

# FutureFrame: AI-Driven Assistance for Human-Design Consulting

Alina Hyk & Samuel Jamieson

## Project Artifacts

### GitHub:

[https://github.com/AlinaHyk/FutureFrame\\_AI-Driven\\_Assistance\\_for\\_HumanDesign\\_Consulting.git](https://github.com/AlinaHyk/FutureFrame_AI-Driven_Assistance_for_HumanDesign_Consulting.git)

### Trello board:

<https://trello.com/invite/b/679c0c93930a271a21437a58/ATTI9bbc25b631b93de815ffdf983fc7ba8c32BA2C96/cs362>

### Live web app:

<https://futureframeai-drivenassistanceforhumandesignconsulting-gvuhiel.streamlit.app/>

**Reflections are at the very bottom; after appendix!**

## Abstract

Human Design is a relatively new integrative system that blends concepts from psychology, astrology, the I Ching, Kabbalah, and other disciplines to provide individualized insights into personal growth, decision-making, and interpersonal dynamics. Although it has gained traction among life coaches, therapists, and curious individuals seeking self-discovery tools, the field's recent emergence and limited pool of experts make it difficult for newcomers and practitioners alike to access high-quality, credible resources. Existing materials are scattered across multiple formats—including private lecture recordings, annotated slides, books, and notes—many of which are behind

paywalls or simply inaccessible.

This project addresses these challenges by developing a specialized chatbot that consolidates and streamlines Human Design knowledge into a single, interactive interface. By curating and embedding a broad range of domain-specific documents—including private textbooks and lecture transcripts—this tool provides two main benefits. First, it offers 24/7 assistance to clients who use Human Design consulting services, allowing them to explore or revisit key ideas at their own pace. Clients can prepare for upcoming consultations, clarify details discussed in previous sessions, and deepen their engagement without waiting for the next live meeting. Second, the system equips Human Design consultants with a centralized knowledge base, ensuring that they can retrieve and reference crucial information quickly and accurately during sessions. This approach saves time, enhances the quality of guidance, and enables professionals to focus on personalized interactions with their clients.

In addition, by integrating broader psychological insights already present in advanced large language models, the chatbot can provide a more holistic perspective when addressing user queries. Rather than replacing expert consultants, the goal is to support them by handling frequently asked questions, aggregating relevant teachings from disjointed sources, and offering contextually rich explanations. By making trustworthy Human Design information accessible and navigable, this project aims to expand the reach of a niche but valuable framework for self-discovery and interpersonal understanding.

## **Goal**

This project's overarching objective is to streamline and enhance the consultation and learning process within a Human Design-focused private firm by providing a centralized, intelligent assistant. The primary goals include:

### **1. Comprehensive Knowledge Repository**

- Aggregate and unify a wide variety of Human Design materials—books, lectures, notes, as well as audio/video recordings transcribed into text—under a single virtual assistant platform.
- Host these data in a secure, scalable environment where they can be easily maintained, updated and assessed.

### **2. Accurate Information Retrieval**

- Implement a robust RAG pipeline that uses embeddings-based similarity search to return the most relevant segments of Human Design material whenever a user poses a question or scenario.
  - Provide traceable references to the source documents, ensuring trustworthiness and reducing misinformation.
3. **Contextual, Multi-Disciplinary Insights**
- Blend psychological and therapeutic principles (e.g., existential therapy, depth psychology) into responses, showing how Human Design concepts complement or diverge from standard psychological practices (This will be achieved through prompt engineering by prompting the model to combine its existing knowledge and expertise in psychology with information on Human Design that can be retrieved from the data corpus used to train the RAG system).
  - Facilitate an environment where these combined insights can aid in personal development, self-awareness, or professional consultation scenarios (also achieved through appropriate prompt engineering .
4. **User-Centric Design**
- Provide a simple, intuitive interface that consultants and clients can use without specialized technical knowledge.
  - Offer a smooth user journey, from query input to final output, ensuring minimal friction and maximum clarity.
5. **Scalable & Extensible Architecture**
- Employ cloud-based services to handle potentially large volumes of queries or expansions in data size.
  - Maintain a flexible design so that new content types (like diagrams, images, or additional languages) can be added with minimal structural changes.

## Current Practices

### Software and Platforms for Human Design Analysis

Human Design consultants rely on specialized software to generate and interpret clients' **BodyGraph** charts. The official source, **Jovian Archive**, provides two primary tools: the classic **Maia Mechanics Imaging (MMI)** desktop app and the new **Maia Mechanics Online** web platform ([jovianarchive.com](https://jovianarchive.com)). These offer advanced chart calculations (individual, relationship, transit, etc.) and are tailored for professional use. For example, the Maia Mechanics online tool supports multi-language charts, automated

transit charts, and even a birth time accuracy scanner to flag chart changes around uncertain birth times ([jaccunningham.com](http://jaccunningham.com)). Jovian's free web tool **MyBodyGraph** serves as an entry-level platform connected to the same database – it makes it easy to create basic charts and even **build a collection of BodyGraphs for clients or family**, with options to purchase deeper detail reports as needed ([jaccunningham.com](http://jaccunningham.com)).

In addition to Jovian's tools, a number of third-party applications have gained popularity among consultants:

- **Genetic Matrix** – A cloud and mobile-based software suite offering one of the largest varieties of Human Design and astrology charts ([geneticmatrix.com](http://geneticmatrix.com)). It's valued for calculating advanced data like **Variables** (e.g. Determination, Environment, Perspective, Motivation) and providing multiple systems (including tropical vs. sidereal charts) ([jaccunningham.com](http://jaccunningham.com)). Genetic Matrix is considered **intermediate** in complexity (better for seasoned students) and includes features like an easy-to-read info list of chart details and even a library of celebrity charts for reference.
- **My Human Design (Jenna Zoe)** – A modern web and mobile app created by a prominent HD influencer. It provides an extremely user-friendly interface for newcomers, with simplified language and quick lookup of key aspects like Type, Profile, Cognition, and Variables ([jaccunningham.com](http://jaccunningham.com)). However, it intentionally uses renamed concepts and “mutated” terminology that differ from the classical lexicon ([jaccunningham.com](http://jaccunningham.com)), which can lead to confusion if mixed with traditional materials.
- **BodyGraph Chart** (BodyGraphChart.com) – A platform and widget for practitioners to integrate chart tools into their own websites ([jaccunningham.com](http://jaccunningham.com)). It allows consultants to **embed a Human Design chart calculator** on a site (WordPress, Wix, etc.) ([woobro.com](http://woobro.com)) and even generate automated reports. Notably, it offers an “**Automated Human Design Reports**” feature that lets consultants create and sell customized written reports effortlessly ([jaccunningham.com](http://jaccunningham.com)). This kind of tool helps practitioners scale their services by providing clients with computer-generated chart breakdowns as a supplement to live readings.
- **BG5 Business Software** – A specialized tool for Career and Business Design analysis, aligned with the BG5 system (an offshoot of Human Design for business). It enables certified **BG5 consultants** to easily create and analyze career-oriented charts and graphs ([bg5businessinstitute.com](http://bg5businessinstitute.com)). The software integrates the vast amount of information relevant to a person's professional design in an intuitive interface ([bg5businessinstitute.com](http://bg5businessinstitute.com)).
- **Niche and Mobile Apps** – A variety of smaller tools cater to specific needs. For example, *Pure Generators* (by Rachel Lieberman) is a colorful free web tool geared towards Generator-type clients, including a custom **quick-start guide for those with Sacral definition** ([jaccunningham.com](http://jaccunningham.com)). It shows all basic chart info (including Variables) and is ideal for Generator/MG clients to get tailored advice ([jaccunningham.com](http://jaccunningham.com)). Another example is the **Neutrino Design App**, a mobile app (iOS/Android) which is handy for tracking daily transits and lunar cycles – popular among Reflectors and those monitoring shifts over a month ([jaccunningham.com](http://jaccunningham.com)). Neutrino allows saving multiple charts and has substantial free content, making it a useful pocket tool for consultants and enthusiasts alike. These niche apps reflect a

trend of segmenting tools by audience or use-case (e.g. type-specific guidance, portable transit monitoring), adding to the ecosystem of resources available.

## Information Access and Organization in Practice

Human Design consultants typically handle large amounts of detailed information for each client – from the basic chart components (Type, Strategy, Authority, Profile) down to specific Gates, Channels, and planetary placements. Modern software platforms greatly assist in **organizing and accessing this information** during analysis. For instance, most tools like MyBodyGraph or Maia Mechanics allow users to save multiple charts in a personal library and load them on demand, which is essential for tracking client histories. A consultant might keep a library of client charts within these applications, often tagged with names or notes, enabling quick retrieval for follow-up sessions. The ability to generate **connection (composite) charts** is another important feature – software can overlay two individuals’ designs to reveal relationship dynamics, a common service in Human Design relationship readings. Similarly, transit and cycle charts (e.g. Saturn return, solar return) can be automatically produced by these tools to help consultants prepare context for timing and life cycle analysis. To utilize the information effectively, many consultants leverage built-in reference content in these platforms. **Genetic Matrix**, for example, presents information in a very visual and list-oriented way that makes it easy to find specific details that otherwise “*can be challenging to find without certain software,*” as one reviewer noted ([jaccunningham.com](http://jaccunningham.com)).

Many professionals also create their own **knowledge management systems** outside of the charting software. This can include personal notes, keyword “cheat sheets” for quick recall of keynotes (traits associated with each Gate, Channel, etc.), or using general tools (like Evernote/Notion) to catalog insights from past readings. However, the trend is that newer platforms are incorporating more of these needs internally. As mentioned, some services auto-generate written reports – consultants can use these as a starting template and then personalize or elaborate on them. The BodyGraphChart system’s **automated reports** feature is explicitly marketed to streamline operations and client engagement ([jaccunningham.com](http://jaccunningham.com)). In practice, a consultant might generate a comprehensive PDF report for a client’s design, then discuss highlights in a live session. This automation frees time while ensuring no key detail is overlooked, since the report will cover all centers, channels, etc., by default.

Another aspect of information management is **multi-modal learning resources**. Consultants often deepen their expertise by studying authoritative texts or courses and then integrating that knowledge into practice. For example, they may reference *The Definitive Book of Human Design* or course materials when preparing a difficult chart. Some software (like 64keys) even provides integrated training libraries – 64keys had developed a “most comprehensive Human Design training in video form” (particularly in German) by 2019 for its users ([64keys.com](http://64keys.com)). This blurs the line between a tool and an educational resource: the platform itself teaches the consultant via built-in media, while also providing the analysis functions.

In summary, today’s consultants utilize a blend of **purpose-built HD software** (for calculation, visualization, and record-keeping) and **supplementary resources** (both within those tools and externally)

to access and organize the vast information in a Human Design chart. The software acts as a real-time reference guide and data manager during client sessions, letting the practitioner focus on synthesis and guidance rather than manual lookup of chart data.

## Client Interaction and Emerging Market Platforms

Traditionally, Human Design consulting has been a one-on-one service delivered either in person or via online meeting platforms (Zoom being particularly common). Consultants manage appointment bookings and client communications through general-purpose tools – email, calendar schedulers like Calendly, and so on. However, as the field grows, there are now **dedicated platforms** designed to facilitate client interaction specifically for Human Design professionals.

One notable development is **HumanDesign.ai**, a UK-based AI-driven platform launched in recent years. While primarily an AI analysis tool (discussed more in the second section), it doubles as a professional hub for consultants. It offers a suite of features such as a **professional profile showcase** (where a consultant can list their qualifications and testimonials) ([humandesign.ai](https://humandesign.ai)). It essentially functions as a marketplace: clients using the platform can find and connect with listed Human Design experts, facilitated by the platform's understanding of what the client might be looking for. In short, HumanDesign.ai not only generates chart insights but also helps match users with coaches/analysts, streamlining client acquisition for practitioners ([humandesign.ai](https://humandesign.ai))

Beyond specialized platforms, many consultants interact with clients through **social media and community forums**. Another trend is the creation of **blended services and communities**. For example, consultants might run private membership groups (on platforms like Mighty Networks or Patreon) where they regularly post transit updates or design tips, and members can ask questions about their charts. This community-oriented model is facilitated by tools but driven by the consultant's expertise. It's a move away from just one-off readings toward ongoing interactive guidance.

Finally, **language localization and global reach** are improving in client interactions. Historically, Human Design material was English-centric (stemming from the original teachings), but today's tools support multiple languages – Maia Mechanics online is available in over 10 languages ([jovianarchive.com](https://jovianarchive.com)), and HumanDesign.ai supports 42 languages for chart insights ([trendhunter.com](https://trendhunter.com)). This allows consultants to work with clients across continents more easily, each viewing charts or reports in their preferred language. It opens up a truly international client base and has led to cross-cultural growth of Human Design coaching.

