

To exit full screen, press **Esc**

# Retrieval-Augmented Generation (RAG) Chatbot System for Clinical Psychology

## 1. Fundamentals of RAG (Retrieval-Augmented Generation)

### 1.1 Concept and Rationale

Retrieval-Augmented Generation (RAG) is an advanced AI architecture combining pretrained Large Language Models (LLMs) with a retrieval system over external knowledge sources. Unlike conventional LLMs, which rely exclusively on their internal weights learned during pretraining, RAG can query an external library of structured and unstructured documents to generate grounded and verifiable responses.

Rationale for Psychology Applications:

- **Factual Accuracy:** Mental health guidance requires absolute fidelity to clinical sources. RAG ensures outputs are sourced from validated guidelines, DSM-5 criteria, and therapist notes.
- **Dynamic Knowledge:** Psychology is a rapidly evolving field. New treatments, clinical trials, and session notes are generated continuously. RAG allows immediate integration without retraining the LLM.
- **Traceability:** Every response can be mapped to the exact paragraph, document, or study it was based on—critical for legal and ethical compliance.
- **Patient Safety:** Reduces hallucination, which is vital when discussing interventions or interpreting patient behavior.

### 1.2 Standard LLM vs. RAG-Based Chatbot

| Feature             | Standard LLM            | RAG-Based Chatbot                              |
|---------------------|-------------------------|--|
| Knowledge Base      | Internal weights only   | Internal weights + curated external sources    |
| Factual Reliability | Prone to hallucinations | Grounded in actual documents                   |
| Update Cycle        | Fixed at training       | Dynamic; updates instantly via added documents |

| Feature              | Standard LLM               | RAG-Based Chatbot                                |
|----------------------|----------------------------|--|
| Traceability         | None                       | Exact reference to source document and paragraph |
| Clinical Suitability | Risky for decision support | Suitable for clinical decision-support roles     |

Analogy:

- Static LLM: Like a clinician trained years ago without access to recent research or notes.
- Dynamic RAG Bot: Like a clinician with instant access to all patient history, treatment protocols, and recent research—capable of providing contextually accurate and current guidance.

### 1.3 Why RAG is Superior to Fine-Tuning for Clinical Data

Fine-tuning adjusts the internal weights of a model to a domain-specific dataset. However, in psychology:

- Risk of Hallucination: Fine-tuned models still generate plausible but potentially inaccurate answers. RAG prevents this by restricting outputs to retrieved documents.
- Citation and Auditability: RAG inherently links responses to sources, which is essential for clinical verification and audits.
- Data Management & Compliance: Patient deletion requests (HIPAA/GDPR) are simple—removing documents from the database is sufficient. Fine-tuning would require retraining to remove sensitive data.
- Rapid Integration of Knowledge: New clinical guidelines, research, or therapeutic techniques can be immediately included by updating the document repository, without retraining.

### 1.4 Real-World Use Cases in Psychology

1. Session Summarization: Transform raw therapist notes into structured SOAP (Subjective, Objective, Assessment, Plan) formats for clinician review.
2. Evidence-Based Intervention Suggestions: Based on retrieved patient history and DSM-5 criteria, the chatbot can propose recommended therapy techniques.

3. Risk Identification: Scan longitudinal patient data to flag early signs of self-harm, depression, or PTSD indicators.
4. Guideline Integration: Automatically cross-reference treatment protocols with recent research papers to suggest updated methods.

## 2. How RAG Works – End-to-End Architecture

RAG is structured as a Retrieve-then-Generate pipeline, ensuring that all responses are grounded, traceable, and contextually relevant.

### 2.1 Inference Pipeline Detailed Flow

1. User Query Input: A clinician queries the system, e.g., *"What CBT techniques are recommended for adolescents with social anxiety?"*
2. Semantic Embedding:
  - The query is converted into a vector representation using embedding models such as text-embedding-3 or domain-specific models like BioBERT.
  - Each vector encodes the semantic meaning of the query, not just its literal words.
3. Vector Retrieval:
  - The vector database (Supabase Vector, Pinecone, FAISS, or Weaviate) is queried to find document chunks whose embeddings are most similar to the query vector.
  - Hybrid retrieval may combine cosine similarity with keyword-based search to ensure critical terms (e.g., "Fluoxetine") are matched exactly.
4. Context Injection:
  - Retrieved chunks are inserted into the LLM prompt. Example:
5. Using only the following therapist notes and DSM-5 excerpts: [Retrieved Context], summarize the recommended CBT techniques for this patient.
6. Grounded Response Generation:
  - The LLM produces a response strictly based on the injected context.
  - Low temperature (0.0–0.1) settings reduce creative but potentially unsafe hallucinations.

- Each output can reference source documents for clinician verification.

