**◎ ChatGPT**

# Therapy Chatbot Implementation Plan

## Introduction and Objectives

Developing a therapy chatbot requires combining advanced AI with proven psychological frameworks. The goal is to create a text-only conversational agent that can assume different supportive personas (e.g. professional therapist, friend, coach, humorous companion) while remaining grounded in Cognitive Behavioral Therapy (CBT) principles. To ensure accurate and helpful responses, the system will use Retrieval-Augmented Generation (RAG) – integrating a domain-specific knowledge base with a Large Language Model (LLM) via the Gemini API. This approach tailors the chatbot's knowledge to therapy-specific content, preventing hallucinations and improving relevance [1] [2]. Long-term user memory will be handled via Pinecone (a vector database) so the chatbot can recall past conversations across sessions, enabling more personalized and context-aware interactions [3]. Safety is paramount: if a serious mental health crisis is detected, the bot will display a pre-defined help message or "flashcard" urging the user to seek professional help, rather than attempting to handle it itself. Data privacy and compliance (while not in the MVP scope) are noted for future implementation.
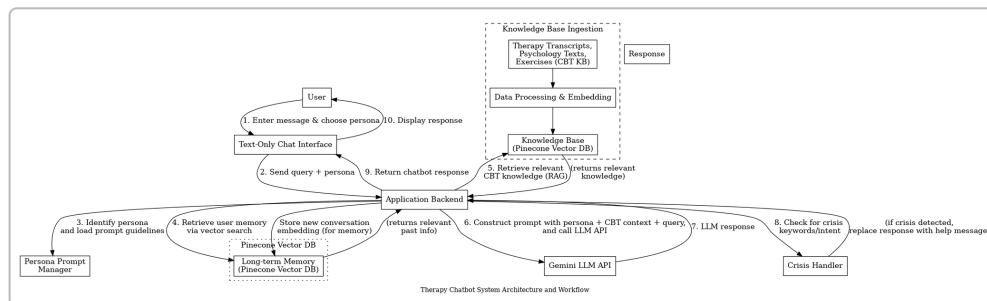
## System Architecture and Workflow



*Figure 1: High-level system architecture and workflow for the therapy chatbot. The user's query and selected persona are processed by the backend, which retrieves relevant context from a vector database (Pinecone) and uses the Gemini LLM API to generate a response grounded in CBT knowledge. A crisis handler checks the output for safety before returning the answer to the user.*

Figure 1 illustrates the end-to-end workflow of the therapy chatbot system. When a **user** enters a message and chooses a persona, the **text-only chat interface** sends this to the **application backend** (Step 1–2). The backend orchestrates several steps in sequence: 1. **Persona Selection:** The backend consults the **Persona Prompt Manager** to identify which persona profile the user chose and loads the corresponding prompt instructions (Step 3). These instructions define the chatbot's style and tone (therapist, friend, coach, or humorist) via prompt engineering. 2. **Context Retrieval (RAG):** The user's query is embedded into a vector representation and used to perform similarity search against two Pinecone indexes – the **long-term memory store** (past user conversations) and the **CBT knowledge base** (therapy domain knowledge). This yields any relevant prior conversation snippets (personal context) and relevant knowledge documents (factual context) (Steps 4–5). Using Pinecone as a "long-term memory" allows the chatbot to bring in

1

pertinent details from earlier sessions or user history without hitting context length limits [4] [3]. Meanwhile, retrieving from the CBT knowledge base ensures the response is backed by *real therapeutic data* rather than just the LLM's training [3]. 3. **LLM Response Generation:** The backend then constructs a prompt for the LLM. It combines (a) the persona-specific system prompt (defining the role and style), (b) the retrieved context (memory and knowledge snippets), and (c) the user's latest query. This full prompt is sent to the **Gemini LLM API** for completion (Step 6). Using RAG, we attach the retrieved text as additional context in the prompt so that the LLM's answer is grounded in the provided knowledge [5]. For example, the prompt might include a brief CBT technique description from the knowledge base if it's relevant to the user's issue, ensuring the LLM can incorporate that into its advice. The Gemini model processes this and returns a response (Step 7). 4. **Crisis Handling:** Upon receiving the LLM's response, the backend passes it (and/or the user's message) through a **crisis handler module** (Step 8). This module checks for indications of serious mental health crisis (e.g. explicit suicidal ideation or self-harm intent). If such a crisis is detected, the system will override the LLM's reply with a pre-crafted **emergency help message** (a "flashcard" recommending the user seek help from a mental health professional or crisis hotline). This safety response is displayed to the user instead of a normal chatbot answer. (If no crisis is detected, the LLM's original response is used.) 5. **Response Delivery and Memory Update:** The **chat interface** then returns the final chatbot response to the user (Step 9–10). In parallel, the conversation is logged: the backend will embed the latest user query (and possibly the bot's answer) and upsert it into the **Pinecone memory index** (as a new vector) to enrich the long-term memory for future interactions. Storing conversations as embeddings in Pinecone ensures the chatbot can later retrieve relevant parts of the dialog history when needed, enabling continuity across sessions.

This architecture is designed for modularity and scalability. Each component (persona manager, retrieval module, LLM interface, etc.) can be developed and improved independently. The use of a vector DB for both knowledge and memory provides flexibility in expanding the knowledge base or the amount of remembered context without being constrained by the LLM's fixed context window [6] [7]. The **Gemini API** is used for cost-efficiency – presumably leveraging Google's generative model which may offer favorable pricing – while still providing strong performance. By grounding Gemini's output with specific retrieved info, we mitigate the risk of generic or incorrect advice, focusing the chatbot on evidence-based therapeutic content.