## Training and Operational Trends for Consultants

The pathway to becoming a Human Design consultant has expanded and diversified. In the past, most professionals were formally trained and certified through the **International Human Design School (IHDS)** or licensed national organizations, following a rigorous multi-level curriculum set by the founder Ra Uru Hu's standards. That route still exists (IHDS offers analyst and teacher certification programs), but

now there are alternative educational avenues and a surge of independent teachers. For instance, some well-known figures have created their own courses: Karen Curry Parker's **Quantum Human Design** certification (which introduces a new vocabulary for HD concepts) or Chetan Parkyn's training programs accompanied by his proprietary software and app ([evolutionaryhumandesign.com](http://evolutionaryhumandesign.com)). These alternatives often aim to make the learning process more approachable or tailored to specific applications (such as coaching, business, or therapy).

A notable trend is the **modernization of language and teaching styles**. The Quantum Human Design movement and Jenna Zoe's approach (as seen in the MyHumanDesign app) both exemplify an effort to **"change traditional language"** and present concepts in a more relatable, upbeat tone for contemporary audiences ([jaccunningham.com](http://jaccunningham.com)). For example, terms like "Generators" and "Projectors" might be renamed to "Alchemists" or "Advisors" in certain schools. While the core knowledge remains similar, this trend indicates consultants are learning in two slightly diverging dialects of Human Design now – the classic and the rebranded. New consultants often have to navigate these differences, especially if they consume content from social media where mixed terminology is present.

In terms of skill development, many consultants are now **blending modalities**. It's common to find Human Design practitioners who also study the Gene Keys, astrology, or coaching methodologies (like NLP or positive psychology) and incorporate those into their practice. This holistic learning trend means the consultant toolkit is broader, and there's a demand for tools and platforms that accommodate multi-disciplinary integration. From a training perspective, some institutes (like BG5 Business Institute for career consulting, or health-focused HD courses for well-being practitioners) address specific niches, so consultants pick tracks aligned with their intended niche.

Continuous learning is also supported by **online communities and resources**. There are active forums (e.g., the 64keys user group on Facebook [64keys.com](https://www.facebook.com/64keys)) and Reddit communities where consultants and enthusiasts discuss case studies, share chart interpretations, and crowdsource knowledge. The openness of information on the internet – including many free YouTube lectures, podcasts, and articles – has lowered the barrier to entry for learning Human Design. As a result, some consultants today are effectively self-taught or piecemeal-taught via these resources, rather than through a singular certification program. They operate as independent consultants based on reputation and results rather than formal credentials.

Operationally, the rise of **remote consulting** has been a huge trend. Especially after 2020, most Human Design readings shifted to Zoom or phone calls by necessity, and this has largely remained due to convenience. Consultants often provide recorded sessions (audio or video recordings of the reading) so clients can re-listen – a practice made easier with today's tech. Some also offer written follow-ups or PDF summaries. This shift has normalized things like clients being in different countries/time zones and has pushed consultants to refine their online communication skills (screen sharing a chart, using digital pointers or annotation to highlight parts of the BodyGraph during a call, etc.).

Lastly, a very recent development is the incorporation of **AI tools to assist consultants**. Where we are beginning to see practitioners using AI-driven aides – for example, an analyst might use an AI chatbot (trained on Human Design material) to quickly pull up obscure information ("What does Gate 34 line 5

signify again?”) during a session, or to generate first drafts of client reports which they then personalize. The HumanDesign.ai platform explicitly markets its AI “Bella” as a tool that even non-certified enthusiasts can use to “*read and interpret charts with confidence*,” effectively lowering the learning curve for newcomers ([humandesign.ai](https://humandesign.ai)). While some analysts may stick to traditional methods, the up-and-coming generation of consultants is comfortable leveraging technology, including AI, to enhance their practice. This segues into the next section, which explores how AI chatbots and retrieval systems are being designed – knowledge that is increasingly relevant to the Human Design field as it embraces these innovations.

## AI Chatbot Development: Current Practices and Relevance to Human Design

### Retrieval-Augmented Generation (RAG) and IR-Based Chatbot Architectures

Modern AI chatbots often utilize a technique called **Retrieval-Augmented Generation (RAG)** to improve the quality and accuracy of their responses. In essence, **RAG merges natural language generation with information retrieval** ([superannotate.com](https://superannotate.com)). Instead of relying solely on a pre-trained language model’s memory, the chatbot first *retrieves relevant contextual information* from an external knowledge source and provides that to the model along with the user’s query. The language model (such as GPT-3.5 or GPT-4) then generates an answer using both its trained knowledge and the retrieved context. This approach is invaluable when up-to-date or specialized information is needed – scenarios where the base model might not have learned those details during training. An **IR-based architecture** for a chatbot follows a similar principle, treating the task as an open-book exam: the chatbot has a “lookup” step (information retrieval) and a “compose answer” step (generation). For example, **Bing Chat** (Microsoft’s GPT-4 powered search assistant) uses a form of RAG. It performs live web searches and feeds the results into an orchestration system with GPT-4 to generate answers grounded in current information. This means when a user asks Bing Chat something factual or timely, the system retrieves relevant webpages, and the GPT model then crafts a response citing that content – greatly reducing hallucinations and increasing accuracy.

The typical components of an IR-augmented chatbot include: a **document store** or knowledge base (which could be a set of webpages, a database of documents, or even a vector database of embeddings), a **retriever** mechanism to find the most relevant pieces of information for a given query, and the **generator** (the large language model) that takes the question + retrieved text to produce a final answer ([help.openai.com](https://help.openai.com); [superannotate.com](https://superannotate.com)). This architecture has become a standard for any domain-specific chatbot, from customer support agents to scholarly assistants. It addresses one of the key weaknesses of standalone LLMs: the tendency to **hallucinate** or confidently fabricate information when they lack knowledge. By “*bridging the gap between generating text and using real-world knowledge*,” RAG ensures the chatbot can reference an authoritative source rather than guess ([superannotate.com](https://superannotate.com)). The retrieved data effectively grounds the model’s response in facts.



In practice, implementing RAG often involves **semantic search** with embeddings. Instead of keyword matching, the chatbot converts the user's query into a vector and finds conceptually similar content in a vector database ([help.openai.com](https://help.openai.com)). This allows it to fetch information that may not contain the exact query words but is contextually relevant (important for natural language questions). The retrieved text snippets are then appended to the prompt for the LLM. For example, if building a chatbot to answer questions about the Human Design System, one would feed it a curated knowledge base (definitions of all Types, Gates, etc., possibly chunks of official texts). When asked "What does having Gate 57 in my design mean?", the system would retrieve the section about Gate 57 from the knowledge base and include it when formulating the answer, ensuring accuracy and richness of detail.

The RAG approach is widely adopted in **GPT-based chatbot frameworks**. Platforms like LangChain provide templates to easily connect a GPT model with a retriever (such as a Pinecone or FAISS vector index) and build a conversational chain that keeps context. Tech companies have also integrated this: OpenAI's own ChatGPT now allows a "knowledge retrieval" plugin or feature for custom deployments, which essentially automates RAG for a given set of documents ([help.openai.com](https://help.openai.com)). Similarly, many enterprise chatbot solutions use RAG to let the bot answer questions about company-specific policies, product docs, etc., which the base model wouldn't inherently know. In summary, **information retrieval (IR)** combined with **generation** is the backbone of current advanced chatbot design, enabling domain-specific expertise and up-to-date knowledge in conversational AI.

## GPT-Based Chatbot Methodologies and Applications

At the core of many AI chatbots today are **GPT models** (Generative Pre-trained Transformers), which excel at understanding prompts and generating human-like text. GPT-based chatbot methodologies revolve around a few key strategies:

- **Prompt Engineering and System Instructions:** Since GPT models can be steered by the prompt, developers craft specific system messages or few-shot examples to give the chatbot a persona or style. For instance, one might instruct the model "You are a professional life coach" or provide dialogue examples to shape how it responds. This method, often used in custom ChatGPT "personas", can produce specialized behavior without changing the model's weights. In the context of Human Design, one could prompt-engineer GPT-4 to act as a "Human Design analyst" by feeding it definitions of the chart and asking it to respond with that perspective.
- **Fine-tuning:** In some cases, a GPT model is fine-tuned on domain-specific Q&A pairs or dialogues. For example, if we had a large dataset of Human Design reading transcripts, we could fine-tune a model like GPT-3 on that, so it learns the style and content. However, fine-tuning is less common with the largest GPT-4-scale models (which are often used as-is with retrieval, due to the cost and complexity of fine-tuning).
- **Memory and Conversation State:** GPT chatbots typically maintain a conversation history to allow for context carry-over. The model input each round can include the last several turns of dialogue, enabling the bot to "remember" what's been discussed. Some advanced systems also implement *long-term memory* beyond the immediate window, such as storing user preferences or facts in a database that gets re-injected when relevant. For example, an AI life coach bot might

save the user's stated goals and personal details separately, and each session fetch those facts to remind the GPT model ("User's goal: exercise 3 times a week; user struggles with motivation"). The **Summit AI Life Coach** is an illustrative case: it lets users personalize the coach's style (e.g. tough love drill sergeant vs. empathetic friend) and the coach "*remembers what you tell it,*" leveraging stored data to simulate long-term memory ([summit.im](https://summit.im)). This dramatically improves the relevance of coaching advice over time.

- **Tool Use and Multi-step Reasoning:** Some chatbot frameworks allow the model to invoke external tools – for example, calculators, calendars, or search – if needed (this is seen in systems like OpenAI's Plugins or the ReAct paradigm). A GPT-based assistant might decide to lookup today's transits via an API if asked about "today's Human Design transits," combining retrieval with real-time tool use. While not always necessary for Human Design (since charts can be precomputed), this methodology is part of the state-of-the-art for making chatbots more capable.

For **Human Design applications**, GPT-based methodologies are particularly relevant because the knowledge is complex and nuanced. A GPT-4 chatbot given the entire corpus of Human Design knowledge (either via training or retrieval) can interactively explain a person's chart, answer follow-up questions, and tailor the explanation to the user's level of understanding. The generative ability of GPT means it can articulate insights in a conversational, personalized manner – similar to how a human consultant would phrase things differently for a novice vs. an advanced client. Moreover, GPT can handle **open-ended questions** ("What career strengths do I have according to my design?") by synthesizing multiple data points (Type, defined centers, profile lines relevant to career, etc.), something that static reports or lookup tables alone can't easily do.

One challenge GPT chatbots must manage is **factual consistency**. In a domain like Human Design, it's critical that the chatbot's output stays true to the established meanings (e.g., not swapping Gate numbers or inventing strategies). That's why the RAG approach is so useful here – by retrieving the official descriptions or expert-written interpretations, the GPT model's answers will align with authoritative sources. If a user asks about an obscure aspect like "What does Color 4 in my Motivation mean?", a properly designed bot would retrieve the relevant paragraph from a Primary Health System manual and then have GPT explain it in simpler terms. Without retrieval, even a powerful model might hallucinate an answer because such niche info might not have been in its training data.

Finally, it's worth noting the **evaluation of GPT-based chatbots**: Many are now built with some form of feedback or rating loop to refine their performance. In live systems, developers monitor outputs for errors or user dissatisfaction and adjust prompts or data accordingly. The iterative nature of prompt engineering means these bots improve over time, either through formal updates or informal learning (some systems fine-tune on conversations to better align with user needs, within the bounds of privacy policies).

In summary, GPT-based chatbot methodology combines the strengths of large language models (fluent natural language generation and reasoning) with strategies to ground them in the correct information (via retrieval, prompt instructions, and memory). The result is a conversational agent that can provide **reliable**,

**context-specific, and user-aware responses** – an approach highly applicable to something like a “virtual Human Design consultant.”

