Agenda

1. What is RAG, why do LLMs need RAG?
2. (Recap) What is embedding? and how does RAG use it?
3. How to implement RAG and use it for … retrieving data.

1. Create the tools scheme

2. Provide this to model

3. Create the tools

4. Hit model —> tools it wants to uv

5. Backend: Execute tools —> Return response to mode

6. No tools → Model gives output

my AI when I ask for a source

"I made it up"

CUSTOMER SERVICE

5000 pdf, doc ; query : " How do I paint on transparent pixels "

" Content Window " → 2M

‹ Models

**5.2**  **GPT-5.2**  Default ⌄  ⧉                                    Compare    Try in Playground
The best model for coding and agentic tasks across industries

| REASONING | SPEED | PRICE | INPUT | OUTPUT |
| --- | --- | --- | --- | --- |
| ●●●●● | ⚡⚡⚡ | $1.75 · $14 | ⊤ ⊠ ▨ ▨ | ⊤ ▨ ▨ ▨ |
| Highest | Medium | Input · Output | Text, image | Text |

GPT-5.2 is our flagship model for coding and agentic tasks across industries. Learn more in our latest model guide. Reasoning.effort supports: none (default), low, medium, high and xhigh.

✦ 400,000 context window
⤷ 128,000 max output tokens
▢ Aug 31, 2025 knowledge cutoff
♢ Reasoning token support

$\dfrac{50K}{4 \text{ bks}}$ → $\underline{S_1}$ ; $\dfrac{3 \text{ bks}}{4 \text{ bks}}$ → $\underline{S_2}$

$80 -$           $\underline{65}$

$S_1 \gg S_2$

Satisfaction score.

Query → [ GPT ] → Words

Query → [ Embedding model. ] → Embedding.

Embedding - gamma → 768
                     $\underbrace{\qquad}$
                     1D vectors.

5000 PDF            STEP1                          STEP-2

PDF 1  →  [ Embedding mod ] → [ 1D —768 ]    Query
                                                |
PDF 2  →  [ Embedding mod ] → [ 1D —768 ]       ↓
                                           [ Embedding mod ]
PDF 3  →  [ Embedding mod ] → [ 1D —768 ]       |
                                                ↓
  ⋮                                        [ 1D — 768 ]
  ⋮

PDF 5000 → [ Embedding mod ] → [ 1D —768 ]

STEP-3                                    → FAISS INDEX

   Cosine- Sim ( PDF- Vector, Query )

          PDF 3 → highest sim Score

↳ LLM as Content + Query → Answer.

Sent

PDF1 →

Chunk1

Chunk2

Chunk3

Name of perp is …..
w1 — w300

w-250 — 550
Gabbar Singh…..

w — 500— 800

Score.

Relevance

```
# Margin heuristic to decide how many *base* hits to expand
top1 = float(D[0][0])
top2 = float(D[0][1]) if len(D[0]) > 1 else 0.0
margin = top1 - top2
init_topk = 1 if (top1 >= 0.35 and margin >= 0.05) else min(3, topk)
```

— Relivence Score for 2nd document

0.95
top1

0.5
top2

; Margin: top1 — top2
= 0.45