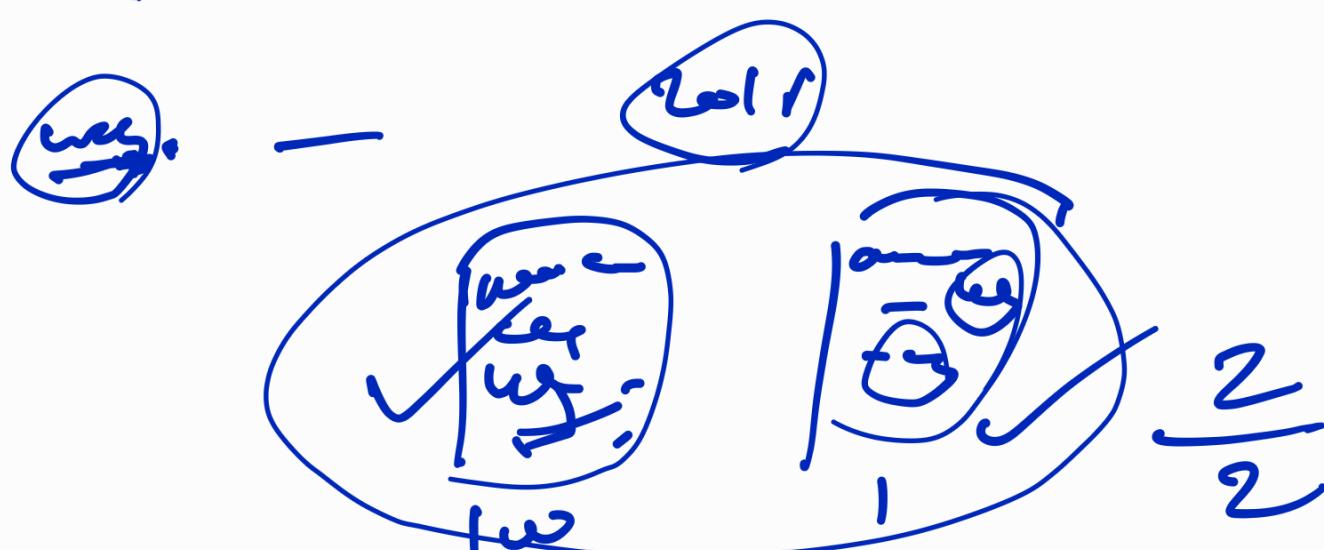
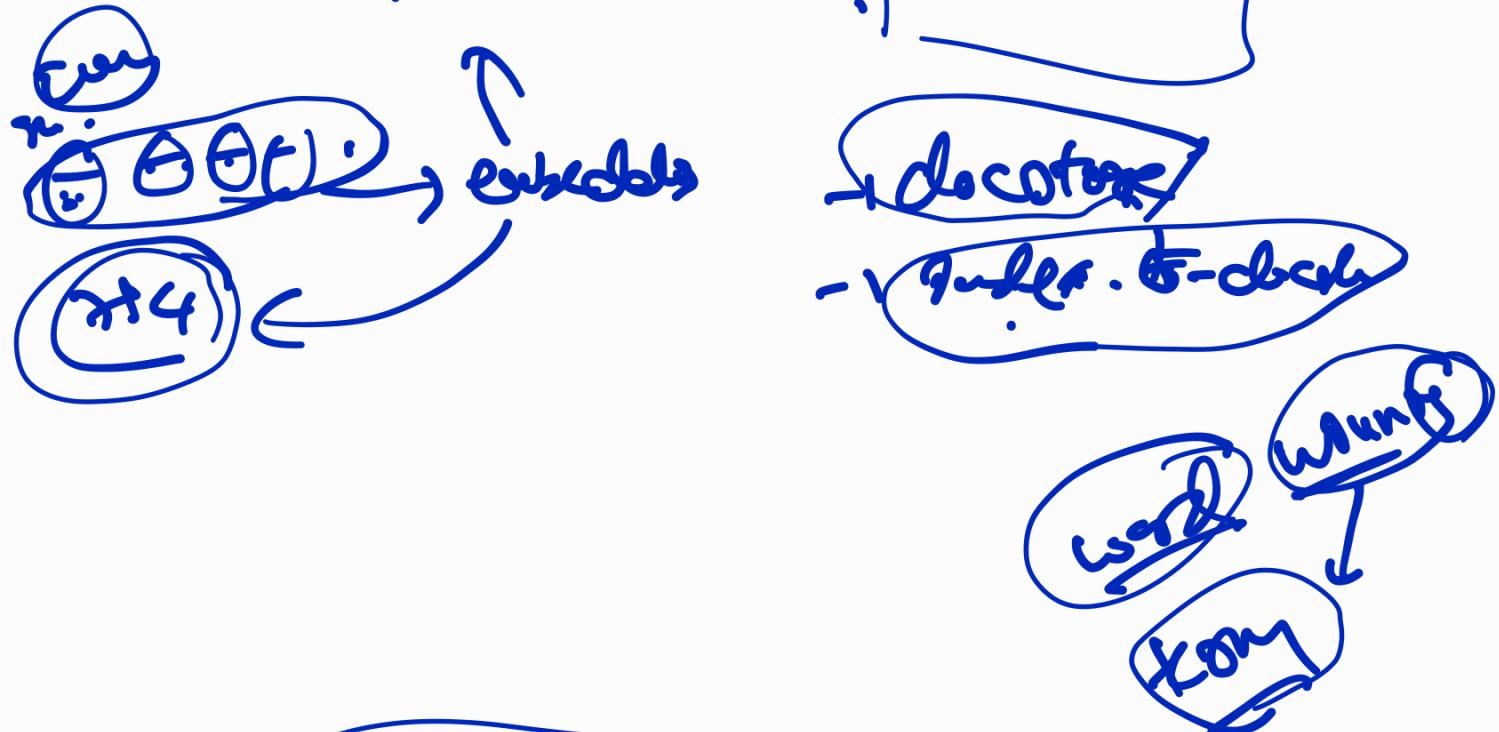