## AI Chatbots for Personalized Guidance and Human Design Use Cases

A number of AI chatbot implementations in recent years mirror the role of a personal coach, advisor, or analyst – roles very akin to a Human Design consultant’s function of providing personalized guidance. These case studies demonstrate how the technologies discussed are applied in practice, and they foreshadow how Human Design consulting can be augmented or complemented by AI:

- **HumanDesign.ai’s “Bella”:** This platform is directly aligned with Human Design consulting. *Bella* is the AI assistant that users (and even coaches) can interact with to get chart insights. A user can “ask *Bella* any question” about their chart or even someone else’s chart and get a tailored answer ([humandesign.ai](https://humandesign.ai)). For example, one might ask “What does it mean to be a 5/1 Manifestor with Emotional Authority?” and *Bella* will generate a detailed yet digestible explanation drawing on the relevant parts of the design. Under the hood, *Bella* likely uses a combination of a Human Design knowledge base (to retrieve specific keynotes or descriptions) and a GPT-like model to formulate the response. HumanDesign.ai advertises this as an “*Intelligent AI Bodygraph Tool*” that provides personalized insights on demand. For professionals, this means they can use *Bella* to augment readings or check interpretations. The platform’s success (17k users as of early 2025) shows receptiveness to AI-driven HD analysis ([markets.businessinsider.com](https://markets.businessinsider.com)). It essentially democratizes some of the analysis process: users can get instant answers instead of waiting for a live session, and budding consultants can learn by observing the AI’s output as a model.
- **Astrology and Coaching Chatbots:** Parallel to Human Design, astrology has seen AI integration through projects like “*Astrology Works*”, which introduced **Astra**, a personal AI astrologer. *Astra* offers “instant answers and guidance on various aspects of life, including self-development, career, relationships, and health,” all through an astrology lens ([trendhunter.com](https://trendhunter.com)). It uses AI to interpret one’s birth chart and provide advice, demonstrating that ancient esoteric systems can be translated by machine learning into conversational guidance. The fusion of AI with astrology was noted to “leverage machine learning algorithms to interpret astrological data and provide personalized insights,” increasing both accuracy and accessibility of astrology readings. This is highly analogous to what an AI in Human Design would do. Another example is **Astrology Birth Chart GPT**, a custom GPT-4 based bot that emerged via OpenAI’s tools, allowing users to input birth details and ask questions to an “expert astrologer” GPT ([blog.hubspot.com](https://blog.hubspot.com)). Users found that it could produce natal chart interpretations quite convincingly, albeit with the occasional mistake that a seasoned astrologer might catch. These examples indicate that people are already using GPT-based bots for personalized metaphysical guidance, which reinforces the viability of similar approaches for Human Design.
- **AI Life Coaches and Therapy Bots:** Outside of esoterics, there are numerous AI chatbots aimed at personal development and mental well-being. **Replika** is a well-known AI friend app that many use for emotional support; **Woebot** and **Wysa** are therapy-oriented chatbots using cognitive-behavioral techniques. More directly in coaching, **Summit: AI Life Coach** (as mentioned) provides goal tracking and motivation. Summit’s AI coach is available 24/7, holds

users accountable to their goals, and can even adopt different coaching styles to suit the user ([summit.im](https://summit.im)). It highlights the advantage of AI in personalized support: consistency and customization. One testimonial noted that using Summit's coach was superior to a generic chatbot because, unlike vanilla ChatGPT, it could reference the specific context of the user's goals and break them down into steps. This is achieved by integrating the user's data (goals, progress) into the AI's knowledge. Translating this to Human Design, an AI could remember a user's design details and life events discussed, enabling a continuity across sessions that feels very personal. For instance, an AI coach that knows a user is a Manifesting Generator might continually remind them to respond to their gut feelings when new opportunities arise, reinforcing strategy in day-to-day scenarios – essentially “living your design” coaching.

- **Knowledge Retrieval Bots in Enterprises:** On a different note, many companies deploy GPT-based chatbots internally for **knowledge retrieval**, such as answering employee questions about HR policies or helping engineers find documentation. While not as personal, these are relevant because they demonstrate robust use of RAG at scale. These systems show that given a defined body of knowledge (like a company handbook or, analogously, the *Body of Knowledge of Human Design*), a chatbot can reliably answer pointed questions. The success criterion is accuracy and consistency with the source material, which the retrieval mechanism ensures. The learnings from enterprise Q&A bots can inform the development of HD bots: for example, ensuring that the **voice of the chatbot matches the authoritative texts** (some solutions will even cite the source material in the responses, much as we are doing here). A Human Design chatbot might similarly cite which official lecture or book a particular insight came from, if needed, to build trust with the user that it's not just “making things up.”

**Industry-leading tools and methodologies** for building these chatbots often involve frameworks like the aforementioned LangChain, vector databases like Pinecone or Weaviate for efficient semantic search, and pre-trained LLMs from OpenAI or open-source alternatives (like LLaMA-based models).

In conclusion, the intersection of AI chatbots with the field of Human Design is a frontier that's quickly becoming reality. By using **retrieval-augmented GPT architectures**, an AI can effectively serve as a knowledgeable, interactive **Human Design guide** – whether as a direct-to-consumer “digital reader” or a smart assistant that professionals use behind the scenes. Key technologies like RAG ensure the chatbot stays accurate to the esoteric knowledge, and advancements in GPT conversational abilities ensure the guidance feels nuanced and human-like. We are witnessing the early case studies (HumanDesign.ai's Bella, Astra the astrologer, Summit coach, etc.) that validate this approach. As these technologies mature, it's likely that **AI-driven coaching and analysis tools will become standard companions in the Human Design consulting space**, enhancing how consultants learn (by providing on-the-spot answers and brainstorming interpretations) and how they operate (by scaling their reach and offering new interactive services to clients).

## Novelty

Existing chatbots that serve specialized domains typically rely on shallow keyword matching or ungrounded large language models, resulting in fragile answers that can misinform users. Our approach

integrates a retrieval-augmented pipeline capable of pinpointing the precise segments in a corpus and feeding that data into a transformer model. This hybrid design ensures more stable, verifiable outputs because every generated response stays tied to a known source, reducing guesswork and preventing hollow speculation. Human Design, due to its complex mix of psychology, esoterics, and individualized user profiles, exposes the weak spots of ordinary question-answer systems. Many standard bots ignore nuance in personal authorities, overshadow the interplay between gates, and fail to account for the intricacies of variable-based analyses. By grounding answers in curated text snippets, we maintain a high-fidelity link to original materials rather than simplifying them into generic bullet points.

We face tough constraints around context length, rapidly changing content, and user-specific queries that demand contextual memory. Handling context length means we must carefully throttle the amount of text we inject into a model's prompt to keep within token limits. Tracking updated or newly ingested data involves re-embedding any additions to the corpus and re-indexing them, all without retraining the entire system. Our architecture supports partial updates quickly, letting it scale efficiently across hundreds of lectures, books, and user-provided notes. Another constraint arises from user preferences around interpretive depth and personalized references. Different consultants or clients might ask for either a concise summary or a deep cross-examination of their chart's centers, channels, or psychological overlays. To address this, we permit dynamic retrieval that ranks chunk relevance and merges them in the final prompt, giving the chatbot a chance to answer exactly at the depth requested while citing each source. We also confront the risk of hallucination whenever the language model tries to fill gaps it cannot find in the embedded texts. Our best guardrail is to include a direct instruction in the prompts: if the retrieved context does not support the claim, the model should respond with a polite disclaimer. By referencing exact paragraphs and requiring the model to quote relevant lines, we curb the creation of unverified statements.

Performance under real-time conditions means we must blend minimal overhead for semantic search with the more expensive calls to large language models. We distribute the workload by caching embeddings and limiting the number of retrieved chunks to the top five or ten per query. This prevents an excessive cost blowout while still providing enough context to shape authoritative, reasoned answers that stand up to scrutiny. Security concerns also factor heavily into our design. Storing private notes from paid lectures or personal sessions requires strong access controls on the database and encryption of sensitive data. We segment public domain references from proprietary content, ensuring the retrieval pipeline draws only from authorized subsets of the corpus. This approach respects licensing conditions while maintaining a uniform interface for the user.

By grounding the chatbot's logic in domain-specific documents, referencing relevant source text, and leveraging specialized retrieval strategies, we resolve the all-too-common pitfalls of purely generative models. The system stays comprehensive enough to traverse the varied threads of Human Design theory, yet precise enough to guide novices and veterans alike without resorting to empty generalities. The result is a holistic, verifiably accurate knowledge assistant that operates within tight computational and ethical constraints, making it both reliable and adaptable in the real world.

# Effects

## User Experience

Retrieval-augmented chatbots deliver more contextually relevant and informative responses, leading to a noticeably improved user experience. By querying external knowledge bases in real time, the system can provide answers that are precise and up-to-date rather than generic, which **enhances overall response quality and engagement** ([tonic.ai](#)). Users benefit from interactions that feel more insightful and personalized, as the chatbot can draw on domain-specific facts or recent information to address queries. Furthermore, retrieval augmentation fosters transparency: many such chatbots return cited sources or evidence along with answers, allowing users to verify the information provided ([ibm.com](#)). This capability to double-check facts in context increases **user trust**, as people gain confidence that the chatbot's guidance is backed by authoritative data rather than unsupported generative text. The combination of higher answer quality, active engagement through relevant content, and verifiable outputs makes retrieval-augmented systems far more trustworthy and satisfying to interact with than standard LLM-based bots.

## Accuracy

Incorporating retrieval mechanisms greatly **enhances the accuracy** and factual correctness of chatbot responses. Traditional large language models sometimes rely on incomplete training data or statistical patterns, which can lead to vague or incorrect answers. A retrieval-augmented approach mitigates this by grounding the generation process in pertinent external information. When a query is received, the system fetches relevant documents or knowledge snippets and uses them to inform the response. This ensures that answers are not solely drawn from the model's parameter memory but are explicitly supported by current data. As a result, the chatbot can provide **precise, contextually relevant answers with far fewer errors** ([tonic.ai](#)). For example, in technical domains or Q&A tasks, the model can quote exact definitions, figures, or rulings from a database or text, thereby increasing the likelihood of correctness. Overall, retrieval augmentation lifts the quality of information delivered – the model's outputs align more closely with verified facts, significantly reducing factual inconsistencies compared to an isolated generative model.

## Misinformation Reduction

One of the most critical impacts of retrieval augmentation is the **reduction of misinformation and AI hallucinations** in chatbot outputs. Hallucinations – instances where the AI fabricates information that sounds plausible – are a well-documented weakness of unguided language models. By tethering the generation process to external evidence, retrieval-augmented systems markedly lower the chance of such fabricated answers ([ibm.com](#)). The model is continually “reminded” of real facts from the retrieved data, which constrains it from veering into fantasy or error. Studies and practical implementations have shown that supplying an LLM with relevant context from a knowledge base **helps tame hallucinations and improve the factuality** of its responses ([capestart.com](#)). In effect, external retrieval serves as a correctness

check: the model's output is expected to be supported by the retrieved content, making outright false statements less likely. Furthermore, when a chatbot provides citations or source links as part of its answer, any remaining inaccuracies become immediately apparent to the user, creating an incentive for the system to maintain high truthfulness. This verifiable approach directly combats the spread of misinformation by ensuring that each response can be traced back to reliable data.

## Computational Efficiency

Retrieval-augmented chatbots also influence the system's computational efficiency and resource utilization. Rather than relying on an extremely large model with every possible piece of knowledge encoded in its weights, a RAG system can leverage a more moderately sized model together with a fast external search index. This design can be **more efficient in both development and runtime**. From a development perspective, organizations can avoid frequent model retraining or fine-tuning whenever new data emerges – a process that is computationally expensive and time intensive ([ibm.com](#)). Instead, updating the knowledge source (for instance, indexing new documents) is often sufficient to keep the chatbot's information current. This significantly reduces GPU hours and memory overhead that would otherwise be required to re-train large models, thereby lowering operational costs and simplifying maintenance. At runtime, there's a small overhead to retrieve relevant data, but this can be optimized with high-performance vector databases and caching. In many cases, the overall **processing speed remains comparable or even improves**, since the model has to generate less content purely from scratch once provided with concise relevant context. In fact, retrieval-augmented systems can streamline interactions by quickly pinpointing the needed information: response latency can decrease as the chatbot pulls the answer from a targeted snippet rather than performing lengthy open-ended text generation ([tonic.ai](#)). Additionally, by focusing the model on a narrower context window of retrieved text, the generation task becomes easier and can lead to faster convergence in decoding. In summary, while retrieval augmentation adds an extra step (the search), it often pays for itself through reduced need for large-scale computation elsewhere, resulting in a balanced or improved throughput and a more scalable deployment in practice.

## System Value

Beyond individual queries, retrieval-augmented chatbots provide substantial **strategic value across enterprise, research, and consumer applications**. This architecture has quickly gained prominence as a cornerstone for deploying practical AI assistants in real-world settings. For example:

- **Enterprise:** Businesses are adopting retrieval-augmented chatbots to unlock siloed corporate knowledge and provide **knowledge-based support** for employees and customers. By integrating with internal document repositories and databases, these systems can handle complex queries about company policies, product details, or troubleshooting guides with authority. As evidence of their impact, Retrieval-Augmented Generation (RAG) has seen rapid enterprise uptake – in one industry survey it was reported to dominate with *51% adoption of generative AI architectures, up from 31% in the previous year* ([menlovc.com](#)). Such chatbots deliver reliable 24/7 support,

improve employee productivity by easing information access, and reduce the burden on human support staff.

