

# RAG Based Psychology Chatbot

## Section 1: Fundamentals of RAG (Retrieval-Augmented Generation)

### What is RAG?

- RETRIEVAL: Finding relevant information (like searching in a library)
- AUGMENTED: Adding/Enhancing (making something better by adding to it)
- GENERATION: Creating responses (like writing an answer)

### Simple Analogy

Think of RAG like an open-book exam!

- Student (LLM): Has general knowledge from studying
- Textbook (External Database): Has specific, detailed information
- RAG Process: Student can look up specific facts in the textbook while answering

### Normal LLM vs RAG-Based Chatbot

Normal LLM chatbot (only model Knowledge)	RAG Based Chatbot (Model + External Knowledge)
<ul style="list-style-type: none"><li>• Knows only what it learned during training</li></ul>	<ul style="list-style-type: none"><li>• Base knowledge + Real-time data access</li></ul>
<ul style="list-style-type: none"><li>• Knowledge cutoff date (e.g., Jan 2025)</li></ul>	<ul style="list-style-type: none"><li>• Can be updated instantly</li></ul>
<ul style="list-style-type: none"><li>• Can't access new information</li></ul>	<ul style="list-style-type: none"><li>• Pulls from your specific documents</li></ul>
<ul style="list-style-type: none"><li>• Might hallucinate (make up facts)</li></ul>	<ul style="list-style-type: none"><li>• Facts are grounded in your data</li></ul>

## Real-World Example: Psychology Chatbot

Psychologist asks: "What are the latest treatment protocols for teenage anxiety?"



**Result:** The chatbot gives you the EXACT protocol from YOUR clinic's database, not outdated general information!

## Why RAG > Fine-tuning for Medical/Psychology Data?

Fine-tuning (Training the Model)	RAG (Retrieval + Generation)
<ul style="list-style-type: none"><li><b>Expensive</b> - Requires lots of compute power</li></ul>	<ul style="list-style-type: none"><li><b>Cheap</b> - Just store documents in database</li></ul>
<ul style="list-style-type: none"><li><b>Time consuming</b> (days/weeks)</li></ul>	<ul style="list-style-type: none"><li><b>Fast</b> - Update documents instantly</li></ul>
<ul style="list-style-type: none"><li><b>Hard to update</b> - Need to retrain entire model</li></ul>	<ul style="list-style-type: none"><li><b>Easy to update</b> - Add new documents anytime</li></ul>
<ul style="list-style-type: none"><li>Risk of catastrophic forgetting (model forgets old knowledge)</li></ul>	<ul style="list-style-type: none"><li>Model keeps all its knowledge</li></ul>
<ul style="list-style-type: none"><li>Still might <b>hallucinate</b></li></ul>	<ul style="list-style-type: none"><li>Responses grounded in real documents (<b>less hallucination</b>)</li></ul>

### **Small note:**

Imagine you have 1000 patient case studies that change weekly. Which approach would you choose and why?

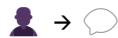
**RAG!**

Because you can update your database daily without retraining. Fine tuning would require retraining every week (impossible and expensive).

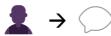
## **Section 2: How RAG Works - End-to-End Architecture**

### **1 . USER ASKS A QUESTION**

Example: "What are the symptoms of teenage depression?"



→



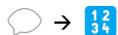
### **2 . QUERY IS CONVERTED TO EMBEDDINGS**

#### **What happens here?**

- Your text question becomes numbers (vectors)
- Think: Converting words into coordinates on a map
- Similar meanings = Similar numbers

Before (Text): question = "What are symptoms of teenage depression?"

After (Embedding - Simplified): question vector = [0.23, 0.87, 0.45, 0.12, ...] # 1536 numbers



→



## Why embeddings?

- Computers can't understand "depression" but they can compare numbers! Words with similar meanings have similar number patterns.

## 3 . SIMILAR CONTENT RETRIEVED FROM VECTOR DATABASE

### What's happening

- Your question embedding is compared to ALL document embeddings in database
- Database finds documents with similar number patterns
- Returns top 3-5 most relevant documents

#### Your Database Has:

- 📄 Document 1: "Depression in Adolescents" [0.25, 0.89, 0.44...] ✓ 95% match!
- 📄 Document 2: "Anxiety Disorders" [0.67, 0.23, 0.78...] 60% match
- 📄 Document 3: "ADHD Treatment" [0.12, 0.45, 0.90...] 30% match

→ Database Returns: Document 1



## 4 . RETRIEVED CONTEXT ADDED TO PROMPT

### Building the final prompt

What the LLM Actually Sees:

#### SYSTEM:

- You are a psychology assistant.