## Key Components and Modules

To implement the above workflow, we will build the system as a collection of key components/modules, each with specific responsibilities:

- **1. User Interface (Text Chat Frontend):** A simple web or mobile chat interface that allows users to interact via text. It should present a menu or option for persona selection at the start (or even dynamically during chat) and then display the conversation messages. Tasks for this module include designing a clean chat UI, implementing persona selection (e.g. a dropdown or buttons to choose Therapist, Friend, Coach, Funny Companion), and ensuring messages are sent to/back from the backend seamlessly. The UI can be minimal for MVP (even a command-line or web form) but must clearly show the persona's "name" or style so the user knows which mode they are talking to.

- **2. Application Backend:** The core orchestration layer handling conversation logic. This will likely be a server (e.g. Node.js, Python FastAPI, etc.) that exposes an endpoint for the chat interface to send messages. Key sub-modules of the backend are:

- **Persona Prompt Manager:** Stores and manages the definitions of each persona. For each persona, we will craft a distinct system prompt or set of instructions that characterize the chatbot's style, tone, and behavior. For example, the "Therapist" persona prompt may include: *"You are a licensed therapist using CBT techniques. Speak in a professional, empathetic tone, ask gentle probing questions and guide the user to challenge negative thoughts."* The "Supportive Friend" persona prompt might be: *"You are the user's close friend: informal, comforting, using upbeat encouragement and empathy (but still subtly applying CBT ideas)."* We will create similar profiles for the "Coach" (motivational, goal-focused, action-oriented) and "Funny Companion" (light-hearted, using humor appropriately to cheer the user up) personas. The Persona Manager's task is to inject the correct persona instructions into the LLM prompt depending on user choice. This can be done by prepending a system message for the LLM [8] . (During development, we will test and iterate these prompts to ensure the persona "voice" comes out clearly in the responses.)
- **Conversation (Dialogue) Manager:** Orchestrates each turn of conversation. It keeps track of the conversation state (possibly short-term context of the current session) and interfaces with the memory and knowledge base. On each user query, it triggers the retrieval of relevant memory and knowledge (calls the Vector DB), constructs the full prompt (persona instructions + retrieved context + user query), calls the Gemini LLM API, and then processes the response (including handing off to Crisis Handler if needed). It also triggers storing the new conversation entries into Pinecone memory. Essentially, this module coordinates Steps 3–9 from the workflow. We might implement this logic using a framework like LangChain for convenience [9] , which can simplify chaining the retrieval and LLM calls, though a custom implementation is fine for full control.
- **Gemini LLM API Client:** A utility module that handles the HTTP requests to the Gemini API for text completion. It will format the prompt in the API's expected format (e.g. including system message for persona, user message, retrieved context as additional messages or in the system prompt). We need to manage API keys securely and handle rate limits or errors gracefully. Since cost is a concern, this module should also support passing parameters like temperature, max tokens, etc., to balance quality vs. cost. (For instance, use a moderate temperature for creative but safe responses, and possibly limit tokens to a reasonable length per response to control costs.)
- **Vector Database Interface (Pinecone Client):** This module manages all interactions with Pinecone. It will handle two primary indexes:
  - *Knowledge Base Index:* storing vector embeddings of the therapy knowledge base documents (transcripts, articles, exercises). This client provides methods to query the index with a new query embedding and fetch top-$k$ similar documents for context.
  - *Long-Term Memory Index:* storing vector embeddings of user conversation snippets (with metadata like user ID, persona, timestamp). The client can query this index with the current user query's embedding to find semantically similar past conversations (or important facts the user shared) to remind the chatbot. It also provides methods to add new embeddings to the index as the conversation grows (with appropriate batching or real-time upserts).
  - The Pinecone interface must be configured with appropriate index dimensions (matching the embedding model output size) and metadata schema (e.g. storing text chunks and tags). Pinecone's advantage is it scales to lots of embeddings and has no fixed context length, so even as memory and knowledge grow, queries remain fast [7] .
- **Crisis Handler:** A safeguard component that inspects messages for crisis-related content. For MVP, this can be a simple rule-based checker (e.g. scanning the user message and LLM response for certain keywords or phrases like "suicidal, want to die, end it all," etc., or using sentiment scores). If a high-risk situation is flagged, the handler can override the bot's reply with a prepared **emergency message**. For example: *"Important: I'm detecting you might be feeling hopeless or considering self-harm. Please remember you are not alone. It might help to reach out to a mental health professional or call a*

*crisis line. You deserve support. (If this is an emergency, please seek immediate help.)"* This will be shown as a special notice (e.g. in a highlighted style) to the user. The conversation could either be ended or gently steered to ensure the user is safe. We explicitly will **not** attempt to continue therapy in such cases, aligning with best practices that AI should defer to human help in crisis [10] . (In future, we might integrate an actual crisis line API or alert mechanism, but for now a static referral message suffices.)

- **Logging & Analytics (optional for MVP):** It's wise to log interactions (with user consent) for debugging and improvement. This module would record user queries, bot responses, retrievals made, flags triggered, etc. Given privacy will be addressed later, for now logs can be basic and stored securely. We will *not* use any data for other purposes without consent. In future, analysis of logs can help refine the model and prompts.

- **3. Knowledge Base (Content Repository):** The collection of documents and data that form the grounding knowledge for the chatbot. For our therapy chatbot, the knowledge base will include:

- *Therapy Session Transcripts:* Real or simulated transcripts of therapist-client conversations (especially using CBT techniques). These provide examples of how a therapist responds to various situations. Having them in the knowledge base means the chatbot can retrieve snippets of similar cases or therapist explanations to mimic effective responses. (Datasets like the HOPE corpus with ~202 counseling transcripts [11] or others from research can be used to source this content.)
- *Psychology Textbook Excerpts and Research Papers:* Authoritative descriptions of mental health conditions, CBT theory, cognitive distortions, therapeutic principles, etc. These ensure the bot has factual, evidence-based material to draw from. For example, definitions of CBT techniques, explanations of why certain thoughts occur, or research-backed coping strategies.
- *Therapeutic Exercises and Worksheets:* Guides for practical exercises (thought records, mood tracking, breathing exercises, behavioral activation tasks, etc.). These are very useful for a CBT chatbot to provide concrete activities that can help improve the user's mood or skills. By including them, the bot can "recommend" an exercise when appropriate and even walk the user through it step by step, improving the chat's tangible benefit.