- **Research:** In academic and scientific research, retrieval-augmented chatbots serve as intelligent assistants that can traverse large bodies of literature. They help scholars and analysts retrieve relevant papers, data, and prior findings on demand. This capability accelerates the research process by synthesizing information from multiple sources. For instance, a researcher can ask the chatbot for a summary of recent findings on a topic, and the system will pull facts from scholarly databases to generate a concise report. RAG-based tools have been used to **create summaries of research papers and documents**, aiding in knowledge sharing and keeping researchers informed of the latest developments ([signitysolutions.com](https://signitysolutions.com)). By connecting experts with up-to-date, domain-specific information, these chatbots promote collaboration and interdisciplinary insight that would be hard to achieve manually at the same speed or scale.
- **Consumer Applications:** Retrieval-augmented chatbots are also transforming consumer-facing services. A prominent example is in search and personal assistants – modern search engine chat modes (such as Bing’s AI-powered chat) combine large language models with web retrieval to provide answers that are both conversational and grounded in real-time information. This means a user’s question about current events, travel, or technical advice is answered with the help of live data rather than stale training content, greatly increasing relevance. Major AI systems have adopted this approach; for example, Microsoft’s Bing chatbot explicitly uses retrieval augmentation as an underlying technique to ensure answers come with up-to-date web results and citations ([cpestart.com](https://cpestart.com)). Similarly, personal AI assistants on smartphones can fetch answers from the internet or a knowledge base before responding, which improves reliability for everyday queries. The result is a more powerful and trustworthy assistant for consumers, combining the **convenience of conversational AI with the confidence of verified information**.

In summary, retrieval-augmented chatbot systems extend the capabilities of conversational AI far beyond what standalone neural models can achieve. They ground language generation in reality, yielding interactions that are accurate, efficient, and credible. This not only elevates the user experience on a case-by-case basis but also unlocks broader use cases – from corporate knowledge management to scientific discovery and daily personal assistance – thereby underscoring the profound value of retrieval augmentation in the evolution of AI-driven communication systems.

## **As for the benefits of our tool and its applications specifically:**

Traditional approaches in information retrieval and generative AI often fall short due to their heavy reliance on basic keyword matching or fully generative models without any grounding mechanisms. These methods frequently result in answers that lack reliability or contain inaccuracies, making them unsuitable for specialized domains where precision is critical. Keyword matching, while straightforward, often fails to capture the nuances and complexities inherent in specialized or nuanced topics. Conversely, unguided generative models, despite their flexibility, tend to hallucinate details, inadvertently fabricating responses due to their probabilistic nature. In contrast, our system prioritizes grounding its answers in concrete, specific, and relevant excerpts taken directly from established authoritative sources. This deliberate grounding method ensures each response is traceable back to the original source material, significantly



increasing transparency and trustworthiness. Rather than leaving responses to chance or guesswork, our method explicitly references authoritative texts, enabling users to verify information independently and reinforcing user confidence in the accuracy and reliability of the provided responses.

Addressing the intricacies and multifaceted nature of Human Design particularly emphasizes the need for such a precise approach. Human Design combines diverse elements such as psychology, astrology, quantum mechanics, and personalized profiling, which together create a complex web of information and interpretation. Standard keyword-driven retrieval methods typically fall short due to their inherent simplicity; these methods can identify terms but lack contextual understanding. Meanwhile, pure generative methods lack any direct linkage to source texts, leading them to generalize or misunderstand nuanced relationships between concepts. Our solution specifically mitigates these issues by identifying exact passages and excerpts that speak directly to the user's query. By pulling precise snippets from authoritative materials, our system respects and preserves the original intent and meaning of the texts. This strategy not only maintains integrity but also ensures nuanced concepts such as personal authorities, gate interactions, and variable-based analyses within Human Design are accurately represented without distortion or oversimplification.

To facilitate this level of precision, we have systematically addressed multiple technical challenges. Firstly, we tackled the limitation of context length inherent in modern language models. Given the finite nature of token limits, careful management of prompt sizes is crucial. We have implemented meticulous strategies to ensure optimal context inclusion without exceeding these constraints. This entails selective prioritization of content, ensuring only the most relevant and impactful information enters each prompt. By refining this selection process, our system optimizes its effectiveness within these strict technical boundaries.

Secondly, we recognized the importance of handling dynamic, constantly evolving content. Traditional retrieval systems often require extensive manual effort to update or rebuild indexes when content changes, proving cumbersome and inefficient. To overcome this limitation, our approach integrates advanced embedding and indexing techniques capable of incremental updates. This methodology allows the system to efficiently re-embed new or modified content and quickly re-index materials, seamlessly incorporating updates without extensive downtime or significant manual intervention. Additionally, managing user-specific queries that require contextual memory adds another layer of complexity. Our system maintains and leverages contextual information from previous interactions, enabling it to generate more coherent, consistent, and personalized responses. This capacity for maintaining continuity is essential in personalized user interactions, particularly within nuanced and intricate knowledge domains like Human Design, where previous interactions and queries significantly impact subsequent information requests.

Our methodology also combines semantic search techniques with sophisticated language model processing to optimize relevance and precision. Through intelligent ranking algorithms, the system assesses text chunks based on their relevance to the user's specific query, ensuring the highest priority is assigned to the most pertinent excerpts. Moreover, our system dynamically adjusts the depth and detail of its responses based on individual user preferences, recognizing that some users seek brief summaries,

while others require extensive, detailed analyses. By dynamically calibrating the detail level of responses, our system can effectively cater to diverse user expectations while maintaining strict adherence to source accuracy.

An essential aspect of our approach is the explicit mitigation of hallucinations, a common pitfall of generative models. Rather than permitting the model to guess in cases of uncertainty, we enforce stringent requirements for direct citations and clear paragraph references. By demanding such explicit grounding, the system significantly curtails instances of fabricated or speculative content, reinforcing overall reliability and user confidence. Performance optimization is another critical consideration, given the computational intensity and associated costs of powerful language models. To maintain efficiency and cost-effectiveness, our system employs embedding caching and selective retrieval strategies. Embedding caching reduces redundant computations, accelerating retrieval processes and minimizing unnecessary model invocations. Additionally, the retrieval of text chunks is strategically limited to only those most relevant to the user's immediate query. This careful management of computational resources ensures the system remains responsive, scalable, and economically viable without sacrificing quality or precision.

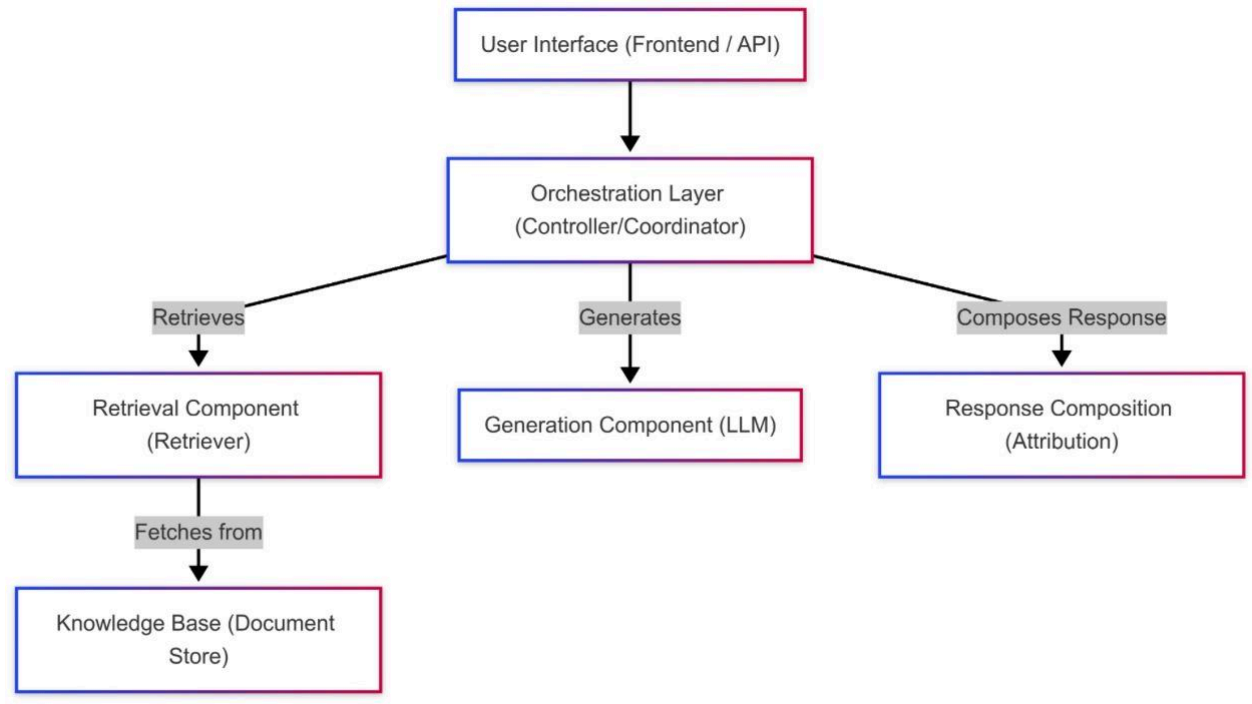
Security remains paramount, particularly given the sensitive nature of personalized and proprietary information. We have implemented rigorous security protocols and robust access controls to protect user privacy and safeguard data integrity. Encryption ensures that private and proprietary content remains inaccessible without proper authorization, while public and private resources are strictly segregated within the system. This separation guarantees secure access to relevant information without compromising overall system usability or user experience.

Collectively, our approach successfully addresses and resolves many common limitations inherent in traditional generative AI systems. By grounding responses explicitly within domain-specific documentation and employing advanced retrieval strategies, we provide a sophisticated, robust, and reliable knowledge assistant. This system is uniquely capable of accurately navigating Human Design's complexity, delivering precise, verifiable insights tailored to individual user requirements. Whether accommodating beginners seeking introductory explanations or experts requiring deep, nuanced analyses, our system effectively eliminates vague generalizations, ensuring every interaction is both meaningful and grounded firmly in authoritative content.

## **Software Architecture**

At a high level, our system's goal is to provide answers to user queries by combining a knowledge base (with relevant information) and a generative model. The user interacts with the system via a user interface (UI), and behind the scenes, multiple components collaborate to retrieve documents and produce a synthesized natural-language response.

Here is a conceptual "blueprint" showing the major areas (like rooms in a building) and how they connect:



## Major Components in the Blueprint

1. User Interface (UI)
  - *Purpose:* Allows the user to enter queries and receive responses.
  - *Architectural Role:* Acts as the *entry point* to the system; it does not hold logic for retrieval or generation, it only captures the user's input and displays output.
2. Orchestration Layer
  - *Purpose:* Coordinates the overall flow (e.g., it receives the query, calls the Retriever, calls the Generator, composes the final answer, and returns it to the UI).
  - *Architectural Role:* The “controller” or “manager” of the system—responsible for directing data among components.
3. Retrieval Component
  - *Purpose:* Finds and returns the most relevant documents or passages in response to a query.
  - *Architectural Role:* Interfaces with the Knowledge Base to search for content that can help answer the question.
4. Knowledge Base

- *Purpose:* Stores the reference documents, domain knowledge, or any form of data from which answers can be derived.
  - *Architectural Role:* Serves as the “source of truth” for the system’s factual information.
  - *Note:* Internally, it may contain indexing structures (e.g., vector indexes) and raw document storage, but at the architectural level, we only need to know it *provides and manages data*.
5. Generation Component (LLM)
- *Purpose:* Produces a human-readable, synthesized answer.
  - *Architectural Role:* Given a user query and relevant context (retrieved documents), it generates a *coherent textual* response.
  - *Note:* The architectural blueprint does *not* dictate whether this LLM is accessed via a cloud API, a local model, or anything else—only that a generative mechanism exists.
6. Response Composition (Attribution)
- *Purpose:* Takes the raw generated text and adds references/citations to the sources used or applies any final formatting.
  - *Architectural Role:* Ensures the answer includes *traceability* back to data sources, fostering trust and transparency.
- 

## 2. Component Responsibilities & Interactions

### User Interface

- Responsibilities:
  - Capture user queries (e.g., through chat input or form).
  - Display system responses, including citations or links.
- Interactions:
  - Sends requests to the Orchestration Layer (e.g., via a network call).
  - Receives final answers (possibly in JSON format) and renders them for the user.

### Orchestration Layer

- Responsibilities:
  - Manage the core sequence: receive a query, invoke retrieval, feed the retrieved documents to the generation component, then invoke the response composer.
  - Handle *errors or fallback behaviors* (e.g., if no documents are found or if generation fails).
- Interactions:
  - Accepts requests from the UI.
  - Dispatches calls to the Retrieval Component (to fetch relevant documents).
  - Dispatches calls to the Generation Component (to produce the answer).
  - Passes the generated answer and document references to the Response Composition for final formatting.
  - Sends the composed response back up to the UI.

## Retrieval Component

- Responsibilities:
  - *Interpret* the user's query in a way suitable for searching the Knowledge Base.
  - *Select* top relevant documents or passages from the Knowledge Base.
- Interactions:
  - Communicates directly with the Knowledge Base (e.g., searching an index or retrieving full documents).
  - Returns the relevant documents to the Orchestration Layer or directly to the Generation Component.