#### CONTEXT FROM DATABASE:

- Depression in teenagers often manifests as:
  - Persistent sadness or irritability
  - Loss of interest in activities
  - Changes in sleep patterns
  - Difficulty concentrating
  - Withdrawal from friends and family

#### USER QUESTION:

- What are symptoms of teenage depression?

#### INSTRUCTION:

- Answer based on the context provided.

This is the MAGIC of RAG! - The LLM now has your exact document content to base its answer on!

 →  (Combined Prompt)

## LLM GENERATES GROUNDED RESPONSE

LLM reads the context and creates answer:

**Chatbot Response:**

Based on clinical guidelines, teenage depression symptoms include:

- Persistent sadness or irritability
- Loss of interest in activities they once enjoyed
- Changes in sleep patterns (insomnia or oversleeping)
- Difficulty concentrating on schoolwork
- Social withdrawal from friends and family



## RAG Components Explained (The 4 Key Players)

### 🌐 Component 1: EMBEDDING MODEL

**Job:** Convert text → numbers (vectors)

**Examples:**

- OpenAI: text-embedding-3-small
- Google: textembedding-gecko
- Open Source: sentence-transformers

**Real-world analogy:** Like a translator converting English to coordinates on a map

### 🗄️ Component 2: VECTOR DATABASE

**Job:** Store embeddings + Find similar ones fast

**Popular Options:**

- **Pinecone:** Easy, cloud-hosted, paid
- **FAISS:** Free, runs locally, need to manage yourself
- **Supabase Vector:** Postgres + vectors, open-source
- **Weaviate:** Open-source, powerful

**Real-world analogy:** Like a super-smart library that instantly finds books similar to what you're looking for

### Component 3: RETRIEVER

**Job:** Take question → Get question embedding → Query database → Return relevant docs

**What it does:**

- Converts your question to embedding
- Searches vector database
- Ranks results by similarity
- Returns top K documents (usually 3-5)

**Real-world analogy:** Like a research assistant who finds the exact pages in textbooks you need

### Component 4: GENERATOR (LLM)

**Job:** Read context + question → Generate answer

**Popular LLMs:**

- **GPT-4:** Very smart, expensive
- **Claude:** Good for long context
- **Gemini:** Google's model
- **Llama 2:** Open-source, can run locally

**Real-world analogy:** Like a smart student who reads the textbook pages and writes a clear answer

### Small Note:

If your embedding model changes, do you need to re-embed all your documents?

### YES!

Different embedding models create different number patterns. If you change from OpenAI to Google embeddings, your stored vectors won't match your new query vectors. You must re-embed everything.

## Section 3: Building RAG-Based Chat System (Technical Stack)

### Data Sources for Psychology Chatbot

#### What Goes into Your Database?

##### **1. Psychology Guidelines**

- APA (American Psychological Association) standards
- Treatment protocols
- Best practices documents

##### **2. Therapist Notes (Anonymized!)**

- Session summaries (with patient names removed)
- Treatment progress notes
- Therapy technique results

##### **3. DSM-5 Summaries**

- Diagnostic criteria
- Disorder descriptions
- Differential diagnosis guides

##### **4. FAQs / Treatment Protocols**

- Common questions from psychologists
- Step-by-step treatment guides
- Crisis intervention protocols

## Chunking Strategies : Why Chunk Size Matters

### What is Chunking?

Breaking large documents into smaller pieces before storing in database.

#### Problem: What if we DON'T chunk?

Imagine you have a 100-page psychology textbook as ONE document:

-  Too much irrelevant information retrieved
-  LLM gets confused by too much context
-  Slow and expensive (processing 100 pages every time)

#### Solution: Chunk the document!

Original Document (5000 words)



↓ CHUNK ↓

 Chunk 1

(500 words)

 Chunk 2

(500 words)

 Chunk 3

(500 words)

 Chunk 4

(500 words)

### Chunk Size Guide

Chunk Size	Best For	Example
<b>Small (200-500 words)</b>	Precise answers, FAQs	"What is CBT?"
<b>Medium (500-1000 words)</b>	Balanced (most common)	"Explain depression treatment"
<b>Large (1000-2000 words)</b>	Complex explanations	"Full therapy protocol"

 **Rule of Thumb:** Start with 500-word chunks (about 1-2 paragraphs). You can always adjust based on testing!

### Overlapping Chunks (Advanced Tip)

Add 50-100 word overlap between chunks so context isn't lost:

```
Chunk 1: [Words 1-500] + [Words 501-550] (overlap) Chunk 2: [Words 501-550] (overlap) + [Words 551-1000]
```