This Knowledge Base is not a single module per se, but rather a dataset that we will prepare and then load into Pinecone. The tasks involved here are detailed in the next section (knowledge base structuring and preprocessing). In short, we'll collect relevant documents, clean and chunk them, embed them with an embedding model, and populate the Pinecone index prior to running the chatbot.

## Knowledge Base Structuring and Preprocessing

Proper structuring of the knowledge base is critical for effective retrieval. We will follow these steps to build a robust CBT knowledge base:

- **Data Collection:** Gather a diverse set of sources:
- *Therapy transcripts:* We may use publicly available counseling transcripts (from research datasets or anonymous forums) and possibly augment with synthetic dialogues if needed. The focus is on conversations illustrating CBT methods (identifying automatic thoughts, challenging them, etc.).
- *CBT Guides and Text:* Extract key sections from CBT workbooks, psychology textbooks, or reputable websites (e.g., sections on cognitive distortions, CBT definitions, steps of exercises like thought

challenging). Ensure we have content covering common issues (anxiety, depression, stress) and techniques to address them.

- *Exercises and Tips:* Include step-by-step instructions for exercises (for example, "deep breathing exercise," "guided imagery script," "activity scheduling instructions"). These often come from therapy handouts or mental health websites.

- All content must be either public domain, licensed for use, or summarized from sources to avoid copyright issues. (For MVP, even a few dozen high-quality documents will suffice as a starting KB.)

- **Preprocessing:** Once collected, each document will be cleaned and split into chunks suitable for the LLM context:

- Remove any personally identifiable information or irrelevant text (especially in transcripts, ensure anonymity and remove filler chat not useful for therapy).
- Normalize text (fix encoding issues, expand abbreviations if necessary, and add any needed metadata like titles or source tags).
- **Chunking:** Break documents into semantically coherent chunks, typically a few sentences or one paragraph each, such that each chunk (when embedded) captures a single idea. We aim for chunk sizes perhaps in the range of 100 to 300 tokens so that a few retrieved chunks plus user query will fit in the LLM prompt window. For example, a therapy transcript might be split by dialogue turns or by topic shifts, and an article might be split by paragraphs or subsections.
- If certain context is needed to understand a chunk, we can include brief overlaps or context windows between chunks, or store an identifier so multiple related chunks can be fetched.

- Tag each chunk with metadata: e.g., source type (transcript vs. article vs. exercise), title or conversation ID, speaker (if transcript, was this the therapist's statement or client's?). These metadata can help filter results if needed (for instance, we might prefer retrieving therapist statements from transcripts rather than client utterances).

- **Vector Embedding:** Choose an embedding model to convert text chunks into vectors. Since we are using the Gemini API for the main model, we can use Google's embedding models (if available via Vertex AI) or an open-source embedding like SentenceTransformers. The embedding model should ideally be tuned for semantic similarity in the domain of text (OpenAI's **text-embedding-ada-002** or similar works well generally). We will use the same model for all knowledge chunks and for queries to ensure they live in the same vector space [5] .

- Each chunk of text is passed to the embedding model to get a high-dimensional vector representation. We then **upsert** those into the Pinecone Knowledge Base index, storing the vector with the chunk's text and metadata [5] . Pinecone will handle indexing for similarity search.

- **Index Configuration:** In Pinecone, we'll create (at least) two separate indexes:

- `therapy-knowledge-index` : for the knowledge base chunks. Set an appropriate dimension (matching the embedding, e.g. 1536 for ada-002). We might allow metadata filtering and use an approximate similarity metric (cosine similarity is common for embeddings).

- `user-memory-index` : for long-term conversation memory. Similar setup but will store shorter texts (user and bot messages or summary thereof). Metadata will include a user ID to ensure we only retrieve a given user's own history, and maybe persona info.

- By separating them, we can query them independently or together. For example, the system might always do two searches: one in `therapy-knowledge-index` for relevant domain info, and one in `user-memory-index` for personal context, then combine the results. This modularity is useful.

- **Maintenance:** We should plan how to update the knowledge base over time:

- For MVP, it's mostly static data prepared once. But if new research or user feedback indicates gaps, we will add more documents and re-embed as needed.
- Pinecone supports updating/deleting vectors, so we can refine the KB continuously.
- (For long-term, consider building an automated pipeline to ingest new content or user-shared resources, but that's future work.)

Through this process, we ensure the knowledge base is **structured and optimized** for retrieval. When the chatbot is asked something, Pinecone will return the most semantically relevant chunks, which the backend can feed into the LLM prompt. For instance, if the user says "I feel like I'm a failure at everything," the system might retrieve a textbook chunk explaining the *"all-or-nothing thinking"* cognitive distortion and a transcript chunk where a therapist helped a client reframe a failure thought. Those will help the LLM respond with accurate CBT-guided feedback rather than just empathy alone. This method grounds the conversation in established therapeutic knowledge, a core advantage of RAG [5].

## Persona Prompt Engineering Approach

Implementing multiple personas in one chatbot is achieved through careful prompt engineering. Each persona will be defined by a distinct **system prompt** (or prefix prompt) given to the LLM at the start of the conversation (and maintained throughout). Here's how we will approach each persona:

- **Professional Therapist Persona:** This is the default and most crucial persona. The prompt will establish the chatbot as a licensed therapist using CBT. Key characteristics: professional but warm tone, uses reflective listening, asks open-ended questions, and introduces CBT concepts in user-friendly ways. The system prompt might say, for example: *"You are a compassionate, experienced therapist. You use Cognitive Behavioral Therapy techniques to help the user identify negative thought patterns and develop healthier perspectives. Speak in the first person, with empathy and professionalism. Guide the conversation, but let the user express themselves. Provide practical exercises or reframing when appropriate."* This ensures the LLM's responses include CBT methods (like cognitive restructuring, thought logs, etc.) framed as a therapist would present them. Indeed, instructing the model to align with CBT principles can be done explicitly in the prompt [8].