## Knowledge Base

- Responsibilities:
  - Store all domain documents, structured data, or reference information.
  - Provide an *interface* for searching or fetching documents based on query parameters.
  - Potentially maintain specialized indexes (e.g., vector-based, keyword-based) for efficient retrieval.
- Interactions:
  - Receives search queries or document fetch requests from the Retrieval

Component.

- Returns content, metadata, or references to the calling component.

### Generation Component (LLM)

- Responsibilities:
  - Generate textual answers using advanced language modeling.
  - Incorporate *retrieved documents* into its reasoning to produce context-grounded answers.
- Interactions:
  - Receives a user query and the retrieved context.
  - Produces the answer text, returning it to the Orchestration Layer or the Response Composer.
  - May place constraints on input length (context window), which the system must respect.

### Response Composer (Attribution)

- Responsibilities:
    - Format or augment the raw generated text to add *citations*, references, or disclaimers.
    - Optionally filter or truncate content for length or style.
  - Interactions:
    - Receives the generated answer plus the list of documents used.
    - Outputs a final text or structured response that can be displayed to the user.
- 

## 3. Logical Flow of Operations

Below is the typical *end-to-end flow* for a single user query:

1. User Interface
  - User types a question, e.g., “What are the benefits of human design?”
2. Orchestration Layer
  - Receives the question from the UI.
  - Passes the question text to the Retrieval Component.

3. Retrieval Component
  - Interacts with the Knowledge Base to find relevant documents/passages.
  - Returns top matching documents to the Orchestration Layer.
4. Orchestration Layer
  - Provides both the user's original query and the retrieved documents to the Generation Component.
5. Generation Component (LLM)
  - Creates an answer by synthesizing the user query and the retrieved context.
  - Returns the raw answer text to the Orchestration Layer.
6. Response Composer
  - Takes the raw answer and document references.
  - Injects citations or source listings.
  - Returns the final answer text back to the Orchestration Layer.
7. Orchestration Layer
  - Delivers the final composed answer to the User Interface.
8. User Interface
  - Displays the answer (with citations) to the user.

**Details on our implication of this architecture can be found in the appendix.**

## 2. Software Design

Below is a **higher-level design** that shows **what** the main components of our system are, **how** they interact at a more conceptual level, and **how** data is organized—*without* delving into code-level details. This provides the overall “look and feel” of the system from both a user-facing (UI) perspective and a backend (data handling) perspective.

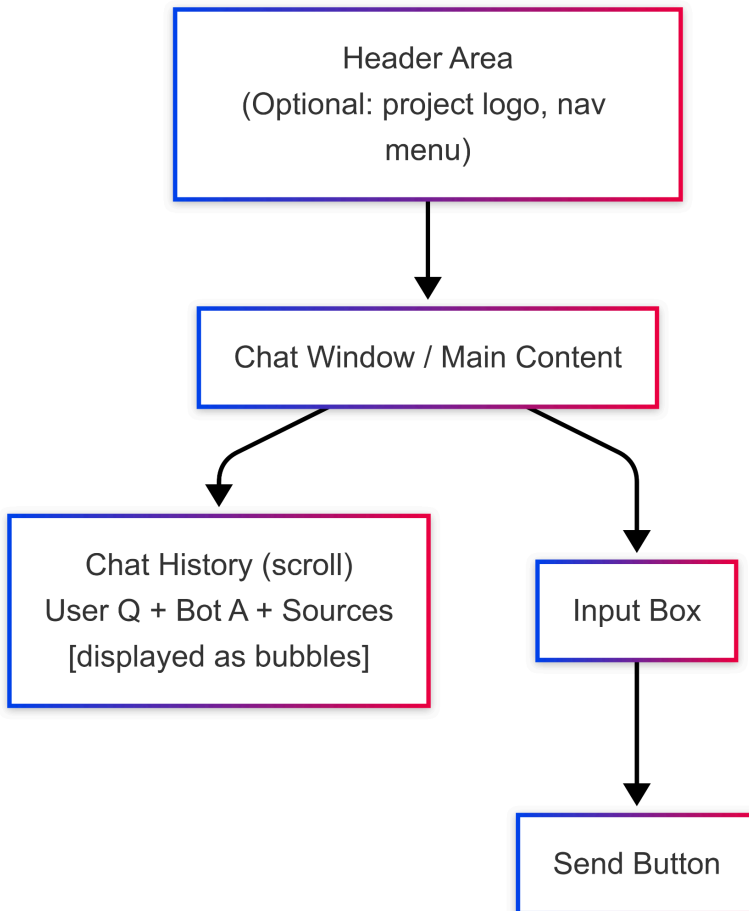
---

### 1. User Interface (UI) Structure

#### Purpose & Layout

- The UI primarily serves as a **chat-style front-end** where users can:
  1. Input their questions (text queries).
  2. View generated answers (with citations).

We can imagine this as a minimal **web application** or **front-end** with a layout like:



## Key UI Elements

1. **Chat History Panel:**
  - Shows a list of **Q&A exchanges** between the user and the chatbot.
  - Each chatbot response may include short references, e.g., “[1], [2]” or a small “Sources” section.
2. **Input Box:**
  - Where the user types their query.
  - Simple text box plus a “Send” button (or press Enter to submit).
3. **Loading / Status Indicator:**
  - Displays when the system is actively retrieving/generating an answer (e.g., a spinner or “Generating answer...”).
4. **(Optional) Additional Buttons/Fields:**
  - Could provide “Clear Chat” or “Download Response” options, but these are **extras** not crucial at this stage.

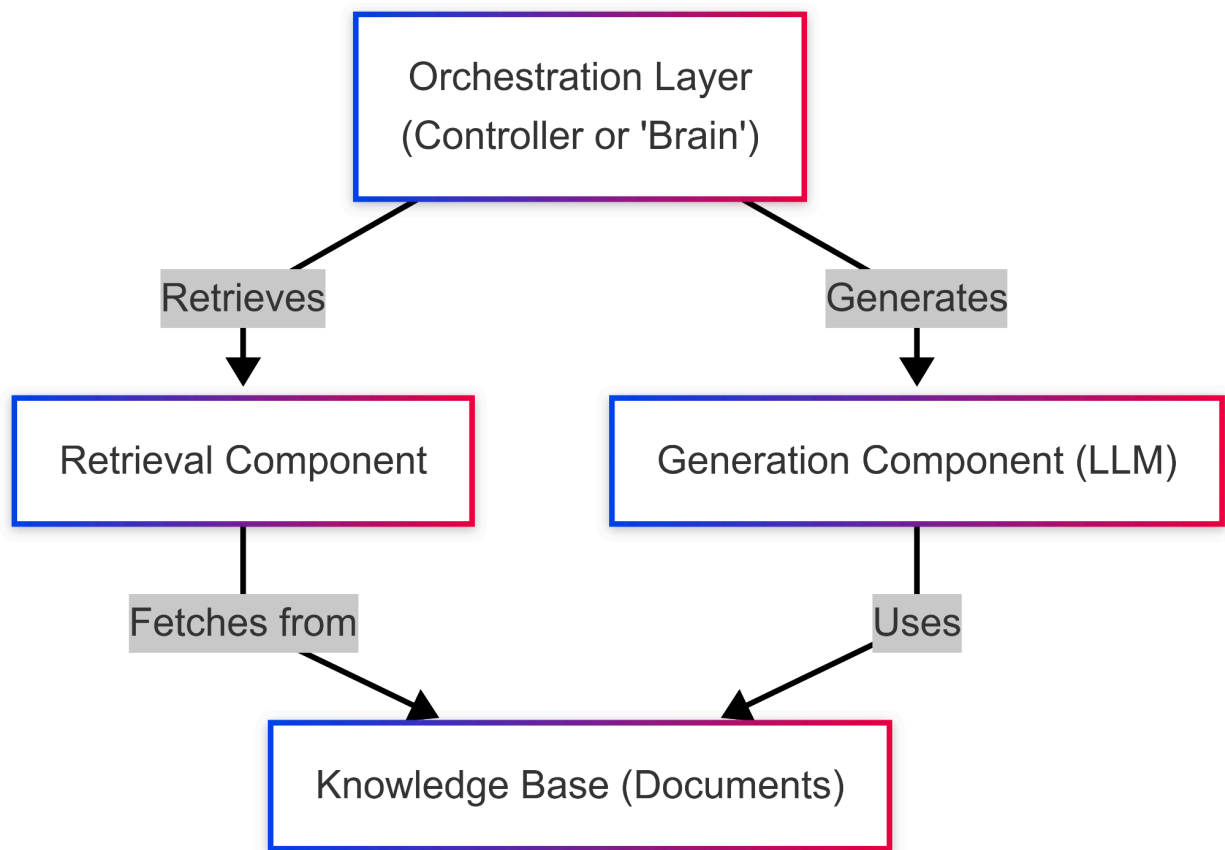


The UI communicates with the backend via a single endpoint (e.g., POST /ask) that **sends the user's query** and **receives** the final answer plus citation data.

## 2. High-Level Backend Design (Conceptual)

### Overview of Backend Components

At a high level, our backend design (the part that powers the chat) can be seen as **three logical layers**:



#### 1. Orchestration Layer:

- Receives the user's query from the UI.
- Coordinates the retrieval of relevant info from the Knowledge Base.
- Invokes the Generation (LLM) for a final, synthesized answer.
- Passes data to the **Response Composition** (for adding citations) before returning the final output to the UI.

#### 2. Retrieval Component:

- Fetches the most relevant documents or passages from the Knowledge Base based on the user's query.

- Maintains or uses an indexing mechanism (e.g., vector-based, keyword-based) to efficiently find relevant content.
- 3. **Generation Component:**
  - Uses a **Large Language Model (LLM)** to produce a human-readable answer.
  - Incorporates the retrieved documents into its prompt or context so that the answer is grounded in actual data.
- 4. **Knowledge Base:**
  - Stores all the domain documents, references, or other data from which the chatbot can draw information.
  - Internally may have two parts:
    1. **Document Store:** Holds the raw text/content and metadata.
    2. **Index** (e.g., vector index): Provides fast lookups for relevant content.
- 5. **Response Composer** (part of orchestration or a small sub-component):
  - Ensures the final answer includes **citations**/references to the source documents.
  - Handles any final formatting (e.g., “Sources: 1. Document A ... 2. Document B ...”).

## 3. Data Organization & Storage

### Document Schema (Knowledge Base)

Each **document** in the Knowledge Base is treated as a record with:

1. **doc\_id** (unique identifier, e.g., a UUID or integer).
2. **title** (short descriptive name).
3. **text\_content** (the actual text).
4. **metadata** (optional key-value pairs: author, date, tags, etc.).

An **index** (vector or otherwise) may store embeddings or text-based tokens separately:

- For vector-based retrieval, each document or chunk is represented by an **embedding vector**.
- A mapping from **doc\_id** → **embedding vector** is kept so we can quickly find the top relevant results.

### Query & Retrieval Data Flow

- When the user sends a query, the system might convert it into a **query embedding** (if using semantic search) or treat it as a keyword set.
- The retrieval component searches the Knowledge Base’s **index** to find a handful of documents most closely matching the query’s content.

### Conversation History (Optional / Extended Feature)

- If the system supports multi-turn conversations, we store recent Q&As or a summary in a session record.
- This may be in a simple table like ConversationHistory:
  - **session\_id**
  - **turn\_index**
  - **user\_query**
  - **bot\_reply**

This helps maintain context across multiple user turns, but it's not strictly mandatory in a *single-turn* scenario.

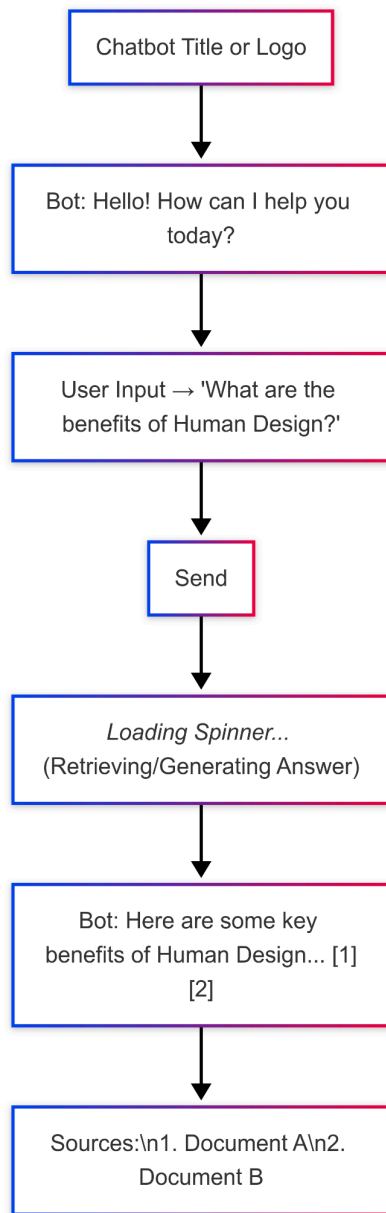
## Logging & Analytics (Optional)

- For debugging or improvement, we might log queries, the retrieved docs, and the generated responses.
- A Logs table could store:
  - **log\_id**
  - **timestamp**
  - **user\_query**
  - **retrieved\_doc\_ids**
  - **final\_answer**

This data can help refine system performance later.