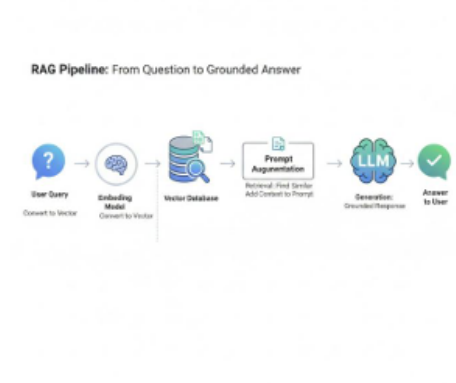
## 2.2 Core Components

- **Embedding Models:** Convert both user queries and documents into high-dimensional vector representations. Domain-specific embeddings improve understanding of clinical terms.
- **Vector Databases:** Efficiently store embeddings and perform semantic similarity searches. Examples: Supabase Vector (pgvector), Pinecone, FAISS, Weaviate.
- **Retriever:** Filters relevant chunks based on similarity and metadata (e.g., disorder type, patient age, session date).
- **Generator (LLM):** Produces text output using the retrieved context. Can be OpenAI GPT-4o, Claude 3.5/4, or Gemini.

## 2.3 Detailed Considerations

- **Chunk Overlap:** Overlapping chunks improve context continuity.
- **Hybrid Search:** Combines vector similarity with exact keyword matches to maintain precision for critical medical terms.
- **Traceability:** Metadata and chunk references ensure every response can be traced to its source.
- **Grounding Verification:** Optional automated checks compare generated responses against retrieved documents to prevent hallucinations.

## 2.4 RAG Flow Diagram



Illustrates Query → Embedding → Retrieval → Context Injection → Grounded Generation, highlighting clinical document sources like DSM-5 summaries and therapist notes.

## 3. Building a RAG-Based Chat System (Technical Stack)

### 3.1 Data Sources

- DSM-5 summaries and diagnostic criteria
- Anonymized therapist session notes
- Evidence-based treatment protocols and FAQs
- Research papers and meta-analyses

### 3.2 Chunking Strategies

- **Chunk Size:** Optimal 500–1,000 tokens per chunk to fit LLM context windows.
- **Semantic Chunking:** Splits occur at natural breaks in clinical topics to maintain integrity.
- **Overlaps:** Small overlaps preserve context for sequential content.

### 3.3 Metadata Schema for Precision Retrieval

- patient\_id (anonymized)
- age\_group (child, adolescent, adult)
- disorder\_type (anxiety, depression, PTSD, etc.)
- severity\_level (mild, moderate, severe)
- session\_date
- source\_type (DSM-5, session note, guideline, research paper)

### 3.4 Technology Stack

| Layer           | Technology   |
|-----------------|--|
| Frontend        | React / Next.js (clinician dashboard)                |
| Backend         | Python (FastAPI) or Node.js                          |
| Orchestration   | LlamaIndex (deep indexing of multi-source documents) |
| Vector Database | Supabase Vector / Pinecone                           |
| LLM             | OpenAI GPT-4o, Claude, Gemini                        |
| Embedding Model | OpenAI text-embedding-3 or BioBERT                   |

- LlamaIndex provides robust connectors for complex document structures like therapy notes with tables or bullet points.
- Supabase Vector ensures HIPAA-compliant storage since it integrates vectors into a relational SQL database.
- Backend can securely manage document ingestion, vector creation, and retrieval while controlling access based on user roles.

#### **4. Psychology-Focused Chatbot Design (Domain-Specific Logic)**

- Role: Assist psychologists in clinical decision-making; do not replace practitioner judgment.
- Key Use Cases:
  - Session summarization into SOAP notes
  - Therapy recommendation support
  - Risk flagging (self-harm, depression, PTSD)
  - Treatment history retrieval for continuity of care

Prompt Design:

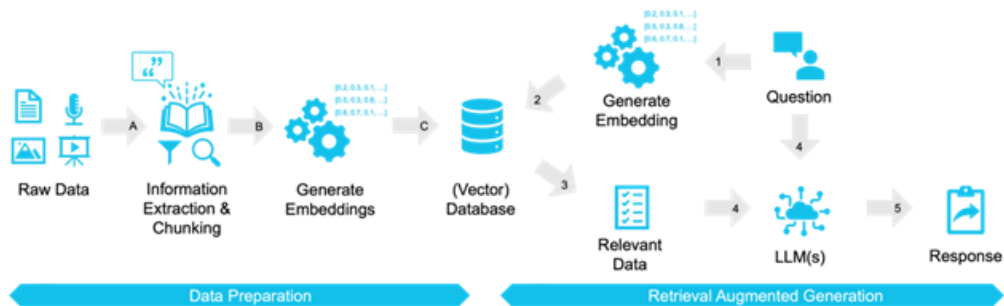
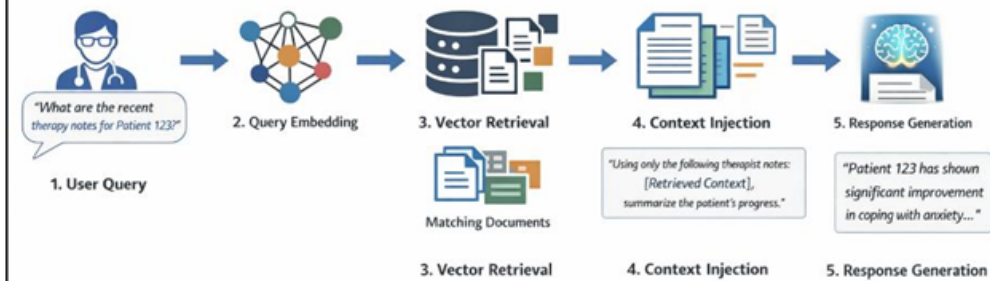
- Professional and empathetic tone
- Explicitly non-diagnostic
- Safe prompt example:

You are a professional clinical assistant. Your tone is empathetic but objective. Never provide a diagnosis. If the context does not contain the answer, state that the information is unavailable.

#### **5. Safety, Ethics, Privacy & Compliance**

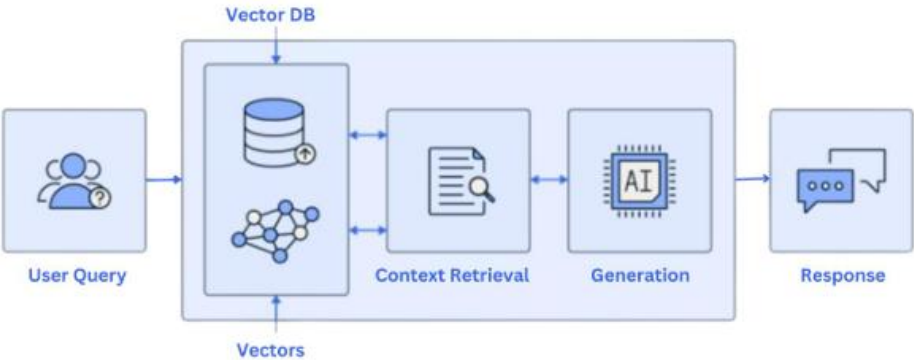
- Data Privacy: All patient data must be anonymized before LLM or vector database ingestion.
- Access Control: Role-based permissions, row-level security, and separation between admin and clinician access.
- Hallucination Prevention: RAG architecture ensures output is strictly grounded in retrieved documents.
- Crisis Management: Hard-coded detection of self-harm or suicidal intent triggers immediate escalation to clinicians and emergency protocols.
- Logging & Audit Trails: Every query, retrieved context, and output is logged for traceability and compliance.

## End-to-End RAG Process





## RAG Architecture



### Chunking Overlap and Metadata Examples

| Chunk ID | Chunk Text (Excerpt)  | Overlap Tokens | Metadata  | Notes  |
|----------|---|----------------|---|--|
| C001     | "Patient reports persistent anxiety over social interactions. CBT exercises include exposure therapy and journaling..." | 50             | patient_id: 12345, age_group: adolescent, disorder_type: anxiety, severity_level: moderate, session_date: 2025-03-15, source_type: session_note | First paragraph of session; overlaps with next chunk to preserve context |
| C002     | "CBT exercises include exposure therapy and journaling. Homework assigned: weekly exposure tasks..."                    | 50             | patient_id: 12345, age_group: adolescent, disorder_type: anxiety, severity_level: moderate, session_date: 2025-03-15, source_type: session_note | Overlaps with previous chunk to ensure no data loss between splits       |

| Chunk ID | Chunk Text (Excerpt)   | Overlap Tokens | Metadata   | Notes  |
|----------|--|----------------|--|--|
| C003     | "DSM-5 criteria for Social Anxiety Disorder: marked fear or anxiety in social situations. Symptoms must persist for at least 6 months..."      | 0              | source_type: DSM-5, disorder_type: anxiety, age_group: adolescent                                  | Reference chunk from DSM-5 guideline; no overlap needed                |
| C004     | "Recent research indicates mindfulness exercises improve anxiety symptoms in adolescents. Incorporate 10-minute daily mindfulness sessions..." | 20             | source_type: research_paper, disorder_type: anxiety, age_group: adolescent, publication_year: 2024 | Overlaps slightly with next chunk for continuity                       |
| C005     | "Treatment protocol recommends monitoring progress weekly and adjusting CBT techniques based on patient response..."                           | 20             | source_type: treatment_protocol, disorder_type: anxiety, age_group: adolescent                     | Overlaps with previous research excerpt to link evidence with protocol |

**Notes:**

- Overlapping chunks maintain context for sequential content.
- Metadata enables precise retrieval for patient-specific queries.
- Supports integration of multiple source types (DSM-5, session notes, research papers, protocols).