- **Supportive Friend Persona:** The prompt for this persona casts the AI as the user's close friend who cares about their well-being. The tone is informal, reassuring, and positive. It likely avoids clinical jargon and speaks more like a peer: *"You are chatting with your best friend who is down. Your role is to listen and support them. Respond with warmth, encouragement, and maybe light humor, as a friend would. You still try to help them challenge negative thoughts or see positives (drawing on CBT ideas) but in a very casual, friendly manner – not as a therapist or authority."* This persona will use more emotive

language (e.g. *"That sounds really hard. I'm here for you. Remember when you accomplished X? I know you have it in you."*) and can share personal-ish sentiments (though fictional). We will prompt-engineer it to be *supportive and uplifting*, while subtly incorporating helpful techniques (for example, encouraging the user to reconsider a negative thought gently or suggesting a fun activity to cope, aligning with CBT behavioral activation).

- **Motivational Coach Persona:** This prompt positions the AI as a life coach or personal coach. Tone: upbeat, goal-oriented, slightly more formal than friend but more cheerleader-like. It focuses on actionable guidance, setting goals, and motivating the user. The prompt might include: *"You are an enthusiastic life coach. You help the user set small goals, challenge their self-doubt, and take positive action. Use motivational language and confident tone. Emphasize strengths and progress. When the user is stuck in negative thinking, help them reframe it constructively. Keep the conversation future-focused and pragmatic."* This persona will likely give concrete suggestions (e.g. *"How about we set one small goal for tomorrow? You got this!"*), aligning with CBT's focus on behavior change and skills.

- **Funny Companion Persona:** For a lighter touch, this persona uses humor to improve the user's mood. The prompt defines the AI as a witty, caring companion who cracks gentle jokes or uses playful language to make the user smile. It must be **very careful** not to make fun of the user or trivialize serious issues – the humor should be self-deprecating or situational, aimed at easing tension. Prompt example: *"You are the user's quirky but caring friend who always finds a way to make them smile. Use light humor, jokes, or funny analogies to help cheer them up, but always be respectful of their feelings. If they mention a problem, you can respond with a touch of wit to lighten the mood while still offering helpful perspective."* This persona can employ techniques like humorous reframing (e.g. turning a negative thought into a silly exaggerated scenario to show it's irrational, a known CBT technique for defusing thoughts). We will test the boundaries of this persona thoroughly to ensure it's appropriate. The prompt will explicitly forbid dark or offensive humor and remind the AI to maintain empathy.

**Implementation:** The Persona Prompt Manager will likely maintain a template for each persona's system message. During a conversation initialization (or when persona is switched), the system message sent to the LLM is set to the chosen persona's template. All personas will share some core instructions as well, such as staying on topic of mental health support, using the provided knowledge base context, avoiding giving medical or medication advice (we might include a disclaimer in prompt about not being a doctor), and not diverging into unrelated areas. We also instruct the model to *stay in character* – e.g. the friend persona shouldn't suddenly speak like a formal therapist. Each response from the LLM should reflect the persona's style guidelines.

We will likely refine these prompts iteratively. A/B testing responses from different prompt phrasings will help dial in the right tone. The beauty is that no model retraining is needed; we simply craft the textual instructions. As an example from a similar project, a system prompt might say: *"You are an AI therapist…your responses should be empathetic, structured, and aligned with CBT principles"* [8] – this clearly steers the LLM's persona. We will do the equivalent for each persona.

Additionally, we might include a few **one-shot or few-shot examples** in the prompt to solidify persona behavior. For instance, showing a short example dialogue of the Funny Companion responding to a sad statement with a light joke. However, examples consume token budget; we will use them only if needed. Given a strong base model like Gemini, a well-written descriptive instruction might suffice.

Finally, we must ensure the persona differences do not contradict the core therapeutic goals. All personas ultimately should promote positive coping, cognitive reframing, and emotional support, just in different styles. We'll maintain a unified "ethical guardrails" instruction across them (e.g., do not produce harmful content, do not give lethal advice, etc., likely provided by the Gemini API's own safety mechanisms too).

## Aligning RAG with the CBT Framework

One challenge is effectively integrating the retrieved knowledge with the conversational CBT approach. We want the chatbot's responses to be both **knowledge-grounded** and **therapeutically sound**. Here's how we plan to align the RAG approach with CBT:

- **Relevant Retrieval:** We will optimize our retrieval to fetch information that directly supports CBT interventions. This might mean tuning the search queries or using metadata filters. For example, if the user message contains cognitive distortions like *"I'm a total failure,"* we could programmatically add a keyword like "cognitive distortion" or "thinking trap" to the search query to ensure we retrieve a relevant explanation from the knowledge base (if the base is tagged accordingly). LangChain or a custom retrieval wrapper can let us modify queries or do hybrid search (keyword + semantic). Ensuring the right info is retrieved is the first step to alignment.

- **Inserting Knowledge into Prompt:** Once we have the top relevant chunks (say the top 2–3 from the knowledge base), we need to include them in the LLM prompt in a *useful* way. Typically, we might attach them either as an extended system message (e.g., "Knowledge: [extracted text]") or as additional conversational turns (some implementations use a format like: system: "Here are some relevant excerpts: …"). We will decide a format that Gemini's API accepts and that the model can easily incorporate. We want the model to treat these as reference material: facts or suggestions it should use in forming its answer. The prompt might say: *"Use the following information to inform your response:"* and then list the snippets. The model should then weave that into its reply, improving accuracy. This method has been shown to give *grounded, relevant responses backed by real data sources* [3] .