## 4. UI Mockups & Interaction Flows

### Example UI Mockup



## High-Level Interaction Sequence

1. **User Inputs Query:**
  - The user types a question like “What are the benefits of Human Design?” and clicks “Send.”
2. **UI → Backend Call:**
  - A request is sent to POST /ask with { "query": "...user's text..." }.
3. **Backend Orchestration:**
  - Reads the query.
  - **Retrieval:** Finds relevant documents in the Knowledge Base (e.g., top 5 matches).

- **Generation:** Feeds user query + retrieved doc snippets to the LLM, gets a synthesized answer.
  - **Response Composition:** Attaches references, possibly in the format: [1], [2] or “Sources: ...”.
4. **Backend → UI:**
- Returns the final composed answer plus citation data (e.g., a JSON with { "answer": "...", "sources": [{id, title}, ...] }).
5. **UI Displays Answer:**
- Shows answer text in a chat bubble and includes the sources or references.

## 5. Notable Design Choices (at a Higher Level)

- **Why a Chat Interface?**
  - Conversational UI is intuitive. Users ask questions naturally, see immediate responses.
- **Why a Separate Retrieval vs. Generation?**
  - Allows modular design. We can **upgrade** the retrieval method or **swap** the LLM without altering the entire system.
- **Why Keep Citation Logic Separate?**
  - Improves clarity. The “raw” LLM answer can be post-processed to add references systematically, ensuring consistent formatting and transparent source attributions.
- **Data Organization**
  - Using a *document-centric approach* with minimal required fields (ID, title, content). Additional metadata is optional but can be leveraged for advanced filtering or display.

**Further details on our implication of this architecture can be found in the appendix.**

## 3. Coding Guidelines

### Python Coding Style Guidelines

We will adopt **PEP 8** as our primary coding style guide for the Python codebase. PEP 8 is the official Style Guide for Python. It covers conventions for code layout, naming, and formatting that improve the consistency and readability of the code. Adhering to PEP 8, therefore, means using meaningful naming conventions, proper indentation, spacing, and other layout practices that make it easier for any developer on the team (or future maintainers) to understand the code quickly. In addition to PEP 8, we have considered the **Google Python Style Guide** as a secondary reference. We will treat PEP 8 as the baseline and use Google’s guide for any areas not explicitly covered by PEP 8.

### Rationale for Chosen Guidelines

The above guidelines were chosen for their **proven value in maintaining code quality**. PEP 8 is chosen because it is the standard. Google's style guide is chosen to supplement because it's well-documented and addresses practical issues.

### **Additional efforts to maintain clear Coding Guidelines and Structure**

We will maintain a short PROJECT\_CODESTYLE document in the repo that highlights our key style rules and tool configurations (for example, stating “We follow PEP 8; use snake\_case for functions and variables, limit line length to 79 characters, etc.”). This document can reference the external guides (with links to PEP 8 and Google style sections) so that developers can read more if needed.

## **4. Process Description – Risk Assessment**

### **1. Integration Complexity of RAG Pipeline:.**

- **Likelihood:** Medium. We have multiple new components (vector DB, embedding model, LLM API) that need to work in unison. Each by itself is complex, and integrating them requires careful engineering. However, given that the lead of the project has fairly extensive experience implementing similar symptoms before, we hope it will mitigate the risks.
- **Impact:** High. If this risk materializes, it could significantly delay development. Integration issues could block core functionality – e.g., if the retriever and LLM aren't interfacing correctly, the chatbot won't work at all or will give nonsense answers.
- **Evidence:** The complexity of RAG-based systems in practice is well-noted. Additionally, since we are essentially orchestrating between different libraries or services, the interfaces might not align perfectly. For example, we might need to format retrieved text carefully to not overflow the LLM's input, handle encoding differences, etc.
- **Risk Reduction/Detection:** To reduce this risk, we are planning early integration tests. Rather than building each component fully in isolation and only integrating at the end, we will create a skeleton end-to-end pipeline as soon as possible (see Github for the first draft of skeleton).
- **Mitigation Plan:** If we encounter major integration hurdles (e.g., the chosen vector DB is not working as expected with our data size, or the LLM API integration has problems), we have a few mitigation strategies. We can allocate

additional expertise or resources to the problem area; for example, if vector search is the issue, we might bring in a simpler solution (like using a simpler similarity search in memory as a temporary fix) while we resolve the complex tool. In the worst case, if integration issues severely threaten timelines, we might **simplify the scope** as a mitigation. This would be last-resort as it might reduce our control, but it's an option to ensure we deliver a working system.

## 2. Data Quality and Knowledge Base Preparation Issues

- **Likelihood:** Low. Building the knowledge base involves gathering documents, cleaning them, and embedding them. If the data is unstructured (e.g., various formats or noisy text), processing it can be tricky. However, there is still a low chance we'll encounter problems like missing data, encoding issues, or simply not having enough quality content on the topics users will ask, because we are working with very extensive and saturated data set and the format of the data implies minimal noise and is not overly complex structure which significantly reduces risks.
- **Impact:** High. The chatbot's usefulness directly depends on the knowledge base content. If data preparation is delayed, it bottlenecks everything. If the data is low-quality or full of errors, the chatbot might give incorrect or confusing answers, failing the project goals. So a failure in this area could render the end product unsuccessful even if the rest of the parts work.
- **Risk Reduction/Detection:** Our approach to mitigate this is to start the data work early. We have planned a data sourcing and preprocessing task as one of the first steps in the implementation schedule. By doing it early, we will quickly see how hard it is and can adjust scope if needed. We'll also set up validation on the data: simple checks like ensuring no document is empty, each has a title, content is readable, etc., to catch issues. Another reduction strategy is prioritizing: we identify a core set of must-have documents and ensure those are processed first. This way, even if the full dataset isn't ready in time, the most important info is in the system. To detect quality issues, we'll do some sample Q&A tests on the knowledge base *before* integration with the LLM. For example, we can manually query the vector index or do simple keyword search to see if relevant docs come up for expected questions – this will reveal if our data is missing key answers or if embeddings aren't working properly.
- **Mitigation Plan:** If we find that the data preparation is lagging or data quality is poor, we have a mitigation in mind. If it's a matter of volume, we will narrow the scope of content and explicitly communicate the limitations in the final system. In the worst case, if the data we have is insufficient, we might pivot the project's focus slightly: for example, turning it into a more general chatbot with some knowledge rather than a comprehensive expert. This would be a scale-back in

requirements to ensure we still deliver something functional. We'll also keep the system flexible: if a question is asked that isn't covered in the knowledge base, the chatbot will respond with a graceful failure ("I'm sorry, I don't have information on that") rather than trying to answer falsely. That mitigates the impact of poor coverage by at least not misleading the user.

### 3. Performance and Scalability Issues

- **Likelihood:** Low to Medium. We anticipate some performance overhead with RAG (embedding search and large model inference). However, our initial target usage might be modest (a demo for this class). On a small scale, it's likely fine. The risk becomes medium if we consider growth or a requirement for real-time responses under high load (what goes beyond the scope of this class). We haven't observed performance yet, so there's uncertainty. Past experience with LLM APIs shows response times in the 1-3 second range for short prompts, which might be okay. The vector search on a few thousand embeddings is usually very fast (milliseconds). Still, integration and longer context might slow it down.
- **Impact:** Medium. If performance is poor (e.g., it takes 10+ seconds to answer or cannot handle concurrent queries), user experience suffers and it could derail stakeholder support for the project. However, if it's only slightly slower than desired, we can sometimes optimize or users might tolerate it initially (hopefully:).
- **Risk Reduction/Detection:** We will address this risk by incorporating **performance testing** in our schedule. Even early on, after the first integrated version, we'll measure how long a typical query takes end-to-end. This will be done in a test environment with timers around major steps (retrieval time, generation time). We'll also do some load testing with concurrent requests (if feasible) to see how the system behaves. By doing this testing throughout development, we can detect any slow component. For example, if we find that generation is taking 5 seconds on average, we might explore optimizing prompt size or using a different model. On the retrieval side, if embedding creation per query is slow, we might use a smaller embedding model or ensure it's running on GPU. We also plan to profile the code to detect any unexpected slowdowns (maybe an inefficient loop in our code, etc.). Another preventive measure is setting reasonable limits: for instance, we might restrict the length of input questions or the number of documents we retrieve (top 3 instead of top 10) to cap the amount of processing the LLM must do. This reduces worst-case latency.
- **Mitigation Plan:** If performance tests show we're too slow, we have several mitigation options. One is to **optimize the components** – e.g., use batching in the vector search, or caching frequent queries. Caching could be very effective: many users might ask similar questions, so caching an answer for popular queries would make repeat responses instantaneous. Another option is to switch to a faster



model. On the vector index side, if we have far more data than expected, we can use approximate search or sharding of the index to keep search quick. As a last resort mitigation, we would explicitly document the limitations (for example, “System may take up to N seconds to respond during heavy load” or limit concurrent usage in the initial release).

#### 4. Accuracy and Hallucination Risk

- **Likelihood:** Medium. Even with relevant documents retrieved, there is a known risk that LLMs sometimes generate plausible-sounding but incorrect statements (hallucinations). RAG is specifically meant to reduce it, but it is not foolproof – the model might misinterpret the context or combine information incorrectly. But, it is worth noting that it largely depends on how well we design prompts and choose the model; with careful prompt engineering, we can lower it, but it cannot be eliminated entirely.
- **Impact:** Medium to High. If the bot routinely gives wrong answers, the project’s credibility is at stake. Users could make wrong decisions based on false information, which is especially dangerous if the domain is something like finance or health. Even one egregious mistake can reduce user trust significantly. Since a core requirement is that the system provides accurate and supported answers, failure here could mean not meeting a primary project goal. It could also have ethical or business implications if misinformation is given (especially given that we work with a real client).
- **Risk Reduction/Detection:** We are taking multiple steps to reduce hallucinations and inaccuracies. First, our prompt design will explicitly instruct the LLM to use the provided information and not assume anything not given. For instance, we will add a line: “If you don’t find the answer in the provided content, say you don’t know.” This kind of instruction can reduce the model’s tendency to just guess. Second, we will perform extensive testing of answer outputs. During development, for a set of sample queries, we will manually verify the answers against the source documents. If we catch instances where the model is veering off, we’ll adjust our approach (either by refining the prompt or by ensuring the retrieval brings more precise info). Another detection method is implementing a self-check: for example, after the model generates an answer, we could feed the question and answer back into a verification step (like another prompt asking “Is this answer supported by the given text?”). We might not fully implement that, but it’s an idea if time permits. Additionally, by always including source citations in answers, we provide a way to trace accuracy. If an answer has a citation, users or testers can check that source. During testing, if we see an answer with a citation that actually doesn’t support the answer, that’s a red flag of hallucination or error. We will also reduce risk by retrieving multiple documents and possibly doing a mild validation:

if the top documents all lack a certain answer and the model still produces one, that answer is likely suspect.

- **Mitigation Plan:** If despite our precautions the chatbot outputs incorrect information, our mitigation focuses on damage control and iterative improvement. One mitigation is to incorporate a feedback loop: when users or testers find an incorrect answer, we log it as a bug and add that question to a test suite. Over time, this builds a set of queries that we know were problematic, and we adjust the system to handle them. For example, if the bot hallucinated an answer for a query that wasn't actually answerable, we might adjust the prompt to encourage an "I don't know" response for that scenario. Another mitigation is to implement a simple post-processor that checks whether the answer contains facts that are directly present in the retrieved text. If not, we could choose not to present that answer. If accuracy issues are frequent and we can't solve them within the project time, a contingency is to limit the use cases in deployment: perhaps we label the system as "beta" or "for internal use only" and include disclaimers that answers should be verified. This is not ideal, but it might be necessary to set user expectations correctly.