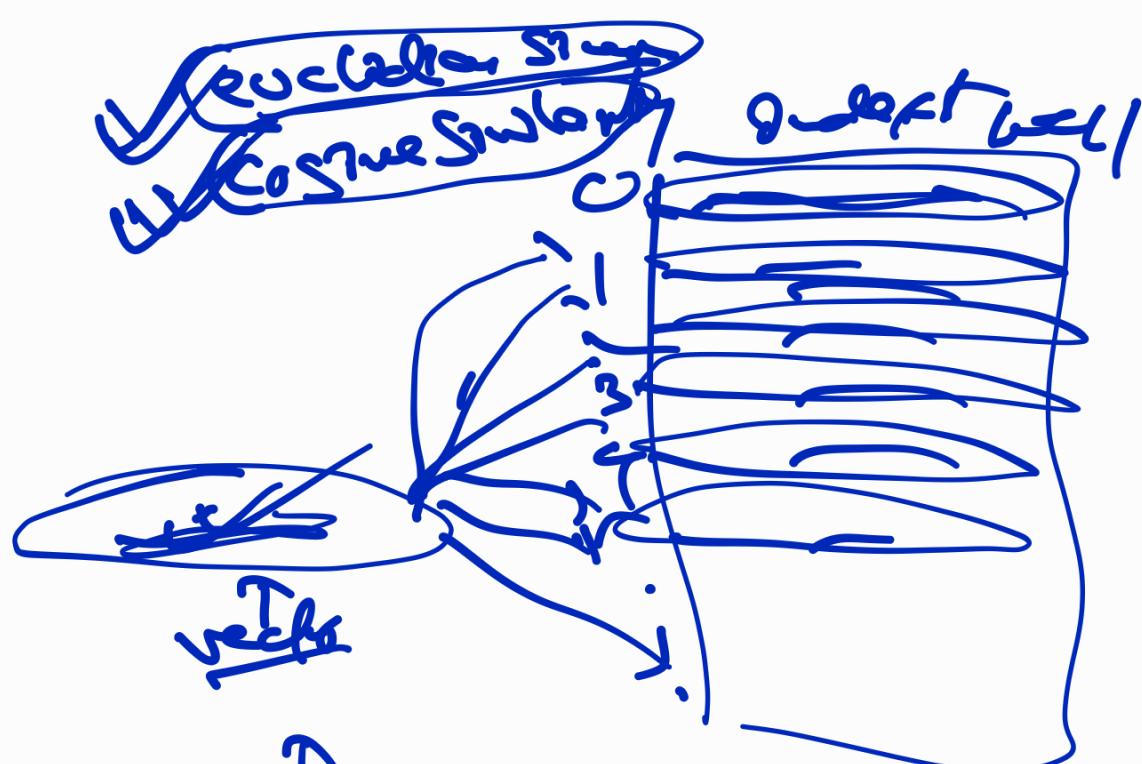