- **Guiding the LLM to Apply CBT:** In the persona system prompt (especially for the therapist persona), we explicitly instruct the model to use CBT techniques. For example: *"…guide the user to identify cognitive distortions and challenge them. Use examples or exercises when helpful."* We can even provide a structured approach (for instance, a common CBT dialogue flow: validate feelings, identify thought, examine evidence, suggest a reframe). By encoding these steps in the instructions, the LLM is more likely to produce a response that not only gives info but follows a therapeutic process. The retrieved knowledge then becomes supporting material for these steps (e.g., the evidence or explanation the bot can share).

- **Example Alignment:** Suppose a user says: *"I made a mistake at work, I'm sure I'll be fired. I can't do anything right."* The pipeline might retrieve a snippet on "catastrophizing" (a CBT term) from an article, and maybe an exercise about challenging negative predictions. The therapist persona prompt already inclines the model to do CBT, so the final answer might look like: *"It sounds like you're really worried and seeing this mistake as a catastrophe. I wonder if this could be an example of 'catastrophizing,' a common thinking trap* [1] [12] *. That's when our mind jumps to the worst-case scenario even if it's not likely. Let's examine this: what evidence do you have that you'll be fired? Mistakes happen to everyone. Perhaps we can come up with a more balanced thought. Also, an exercise I suggest is writing*

*down the situation, your thought, and then a more realistic thought. For instance, 'I made a mistake, but I can learn from it; one mistake doesn't mean I'll lose my job.'"* In this response, the bold term might have come from the knowledge base, and the exercise suggestion from a retrieved therapeutic exercise. The model, guided by CBT principles, integrates them into a coherent helpful reply.

- **Continuous Improvement:** We will test the bot with various user inputs (covering different issues like anxiety, depression, etc.) to see how well the RAG+CBT combination works. If the bot gives too generic advice, maybe the knowledge retrieval isn't being used effectively — we might then refine the prompt to more strongly encourage using the provided info ("*In your answer, incorporate any useful information from the knowledge above.*"). If the bot is too dry or technical (just regurgitating textbook text), we might need to adjust by telling the model to *explain things in simple terms and conversationally*. The balance is key: factual grounding without losing empathetic tone.

- **Tool Use (Future):** In the long run, we could implement specialized tools or functions the LLM can call, such as a "cognitive distortion detector" or "mood tracker." In the CBT game example we saw, the developer had the LLM call a distortion detection tool to analyze journal entries [13] . For our MVP, we likely won't implement complex tools, but we can simulate a simpler version by instructing the model itself to identify distortions or emotions and then retrieve relevant info. If needed later, we could add a classification step (using a smaller model or a rule-based approach) that labels the user's issue (e.g. "distortion: mind-reading") and then use that label to retrieve a specific remedy.

By tightly coupling the **retrieval mechanism** with **CBT-focused prompts**, the chatbot should deliver responses that are not only **well-informed** but also **therapeutically meaningful**. This approach leverages the strengths of LLMs (natural dialogue, reasoning) with the reliability of evidence-based content from our knowledge base [3] . The CBT framework acts as a guiding lens through which all information is filtered and used, ensuring the conversation stays on proven therapeutic tracks rather than random advice.

## Fine-Tuning and Persona Behavior Tuning

While prompt engineering will be our primary method to achieve the desired persona behaviors and CBT alignment, we should also consider fine-tuning or prompt-tuning for improved performance and robustness:

- **Prompt Tuning (iterative refinement):** Initially, we will iteratively tweak the prompts based on testing. This manual "prompt tuning" can greatly improve quality. For example, we might find the therapist persona sometimes gives advice too directly – we could refine the prompt to say "ask the user questions instead of giving direct advice when possible" to encourage a Socratic style. Or if the funny persona isn't funny enough, we add a few example jokes in the prompt. This process is essentially a form of prompt optimization done by the developers.

- **Low-rank Fine-Tuning for Style:** If the budget allows and the base model supports it, we could fine-tune the LLM on a small corpus representative of each persona. For instance, compile example conversations of a therapist doing CBT, a friend offering support, etc., and fine-tune using those as supervised examples. According to recent research, fine-tuning even smaller models on domain-specific dialogues significantly boosts their performance on those tasks [12] . In one study, models fine-tuned on CBT dialogue data showed marked improvement in delivering accurate CBT techniques compared to just instruction-tuned models [12] . This indicates that, in the future, training

a custom model (or using techniques like LoRA adapters) could yield a chatbot highly specialized in therapy style and knowledge. For MVP, however, full fine-tuning of Gemini (which is likely not open for fine-tuning yet) isn't feasible. Instead, we might fine-tune an open-source model (7-13B parameters range) on therapy data as a backup or for offline use, but keep Gemini for the production due to its presumably superior quality out-of-the-box. Fine-tuning could also help the humorous persona (we could train on a dataset of empathetic humor responses).

- **Embedding Persona in System via Few-Shot:** If fine-tuning is not an option, another approach for persona consistency is few-shot exemplars. As mentioned, we can include a short transcript snippet for each persona in the prompt to "demonstrate" the desired behavior. For example, for the coach persona, include a mini dialogue where a user says *"I feel demotivated"* and the coach responds *"I hear you. Let's set one small goal..."* in a very motivational tone. This helps the model see what style to follow. We must be mindful of token cost – perhaps only one example per persona if needed. We can dynamically prepend an example when persona is selected.

- **Continuous Learning from Interactions:** Over time, if the chatbot is used by test users, we will gather transcripts. We can fine-tune the model or the prompts based on real feedback (where did it fail, where did it succeed?). For instance, if users say the friend persona responses feel fake, we adjust language; if the therapist persona missed opportunities to apply CBT in some cases, we strengthen those instructions or add an example.

- **Safety and Alignment Tuning:** We might also need to fine-tune or adjust the model for avoiding certain pitfalls (like being too lenient with suicidal talk or giving inappropriate jokes). If Gemini API provides system-level safety parameters, we'll use those. Otherwise, open-source fine-tuning could incorporate some reward modeling to penalize unsafe responses. That said, a robust solution likely involves an external moderation filter in production.