**Why?** If important info spans two chunks, overlap ensures you don't miss it!

## Metadata Usage : Making Search Super Smart

### What is Metadata?

- Extra information about each chunk that helps filter results

### Example: Document with Metadata

```
{  
  "id": "chunk_123",  
  "text": "Cognitive Behavioral Therapy for teen anxiety involves...",  
  "metadata": {  
    "age_group": "13-18",  
    "disorder_type": "anxiety",  
    "severity": "mild",  
    "language": "english",  
    "source": "APA Guidelines",  
    "date_added": "2024-01-15"  
  }  
}
```

### How Metadata Improves Search

#### Without Metadata:

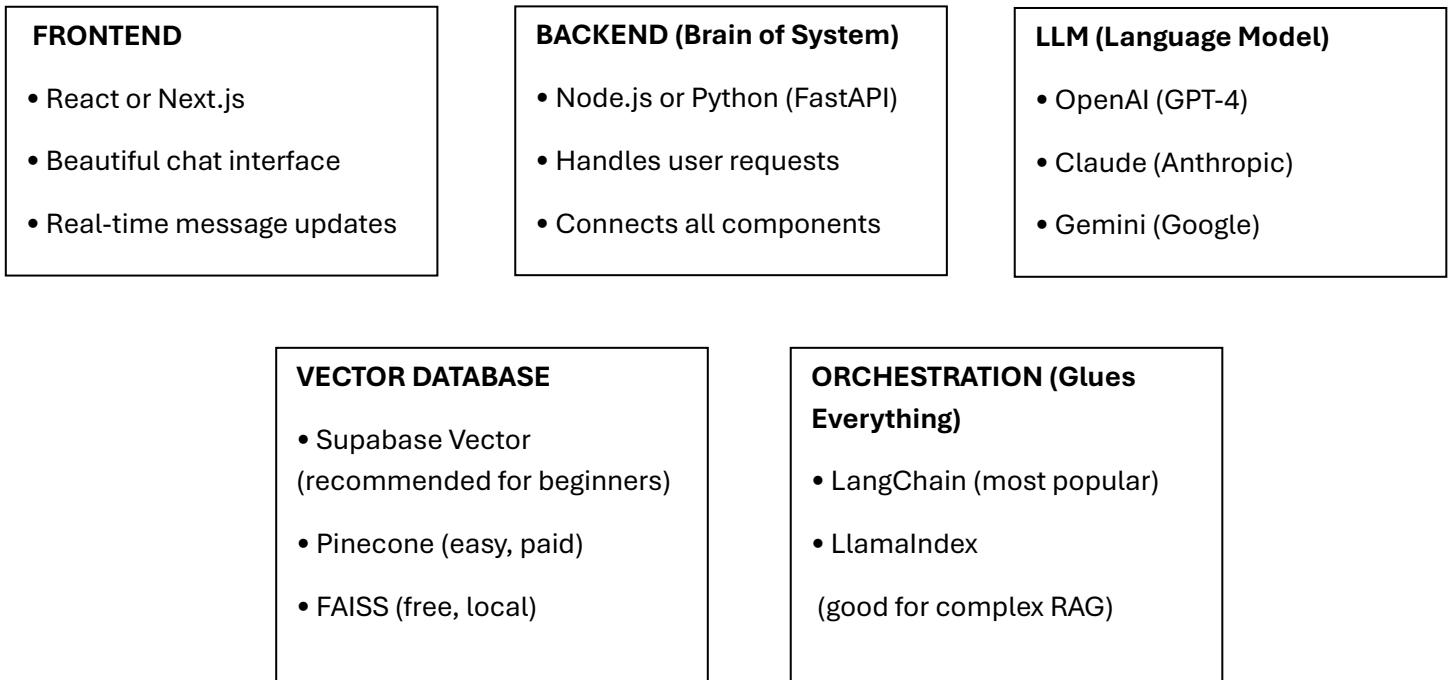
Query: "Treatment for teenage anxiety"  
Returns: ALL anxiety treatments (adults, children, elderly...)  
 You get info about treating 60-year-olds for anxiety!

#### With Metadata:

Query: "Treatment for teenage anxiety" + Filter: age\_group="13-18"  
Returns: ONLY teenage-specific treatments  
 Perfect, relevant results!

## Tech Stack Example

### Chatbot Architecture



### Small Note

**Why would you add "age\_group" metadata instead of just letting the LLM figure it out from the text?**

Because filtering by metadata is FASTER and MORE ACCURATE!

If we have 10,000 documents, searching embeddings + filtering by age\_group="13-18" instantly removes 80% of irrelevant results. The LLM would have to read every document to figure out ages - super slow and unreliable!