**Comparison with Initial Risk Assessment:** In the initial requirements document, we identified several of these risks, but some estimates and mitigation strategies have evolved:

- Integration complexity and technical unknowns were initially recognized; we've now narrowed down our technology choices which lowers uncertainty a bit, but we maintain vigilance via early testing.
- Data preparation risk was initially somewhat underestimated (we love real world data:) We adjusted by starting that earlier and prioritizing data scope (this is a change from initial plan).
- Hallucination/accuracy was known from the start as an AI project risk; our understanding of mitigation (citations, prompt tuning) has deepened after prototyping and research. We're more concrete now in how to handle it versus initial broad statements.
- Performance wasn't a top focus in initial risk brainstorming (we were more concerned with "will it work correctly" than "will it be fast"). After some discussions, we realized user experience matters, so we've added performance as a notable risk to continuously monitor.
- Scope creep was a concern initially given the novelty factor. We've addressed a lot of that by finalizing requirements now. So while initially it was medium risk, now it's lower because we've set clearer boundaries.

## 5. Project schedule

### Week 7: Development & Integration

During this week, the focus will be on integrating all core system components and ensuring they work together seamlessly. The primary objectives include:

#### 1. Finalizing Data Integration Workflow

- Complete the ingestion pipeline to process documents and convert them into embeddings.
- Ensure data formatting and cleaning steps are finalized to avoid inconsistencies in retrieval.
- Verify that the document storage structure allows for efficient updates and retrieval.

#### 2. Implementing the Embeddings Generation Pipeline using OpenAI's Models

- Develop a structured pipeline that takes raw text documents and generates high-dimensional vector representations using OpenAI's text-embedding-ada-002.
- Ensure that all embeddings are indexed in the vector database (Qdrant or FAISS) correctly and are queryable in real time.
- Conduct preliminary tests to validate that embeddings are meaningful and properly represent the document semantics.

#### 3. Initial LLM Integration with OpenAI GPT-4o for Response Generation

- Set up API calls to OpenAI's GPT-4o and ensure proper authentication and rate limit handling.
- Implement a structured prompt format that integrates retrieved documents with user queries for optimal response generation.
- Conduct initial response evaluations to check if the generated answers are coherent, relevant, and correctly grounded in retrieved documents.
- Address any early issues with hallucinations or missing citations by refining the retrieval-prompt interaction.

At the end of Week 7, the system should be capable of taking user queries, retrieving relevant documents, and generating LLM-powered responses, forming a **Minimum Viable Prototype (MVP)** of the chatbot.

## Week 8: Testing & Optimization

This week is dedicated to improving system performance, ensuring reliability, and addressing security concerns. The major focus areas include:

### 1. Conducting Performance Testing

- Measure query response times across different scenarios (short vs. long queries, simple vs. complex retrievals).
- Analyze potential bottlenecks in retrieval, embedding search, or LLM response time.
- Optimize system parameters (e.g., limiting the number of retrieved documents, adjusting LLM temperature settings) to balance accuracy and speed.
- Implement caching mechanisms for frequent queries if necessary to improve efficiency.

### 2. Addressing System Reliability Risks

- Simulate real-world usage conditions, including handling concurrent queries.
- Introduce fail-safe mechanisms for scenarios where external APIs (such as OpenAI's GPT-4o or vector search) experience downtime.
- Implement timeouts and fallback responses to avoid system failures due to delayed LLM responses.
- Ensure robust logging is in place for error tracking and debugging.

### 3. Implementing Security Measures

- Secure API keys and authentication credentials using environment variables or secrets management tools.
- Review and test input sanitization to prevent prompt injection or adversarial attacks.
- If applicable, validate that user queries do not expose sensitive data from documents beyond their intended scope.

By the end of Week 8, the chatbot should be **stable, performant, and secure**, ready for final validation.

## Week 9: Final Adjustments & Presentation Preparation

This final week is dedicated to refining and polishing the system before submission and

presentation. Key activities include:

### 1. **Finalizing Testing Results & Risk Mitigation Strategies**

- Run final validation tests to confirm all use cases are handled correctly.
- Ensure all major bugs and performance issues identified in previous testing are addressed.
- Verify that fallback mechanisms for failure scenarios (e.g., retrieval errors, long response times) function as intended.

### 2. **Preparing Presentation & Demo Materials**

- Create a structured project presentation covering the chatbot's objectives, technical approach, key features, and challenges faced.
- Develop a live demo showcasing how the chatbot processes user queries, retrieves documents, and generates responses.
- Highlight improvements made based on feedback (such as citation handling, hallucination reduction, or UI enhancements).
- Prepare visuals (e.g., architecture diagrams, performance graphs) to explain system components and workflows effectively.

### 3. **Conducting Final Internal Review Before Submission & Presentation**

- Ensure all documentation (User Guide, Developer Guide, Admin Guide) is complete and clearly written.
- Conduct a final team review to test the chatbot in different scenarios and refine presentation delivery.
- If applicable, create a backup demo or pre-recorded session in case of technical issues during the live demonstration.

By the end of Week 9, the project should be **fully functional, well-documented, and ready for submission and presentation** to stakeholders.

## 6. Team Structure

Our project team is composed of 2 members, each with specific roles and responsibilities to ensure all aspects of the project are covered. Below is the team structure and the primary responsibilities of each role:

- **Project Manager, Team Lead and AI Engineer (Alina Hyk)** – Alina acts as the coordinator of the project. Her responsibilities include planning and tracking the schedule, organizing meetings, and ensuring that the team is communicating effectively. She also

interfaces with “stakeholders” (TA and client) to report progress. In addition, Alina contributes to requirements definition and documentation. She has a technical background, so she will also assist in code reviews and integration tasks as needed, making sure the components the team develops fit together per the architecture. Alina is also responsible for the core AI components of the chatbot. This includes selecting the embedding model for document encoding, setting up the vector database, and implementing the Retriever. She will also work on prompt engineering and help integrate the LLM. Alina’s expertise in NLP and ML makes her the point person for anything related to the model’s behavior or the retrieval quality. For example, he will evaluate embedding effectiveness, and adjust if needed. She will also design the evaluation tests for answer accuracy and help with mitigating hallucination. In terms of deliverables, Alina will produce the Knowledge Base and Retrieval subsystem, and contribute significantly to the Generator logic.

- **Backend Engineer (Alina Hyk)** – Alina focuses on the software engineering side of the backend. She will implement the API endpoints (the controller that receives queries and returns answers). She is responsible for the code that connects the Retriever with the Generator and formats the output (Response Composer). She ensures that the system is robust – adding proper error handling, logging, and making the code modular as designed. She will also implement unit tests for the backend components. Security and performance considerations in the backend (like input validation, avoiding injection issues, or making sure asynchronous calls are handled) fall under her purview.
- **Frontend Engineer (Samuel Jamieson)** – Samuel works on the user interface. He ensures that when the user types a question, it’s sent to the backend, and that the answer (with citations) is nicely displayed. He will also implement any user-centric features like a “clear chat” button or loading spinner. He focuses on making the chatbot easy to use and visually clear. He will also conduct usability tests and refine the UI based on feedback.

**Communication** within the team is set up as follows: daily quick check-ins on progress (especially during heavy weeks), a shared Discord channel for questions, and a weekly meeting to review status versus the schedule. Given each person’s focus area, we expect most day-to-day questions will be handled by the person responsible, but any major decisions (like changing a library or adjusting a requirement) will be discussed among the team and TA.

### **Additional Rules**

1. **Respect:** Treat teammates courteously; value diverse perspectives and encourage open dialogue.
2. **Clear Communication and Consistent Involvement:** We expect regular asynchronous communication via Discord. We are also expected to attend weekly

appointments with the TA and stay after them as needed for synchronous work or specific troubleshooting.

3. **Attend Meetings:** We are expected to attend all scheduled meetings to stay informed. If unable, we are expected to notify each other at least 20 hours prior.
4. **Deadlines:** We acknowledge that we are starting on a slightly different timeline than the other teams due to the late start of the project. Therefore, one of our main immediate goals is to catch up on the steps we missed, which will be our focus until the end of week five. More generally, we are expected to follow the suggested timeline of submissions on Canvas and submit them no less than two hours before the deadline. We are also expected to have materials ready for all main presentations at least 24 hours in advance.

## 7. Test Plan & Bug Tracking

A comprehensive test plan will be implanted. We outline below what aspects of the system will be tested, the types of testing to be conducted (unit, integration, system, usability), and our approach to bug tracking.

### Testing Objectives and Scope

We aim to test **functionality, performance, and usability** of the RAG chatbot:

- Functionality testing will verify that each feature (each requirement and use case) works as intended.
- Performance testing will ensure the system responds within acceptable time limits and can handle expected load (e.g., multiple queries in succession).
- Usability testing will ensure the user interface is intuitive and that the answers provided are presented in a clear, useful manner (including source citations).
- Security-related testing (within scope) will include making sure there's no way for a user to maliciously break the system via input.

### Aspects to be tested include:

- **Retrieval accuracy:** Given various query inputs, ensure the correct documents are being retrieved from the knowledge base. We'll test this by crafting queries where we know which document is relevant and checking if it's in the top-k results. This validates the embedding and vector search pipeline.
- **Generative answer correctness:** For a set of known Q&A pairs (where we know the expected answer and source), verify the chatbot's answer matches the expected

information. This covers requirement validation for correctness.

- **Citation presence and correctness:** Whenever an answer is given, ensure that sources are cited. We'll test cases to verify that citations correspond to the documents that truly support the answer..
- **Edge cases in queries:** Test questions that are unusual or tricky: very short queries, very long queries, queries with typos, or multi-part questions. This ensures robustness. We might test a nonsense query to see if the system properly says it can't find an answer.
- **Multi-turn behavior (if implemented):** If we support follow-up questions, we will test scenario sequences. E.g., User: "What is RAG?" (gets answer), then User: "How does it improve accuracy?" (the bot should understand "it" refers to RAG and answer accordingly). We verify context is maintained.
- **User interface flows:** Test the front-end: entering a question, getting an answer displayed. Check that the UI doesn't break if a question is asked while another is loading, etc. Also test any UI elements like the "Clear conversation" button if present and works correctly.
- **System integration:** End-to-end testing from UI to backend to knowledge base and back, to ensure all parts together function with no disconnect.
- **Performance:** Measure response time for typical queries and heavy queries (like long context or multiple documents). Also test with concurrent accesses if possible (maybe simulate 10 users asking at same time) to see if any race conditions or resource contentions occur.
- **Error handling:** Induce some errors to see how the system behaves. For example, temporarily disable the vector index or LLM API (simulate it) and see if our fallback messages appear. Or feed an extremely long gibberish query to see if the system times out gracefully.
- **Security tests:** Although not a primary focus, we will test some injection scenarios. For instance, ask the bot a question that might be trying to manipulate the prompt (like "Ignore previous instructions and do X"). Ensure that our prompt design mitigates straightforward attempts. Also, if the UI is web-based, test for things like script injection in user input (the backend should treat everything as plain text).

## Bug Tracking

We will use **GitHub Issues** for tracking bugs and issues throughout development, as decided in our project management approach. Our usage of GitHub Issues will follow these practices:

- When a team member finds a bug (e.g., a test fails, or something looks wrong in a manual test), they will create a new Issue if one doesn't exist already. The issue will have a descriptive title and a detailed description including steps to reproduce the bug, expected vs actual behavior, and severity.



- We will label the issues appropriately, e.g., **bug**, and possibly severity labels like *critical*, *minor*, etc. We might also use labels like *UI*, *backend*, *data* to categorize bugs by area.
- We will assign the issue to a team member.
- We will also use the issue tracker to manage enhancements or tasks, but we will clearly differentiate bugs from feature tasks.
- Issues will be linked to code commits or pull requests when a fix is implemented (using keywords like “Fixes #123” in commit messages so that GitHub auto-closes the issue when merged).
- We will regularly review open issues in team meetings to ensure none are being overlooked and to prioritize fixes for critical ones.
- GitHub Issues will effectively serve as our bug database and also a way to track progress on resolving them.

Our goal is to resolve all critical and major bugs before final release. Minor issues that are not critical might be noted as future work if they don’t significantly hinder functionality or user experience, but we aim to polish as much as possible.

**Summary of Testing Process:** We start with unit tests as components are built, add integration tests as we glue components together, and finally perform system and usability tests on the whole. By final delivery, we expect to have a robust test suite that can be run at any time to verify the system, and a bug tracker that ideally has zero open high-priority issues (and any remaining issues documented for future improvement).

## 8. Documentation Plan

Documentation is essential for users, administrators, and future developers to understand and effectively use our RAG-based chatbot system. We plan to provide several types of documentation, each tailored to a different audience and purpose. The following outlines the documentation deliverables and how they will be structured:

- **User Guide (End-User Documentation):** This will be a concise, friendly manual for end users of the chatbot. It explains *how to interact with the chatbot* and what to expect from it.
  - *Content:* Instructions on how to start a chat session (e.g., accessing the web interface or sending a message), examples of questions the user can ask, and guidance on phrasing queries for best results. It will include screenshots of the chat interface with explaining interface elements (like where the answer will appear, what the citation links mean, etc.).
  - *Format:* PDF or web page. The user guide will not assume any technical

background; it will be written in simple language.