- **Parameter-Efficient Tuning:** Given cost constraints, if we pursue model tuning, we can use parameter-efficient methods (like LoRA or prompt tuning where a small vector is trained and prepended to prompts). This avoids retraining the full model and can be done on smaller hardware. An interesting idea could be to train a separate small model to generate the appropriate persona style given a draft response. However, that complexity is beyond MVP scope.

In summary, **Version 1 will rely on prompt engineering primarily**, which is flexible and zero-cost (besides tokens). We'll only consider fine-tuning if necessary to achieve naturalness or if planning for a custom model deployment to reduce API usage. The research literature supports that a well-tuned model on therapy data excels at therapeutic tasks [12] , so it's an avenue for future enhancements. For now, the combination of a strong base model (Gemini) and carefully crafted prompts, plus the grounded data, should give us acceptable results. We will keep evaluating and note areas where more training might be needed.

## Version 1 Milestone Plan

To manage the development, we'll break the implementation into milestones for Version 1 (MVP):

**Milestone 1: Project Setup and Base Integration**
*Timeline: Week 1-*

- Set up development environment, repository, and accounts/keys for required services (Google Cloud for Gemini API, Pinecone for vector DB). Confirm access to Gemini API using a test call (e.g., a simple prompt to ensure it returns responses). - Install necessary SDKs (e.g., Pinecone client, any Python libraries or JS libraries for API calls). - Skeleton backend service: create a basic server that can accept a request (a dummy user message) and responds with a static reply, just to ensure the piping from UI to backend to API works. - **Success Criteria:** Environment is ready, and a "hello world" query to Gemini API returns successfully. Pinecone connection established (create empty indexes).

**Milestone 2: Knowledge Base Construction (Initial)**
*Timeline: Week 2-*
- Collect a small but representative sample for the knowledge base. Perhaps ~20 documents (e.g., a few therapy transcripts, a couple of CBT articles, and some exercises). We can start with publicly accessible text (even Wikipedia sections on CBT, if needed, just to have content). - Write the preprocessing script: chunk the documents, generate embeddings (likely using an API call to get embeddings or a local model if available), and upsert into Pinecone. - Test the retrieval by manually querying Pinecone (e.g., use a sample query like "negative thinking about failure" and see if it returns a chunk about cognitive distortions). - **Success Criteria:** We have a working Pinecone knowledge index with content and can retrieve relevant snippets via a simple test script. Documentation exists on how to add more data going forward.

**Milestone 3: Basic Chatbot Logic with One Persona (Therapist)**
*Timeline: Week 3-*
- Implement the core conversation loop in the backend for a single persona (start with the therapist persona as default). This includes: - Composing the prompt with a fixed therapist system message and the user query. - Performing a Pinecone search for knowledge (skip memory for now, since no conversation history yet). - Inserting retrieved text into the prompt. - Calling Gemini API and returning the answer. - Build a minimal front-end to test this (could be a simple web page or even using an API tool like Postman). Input a user question and see the response. - Evaluate outputs for a few scenarios (maybe simulate a user saying "I feel anxious all the time"). Verify that the bot uses knowledge if appropriate. - **Success Criteria:** The chatbot can handle a single-turn Q&A with the therapist persona, producing sensible answers. The integration of Pinecone and Gemini in one pipeline is confirmed.

**Milestone 4: Long-Term Memory Integration**
*Timeline: Week 4-*
- Enable multi-turn conversation memory. Implement the Pinecone memory index usage: after each user message + bot response, store an embedding of that message (and possibly response or a summary). Use the user's ID to segregate memory. - On each new user query, perform a Pinecone query against the memory index to fetch any semantically relevant past conversations. This could be limited to top 1-3 results to avoid prompt bloat. - Merge that into the prompt as additional context (for example: "Earlier you told me: [summary of past issue]."). Ensure the model distinguishes between past context and current query (this can be handled by clearly labeling or by just including it in the system prompt as background info). - Test this by simulating a short conversation: e.g. User says "I have been feeling down this week." (bot responds). Later user says "Also, I fought with my sister." See if the bot can recall the user was feeling down if relevant. While this is basic, it proves the concept. - **Success Criteria:** The chatbot remembers information provided earlier in the conversation or in previous sessions (when using the same user ID). It can incorporate that memory into new responses, making the conversation feel continuous.

**Milestone 5: Multi-Persona Functionality**

*Timeline: Week 5-*

- Implement the persona switching mechanism. On the UI side, allow the user to select a persona (e.g., via a dropdown at conversation start). Pass this selection to the backend with each message (or store it in session). - In the backend, expand the Persona Prompt Manager to include all four persona prompts. Based on the selected persona, prepend the corresponding system prompt. - Adjust any style specifics in how we handle retrieval context for different personas. (E.g., a friend persona might not want to quote a scientific article directly – so we might instruct the model accordingly in that persona's prompt to "explain any reference in your own words". The therapist persona might directly mention a CBT term since that fits their role.) - Thoroughly test each persona with a few example inputs to ensure the style difference is evident. For instance, compare the responses to "I'm feeling stressed at work" under each persona and verify they align with expectations (therapist gives structured CBT advice, friend gives casual support, coach gives motivating action, funny tries humor). - **Success Criteria:** The system correctly produces different styles of responses based on persona choice, without mixing them up. Users can switch persona (likely by starting a new session with a different selection) and get a new style of interaction.