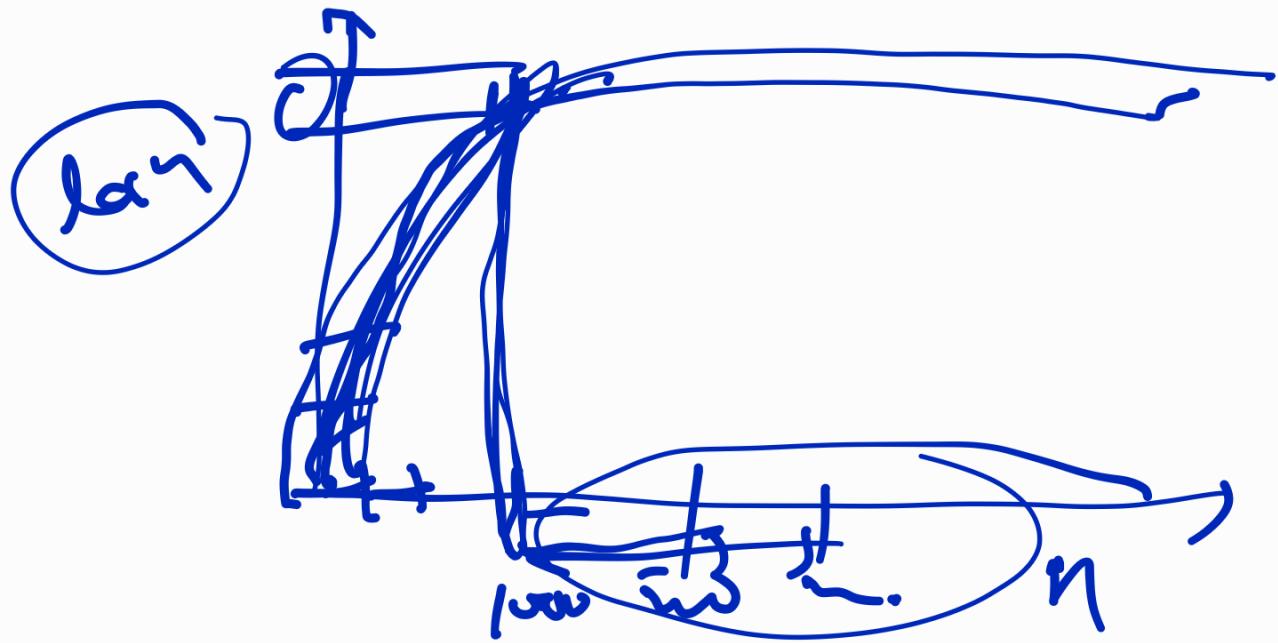


tf-idf = (Term frequency) × (Inverses document frequency)

improved tf-idf = RBM (Restricted Boltzmann Machine)

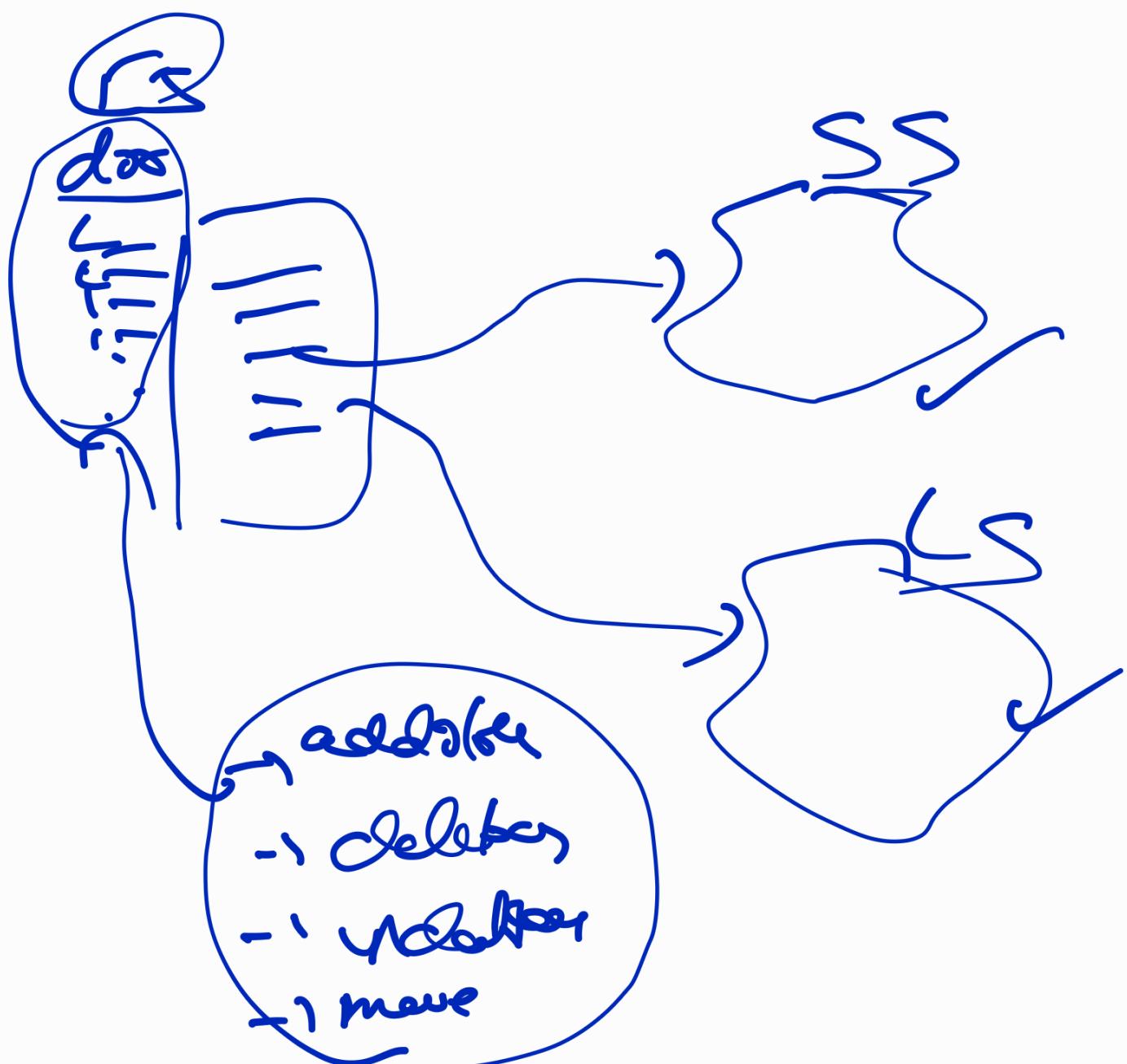






$$\text{let } (\text{freq}, 0.001) \cdot \frac{1}{\ell_0} \cdot c_{\ell_0}$$

$$(\text{TF} \times \text{IDF}) = \frac{\text{Not osc}}{\text{not unosc}} \times \text{IDF}$$



✓ "bat"  
python ->