**Milestone 6: Crisis Handling and Safeguards**

*Timeline: Week 6-*

- Develop the crisis keyword list or detection logic. Possibly use a simple list for MVP (e.g., "suicide, kill myself, not want to live, etc.") and maybe check message sentiment intensity. We might integrate an open-source sentiment model if needed, but keywords likely suffice for now to catch the obvious cases. - Implement the flashcard help message: design it to be very clear and resource-focused. (Optionally, prepare messages for various levels: e.g., one for suicidal ideation, one for mention of abuse, etc. For MVP one generic emergency message is fine.) - Make sure that if a crisis phrase is detected in either user's message or the LLM's draft response, the final output is replaced or prefaced with the help card text. We can simply not show the LLM's normal response in this case, or show it after the card with a big warning not to rely on it. But likely, simplest is to not show the normal response at all if it's a high-risk situation. - Test this by inputting something like "I want to end it all, there's no reason to live." The expected behavior: The bot should *not* attempt a normal empathetic response only – it should show the crisis card message, perhaps alongside a gentle note. This test ensures we don't inadvertently produce something unsafe. (We know from CrisisTalk that mainstream models like Gemini will suggest contacting 988 in such cases [10], but we won't rely solely on the model – we enforce it.) - **Success Criteria:** Crisis messages trigger correctly and the user is presented with a helpful, resource-oriented message. The conversation either stops or is heavily guided towards seeking help. No egregious failures (like the bot continuing a conversation encouraging harm).

**Milestone 7: Frontend Polish and User Experience**

*Timeline: Week 7-*

- Refine the UI to clearly show persona (maybe use avatars or different colors for each persona's messages). This makes it engaging; e.g., a friendly avatar icon for the friend persona, a suited professional icon for therapist, etc. - Add instructions or disclaimers in the UI: e.g., "This chatbot is not a licensed therapist and is for support purposes only. If in crisis, please seek professional help." (Important legally and ethically to set user expectations.) - Ensure the UI can handle multi-turn chats smoothly (auto-scroll, loading indicator while waiting for reply, etc.). - Maybe implement minor features like the ability to **toggle personas within a session** (though that complicates continuity – possibly we ask the user to start a new session to change persona, to avoid confusion). - Cross-browser/cross-device testing to ensure it works on mobile, etc. - **Success Criteria:** The application is user-friendly and visually clear. Persona selection and chat flow are

intuitive. Basic styling and branding (if any) are applied. The MVP is essentially feature-complete by this point.

**Milestone 8: Testing and Iteration**
*Timeline: Week 8 (and ongoing)*
- Conduct more extensive testing with various simulated user inputs. Create test cases for common issues: negative self-talk, anxiety scenario, asking for advice on relationships, etc. Verify the bot's responses are helpful and utilize CBT (at least the therapist persona should). If some responses are off, log those and adjust prompts or retrieval accordingly. - Specifically test boundaries: e.g., what if user asks for medical advice or something out-of-scope? The bot should politely decline or redirect (we may need to add a line in system prompt: "Do not give medical or medication advice. If asked, encourage seeing a doctor."). Similarly test if user tries to get the bot to do something it shouldn't (like unethical things). Rely on base model's guardrails but add ours as needed. - Performance test: measure response time with RAG. Each query involves an embedding call + Pinecone query + LLM call. Ensure this is within acceptable latency (~ couple of seconds). If it's slow, consider optimizations (like reducing number of retrieved chunks, etc.). - Compile feedback from any pilot users if available and fix any glaring issues or misunderstandings. - **Success Criteria:** The system performs reliably across a range of scenarios, and we have confidence in launching the MVP.

At this point, Version 1 should be ready. The MVP will have: - Text chat interface with persona selection. - Multi-persona responses following CBT framework. - Knowledge-grounded answers via RAG (with a curated CBT knowledge base). - Long-term memory across chats using Pinecone. - Basic safety net for crises (static message). - (No deep integration with privacy or external services yet, which is acceptable for MVP as per requirements.)

We will document all configurations and provide usage instructions. From here, we can gather user data (with consent) to plan improvements.

## Additional Suggestions for Robustness and Long-Term Viability

Building a therapy chatbot is a sensitive endeavor. Beyond the initial implementation, here are additional recommendations to ensure the solution is robust, safe, and sustainable in the long run:

- **Thorough Testing and Quality Assurance:** Regularly test the chatbot with diverse user profiles and situations. It's crucial to include edge cases (like extremely depressed users, users with anger, etc.) to see how the bot handles them. Use evaluation metrics or have mental health professionals review some example conversations for quality. Continuously improving the prompt or knowledge base based on these tests will refine the chatbot's effectiveness.

- **Human Oversight and Ethical Guardrails:** Even with disclaimers, some users may rely on the chatbot as if it were a therapist. We should incorporate clear notices that it's an AI and not a substitute for professional help. In future versions, consider an **escalation pathway** – e.g., if the bot detects a very troubled user (not just crisis keywords, but say someone who has been very depressed over many sessions), perhaps prompt them gently to consider real therapy. We might even integrate a button for "Speak to a human counselor" (though connecting to actual therapists is beyond MVP, it's worth planning if this is intended as a real mental health tool). At the very least, ensure the bot **never claims to be human or licensed**. Transparency is key to user trust.

- **Data Privacy and Security:** While not required for MVP, any stored user data (conversation logs, Pinecone memory) should eventually be protected. This means:

- An option for users to delete their data or opt-out of memory storing.
- Encrypting sensitive data at rest. For example, one project encrypted journal entries before storing in a database [14] – we can do similar with user messages in Pinecone (perhaps storing only hashed or encoded representations if possible). Pinecone itself is a managed service that likely ensures data security, but we should still treat mental health data as highly sensitive.

- When we do address compliance, we'll look at regulations like HIPAA (if applicable) or GDPR for EU users. This might influence where we host data and what consents we gather.

- **Scaling and Cost Management:** Using the Gemini API and Pinecone will incur ongoing costs. To keep the service viable:

- Monitor usage patterns. If certain users have extremely long conversations that fill up context, consider implementing a summarization function to compress old chat history (freeing up Pinecone space and reducing retrieval noise).
- Possibly cache certain embeddings or responses. For example, if many users ask similar generic questions (like "What is CBT?"), the bot could serve a stored answer or at least the Pinecone retrieval doesn't need to run fresh every time if cached.
- Evaluate model options periodically. If open-source LLMs become sufficiently capable, migrating to an in-house model could cut API costs. (The fine-tuning route on a smaller model is one such plan – it could allow an on-prem deployment eventually.)

- Pinecone scaling: choose an index metric and size that's cost-effective. We can start with a smaller index (Pinecone has tiers) and scale as needed. Also periodically prune the vector DB (e.g., remove memory vectors older than X months if they become irrelevant, or compress them via summary vectors).

- **Continuous Knowledge Base Updates:** Mental health research evolves, and also the range of user queries might expand. We should plan to update the knowledge base with new content (new therapeutic techniques, or if we notice queries about topics we lack info on – e.g., if many users ask about a specific condition or a new popular self-help trend). A robust process could be: maintain a content pipeline where new documents can be added, embedded, and inserted without downtime. Pinecone's live index update feature will help [15] .

- **User Personalization:** In the long term, beyond just remembering past chats, the bot could tailor its approach to the individual. For example, if a user consistently responds better to a certain style (maybe they prefer the friend persona's approach), we could note that. Or store user's goals/ progress (with consent) to bring up later ("Last time you planned to take a walk daily – how did that go?"). This creates a more **coach-like continuous experience**. Achieving this requires robust memory and perhaps explicit user profiles.

- **Robust Error Handling:** Ensure the system fails gracefully. If the LLM API is down or responds with an error, the backend should catch that and inform the user that the service is unavailable temporarily, rather than just silence. Similarly for Pinecone errors. Maybe have a fallback response

for when knowledge retrieval fails ("I'm here with you, though I'm having a bit of trouble accessing my notes right now." – some graceful degradation).

- **Monitoring and Analytics:** Track key metrics such as: average session length, most common user issues (via simple keyword tagging), how often crisis message was triggered, etc. This data (anonymized) can guide where to focus improvements. For example, if many users ask for "medication advice" which we disallow, maybe we need to adjust the prompt to handle that question better (like always provide a respectful deferral).

- **Community and Expert Feedback:** If possible, involve a mental health professional in reviewing the chatbot's responses periodically. They can catch subtle issues (e.g., tone that might be problematic, or advice that, while not wrong, might not be optimal). This will increase the credibility and safety of the chatbot's guidance. Also, user feedback is crucial – include an easy way for users to flag a response as problematic or unhelpful. This feedback loop will let us refine the system promptly if something goes awry.

- **Future Features:** Looking beyond MVP, there are interesting enhancements:

- **Multi-modal input:** maybe allow users to journal in longer form (which the bot could analyze, as in some CBT apps) or even voice notes (speech-to-text to feed into the system). Our architecture mostly remains the same, just an extra input processing layer.
- **Gamification:** The referenced CBT game example shows a novel way to engage users [16] . While we may not build a game, we could include elements like "achievements" for doing exercises, or a mood tracker that shows progress, to motivate users.
- **Plugin Tools:** Incorporate external tools when safe – e.g. a breathing exercise video link when user is panicking, or integration with calendar to schedule activities the coach persona suggests.
- **Group support persona?** Possibly a future persona that behaves like a support group facilitator, though that complicates single-user chat.

By implementing these suggestions over time, the therapy chatbot can remain **effective, safe, and engaging**. The combination of a well-structured RAG architecture and CBT grounding sets a strong foundation. With careful iteration and ethical oversight, this solution can provide accessible support to many who need someone to talk to. It's important we proceed thoughtfully: always prioritizing user well-being over technological gimmicks. In summary, our plan delivers a version 1 that is comprehensive yet expandable, balancing immediate functionality with a path for future growth and refinement.

**Sources:**

- Wasay Ali, *"Building a RAG Chatbot: LangChain, Pinecone, and Gemini," Medium*, Aug. 30, 2024. (Describes using vector DB and LLM for domain-specific chatbot) [1] [17]

- Pinecone Documentation, *"Chatbots with Pinecone,"* (Explains how vector databases provide long-term memory and grounded context for chatbots) [3] [5]

- Syamsul Zaman, *"Building a Therapeutic RAG-Enhanced AI Agent for CBT Gamification," Medium*, Jun. 10, 2025. (Illustrates a CBT-based AI agent architecture and prompt usage) [8] [14]

- *#CrisisTalk* Magazine, *"The Generative AI Therapy Chatbot Will See You Now,"* Jun. 2025. (Discusses safety approaches for AI mental health chatbots and how they handle crisis language) [10]

- Aseem Srivastava, *"5 Must-Know Mental Health Counseling Datasets for AI Research," Medium*, Jun. 6, 2023. (Lists datasets like HOPE with therapy transcripts useful for knowledge base) [11]

- Talha Tahir et al., *"Fine-Tuning Large Language Models to Deliver CBT for Depression,"* arXiv preprint 2412.00251, Nov. 2024. (Study showing that CBT-specific fine-tuning improved LLM therapy skills) [12]

---

[1] [9] [17] Building a RAG Chatbot: LangChain, Pinecone, and Gemini | by Wasay Ali | Medium

https://medium.com/@wasay.abbs/building-a-rag-chatbot-langchain-pinecone-and-gemini-d6f6b4be1015

[2] [3] [4] [5] [6] [7] [15] Chatbots with Pinecone | Pinecone

https://www.pinecone.io/learn/chatbots-with-pinecone/

[8] [13] [14] [16] Building a Therapeutic RAG-Enhanced AI Agent for Cognitive Behavioural Therapy (CBT) Gamification | by Syamsul Zaman | Jun, 2025 | Medium

https://medium.com/@thesyamsulzaman/building-a-therapeutic-rag-enhanced-ai-agent-for-cognitive-behavioural-therapy-cbt-gamification-d5247b48d7fc

[10] The Generative AI Therapy Chatbot Will See You Now - #CrisisTalk

https://talk.crisisnow.com/the-generative-ai-therapy-chatbot-will-see-you-now/

[11] 5 Must-Know Mental Health Datasets for AI Research | August 2023 | by Aseem Srivastava | Medium

https://as3eem.medium.com/5-must-know-mental-health-counseling-datasets-for-ai-research-dd1a1b9f30b4

[12] arxiv.org

https://arxiv.org/pdf/2412.00251