- *Location:* This will be provided as part of the deliverables (e.g., in the repository's docs folder as UserGuide.pdf, and possibly hosted on a help menu in the app).
- **Administrator Guide (System Admin / Deployment Documentation):** This document is for those who will maintain the system in production .
  - *Content:* Instructions on how to deploy and configure the chatbot system. This includes prerequisites (software dependencies, hardware requirements), installation steps, how to start/stop the server, and how to update the knowledge base (e.g., adding new documents or re-running the embedding process). It will also cover configuration options such as changing the LLM API key, adjusting parameters like number of retrieved documents or switching out the vector database if needed.
  - Additionally, it will have sections on maintenance tasks: for example, monitoring logs for errors, how to retrain or update embeddings when adding new data, and backup procedures for the knowledge base. If any scheduled tasks are needed (like periodic re-indexing), those will be documented.
  - *Format:* A written manual, likely in Markdown or PDF form, included in the repo (AdminGuide.pdf or similar). It will assume the reader has some technical knowledge (comfortable with command line, etc.), as it's aimed at system admins.
- **Developer Guide (Technical Documentation for Developers):** This is aimed at developers who might work on the project in the future, or evaluators who want to understand the internals..
  - *Content:* An overview of the system architecture and design (somewhat a summary of this design document, but focused on implementation). It will describe each major component (retriever, generator, etc.) in terms of code structure: which modules/classes correspond to which components, important algorithms or flows (maybe including simplified class diagrams or pseudocode for core logic). It will also cover how to set up a development environment, run the test suite, and guidelines for contributing (coding conventions, how to add a new feature).
  - This guide will also document any important decisions or configurations in code – for instance, if we have certain hyperparameters or thresholds (like top\_k documents = 5), we'll note where and why. We will include references to relevant sections of code (like “see retriever.py for the vector search implementation”).
  - *Format:* Likely a Markdown file in the repository (e.g., DEVELOPERS.md). We might also use the project's GitHub Wiki for this purpose, which can be easier to navigate online and update.
- **Testing and QA Reports:** Although primarily internal, we will document the test results

and provide a summary in the documentation and final presentation:

## Finalized Requirements

We list the final requirements divided into functional and non-functional:

### Functional Requirements:

1. **FR1: Domain-Specific Q&A** – The chatbot shall answer user questions using the information in the provided knowledge base of documents. It should cover the defined domain topics. This is the primary function: for any question within the domain, attempt to retrieve relevant content and respond with an answer drawn from that content.
2. **FR2: Source Citation** – The chatbot shall provide source references for its answers. For each answer that uses knowledge base content, the system should cite the document(s) it used (e.g., by title or an identifier). This addresses some converse on transparency and backtracking of sources that were brought up in the discussion posts of students to our presentation.
3. **FR3: Unanswerable Query Handling** – If the user’s query cannot be answered based on the knowledge base (no relevant info found or question is out of scope), the chatbot shall respond with a polite statement of inability to answer, rather than guessing or giving incorrect information. Essentially a built-in “I don’t know” response.
4. **FR4: Multi-turn Context** – The chatbot shall support basic follow-up questions in a single session. Specifically, it will allow a user to ask a related question referencing the previous interaction, and the chatbot will use the conversation context to interpret it.
5. **FR5: Knowledge Base Update Mechanism** – The system shall provide a mechanism to update the knowledge base with new or modified documents without major code changes.
6. **FR6: System Feedback & Help** – The chatbot interface shall include a way for the user to understand how to use it and should handle minor user input errors.

### Non-Functional Requirements:

8. **NFR1: Accuracy** – The chatbot’s answers should be factually correct with respect to the content in the knowledge base. (While 100% accuracy may not be measurable, our goal is to minimize incorrect answers/hallucinations. We set a target: in testing, at least 90% of answers for test queries are correct and supported by the documents. This ties to how we evaluate success.)
9. **NFR2: Performance** – The system shall generate a response within an acceptable time frame. We define acceptable as *on average under 10 seconds* for a typical question with a knowledge base of our size (and ideally faster, 3-5 seconds, for great UX).
10. **NFR3: Scalability (Data)** – The system should handle a knowledge base of at least a few

hundred documents without significant degradation in performance. (This requirement ensures our design choices can scale to a reasonable size, as initially planned).

11. **NFR4: Usability** – The chatbot UI must be intuitive and easy to use for a non-technical user. This means the language in responses should be clear, the formatting (such as citations or answer highlighting) should be easy to read, and any UI elements (buttons, text input) should be clearly visible and labeled.
12. **NFR5: Reliability** – The system should be stable during use; it should not crash or hang for typical inputs. Also, it should handle unexpected inputs gracefully (e.g., extremely long questions, or inputs with special characters) without crashing.
13. **NFR6: Maintainability** – The codebase and configuration should be organized and documented such that a developer can modify or extend the system (this is a quality requirement aimed at future development).

## Use Cases (Functional Requirements)

### Use Case 1: Consultant Assistance for Human Design Analysis

Actors: Human Design Consultant, RAG-based Chatbot

Scenario: A Human Design consultant seeks quick and accurate information regarding a specific concept during a session with a client.

1. **Consultant Query:** The consultant asks the chatbot, “Explain the Emotional Authority and its practical application.”
2. **System Response:** The chatbot retrieves relevant materials from its corpus, providing a concise explanation of Emotional Authority, along with references to key texts and transcriptions.
3. **Outcome:** The consultant uses the response to guide the session, saving time and ensuring accuracy.
4. **Extensions:** The chatbot suggests related concepts such as Splenic Authority or Sacral Authority, enriching the session.

### Use Case 2: Client Self-Discovery and Personal Growth

Actors: Human Design Client, RAG-based Chatbot

Scenario: A client exploring their Human Design profile wants to understand how their

specific design elements affect their relationships.

1. **Client Query:** “How does having a Defined Heart Center impact relationships?”
2. **System Response:** The chatbot retrieves text segments explaining the Defined Heart Center’s characteristics, focusing on how it influences interpersonal dynamics and self-worth.
3. **Outcome:** The client gains deeper insights into their relationship patterns and how to leverage their strengths.
4. **Extensions:** The chatbot recommends reflective exercises based on existential psychology to encourage personal growth

### **Use Case 3: Educational Support for New Consultants**

Actors: New Human Design Consultant, RAG-base Chatbot

Scenario: New consultants undergoing training use the chatbot as an educational resource to clarify complex topics.

1. **Consultant Query:** “What are the key differences between Projector and Manifestor types?”
2. **System Response:** The chatbot provides a detailed comparison of the two types, referencing specific sections of Human Design manuals and transcriptions.
3. **Outcome:** The new consultant quickly understands the differences, enhancing their learning process.
4. **Extensions:** The chatbot suggests additional readings and video lectures for deeper understanding.

### **Use Case 4: Multidisciplinary Insight for Therapeutic Sessions**

Actors: Therapists, Human Design Chatbot

Scenario: A therapist integrating Human Design concepts into their practice seeks to combine Human Design with psychological frameworks like existential therapy.

1. **Therapist Query:** “How can the Defined Ajna Center be interpreted in the context of existential therapy?”
2. **System Response:** The chatbot combines Human Design knowledge with psychological insights, offering a holistic interpretation and suggesting therapeutic

approaches.

3. **Outcome:** The therapist provides a more nuanced and personalized session.
4. **Extensions:** The chatbot suggests additional therapeutic exercises aligned with Human Design principles.

### **Use Case 5: Visual and Diagrammatic Assistance (Stretch Goal)**

Actors: Human Design Consultant, RAG-based Chatbot

Scenario: A consultant wants to visualize key concepts using diagrams.

1. **Consultant Query:** “Show me a diagram of the Human Design Centers.”
2. **System Response:** The chatbot retrieves a visual diagram of the Human Design Centers and offers a brief explanation for each center.
3. **Outcome:** The consultant uses the visual aid to explain concepts more clearly to their client.

# **Appendix**

## **Detailed Technical Approach**

This section details the technical architecture and components of the RAG-based chatbot system. The system comprises multiple stages: data ingestion from various sources, embedding generation for semantic search, vector database storage and retrieval, response generation using a language model, and the integration of these components. We also address security measures built into the design. Below is a breakdown of each component and how they work together:

### **Major Components & Functionality**

Our Retrieval-Augmented Generation (RAG) architecture combines information retrieval with a generative model to overcome the knowledge limitations of standalone large language models (LLMs). Our system is composed of several major components, each with a distinct role:

- **User Interface & API** – Captures user queries (via chat UI) and displays the chatbot’s responses. It sends the user’s query to the backend and presents the generated answer

back to the user. This component might be a web interface or messaging integration that the user directly interacts with.

- **Query Processor** – Receives the raw user query and pre-processes it (e.g., cleaning/parsing) if necessary. It also handles understanding context or conversation state if doing multi-turn dialogue. The processed query is then forwarded to the retrieval component.
- **Retriever** – Finds relevant information for the query from the knowledge base. It encodes the query into a vector and searches an index of document embeddings to retrieve the top-matching documents. This is implemented with a vector similarity search (e.g., using cosine similarity on embedding vectors) to identify the most relevant pieces of information. The retriever’s goal is to return a small set (we will experiment with varying top k, but for now we are using 5 as a default) of documents or passages that are likely to contain the answer to the query.
- **Knowledge Base (Document Store)** – A storage of the domain knowledge the chatbot can draw from. This consists of a corpus of text documents, webpages, FAQs, and any relevant data (PDFs, Doc/Docs, images and audio field). All documents in the knowledge base are pre-processed and encoded into embedding vectors, which are stored in an index to enable efficient similarity search. We plan to use a vector database for this index, as it allows fast nearest-neighbor searches in high-dimensional embedding space.
- **Generator (LLM)** – A generative language model that produces the final answer. It takes the user’s query along with the retrieved context documents as input, and generates a coherent response. The LLM is “augmented” with external knowledge by providing those retrieved documents as part of its prompt/context. This ensures the model’s answer is grounded in the up-to-date information from the knowledge base, rather than just its static training data. We will use a pre-trained LLM–GPT–via an OpenAI API.
- **Response Composer** – (Part of the generation step) Integrates the LLM’s output with references to the source documents. For transparency and user trust, the system will cite which documents or sources were used to derive the answer. This component attaches source identifiers (e.g., document titles or IDs) to the answer, so the user can verify information if needed. It may also apply some post-processing to format the answer or truncate it if too long.

The overall functionality is that when a user asks a question, the system retrieves relevant knowledge and feeds it to the LLM to generate a contextually-informed answer. This approach lets the chatbot provide accurate, up-to-date answers (since it can pull in recent or specific data) while still leveraging the fluency of the LLM for natural language responses. By partitioning the system into these components, we ensure a modular design where each part can be optimized or replaced independently (for example, swapping in a new retrieval algorithm or a different LLM

in the future).

## Interfaces Between Components

Each component communicates through well-defined interfaces to ensure modularity:

- **UI to Query Processor:** The user interface sends the user's query (as text) to the backend via a REST API call or WebSocket. The interface can be a simple API endpoint like POST /ask with the query in the request body. The Query Processor receives this request, extracts the query string, and returns a standardized query object or text string for the next stage.
- **Query Processor to Retriever:** The interface here is a function or API call such as retrieve(query) -> list<documents>. The Query Processor hands off the user's query (possibly cleaned or tokenized) to the Retriever component. The Retriever's interface expects a query string or vector and returns a list of relevant documents/passages (each including content and metadata like title or ID).
- **Retriever to Knowledge Base:** The Knowledge Base (KB) provides an interface for searching its indexed documents. For example, the Retriever might call an index method: index.search(query\_vector, top\_k) -> list<doc\_ids>. Internally, the vector index (in the Knowledge Base component) will perform a similarity search and return the IDs of the top-K matching documents. Then the Retriever fetches the full content of those documents via another Knowledge Base interface call (e.g., get\_documents(doc\_ids) -> list<documents>). This separation allows the Knowledge Base to manage storage details while the Retriever handles search logic.
- **Retriever to Generator:** Once relevant documents are retrieved, they are passed to the LLM. The interface could be a call like generate\_answer(query, context\_documents) -> answer. In practice, this means formatting a prompt for the LLM that includes the user's question and the retrieved text (for example, by concatenating them or using a template prompt saying "Using the information below, answer the question..."). The Generator component exposes an interface that accepts this composite prompt and returns the generated response text.
- **Generator to Response Composer:** The raw answer from the LLM is handed to the Response Composer along with the identifiers of the documents that were used. The interface here might be compose\_answer(raw\_answer, source\_docs) -> final\_answer. This component will attach citations or source attributions to the raw answer. The final formatted answer is then produced.
- **Backend to UI:** Finally, the composed answer (which includes the answer text and possibly structured references) is returned to the user interface. Using a REST API, the response to the POST /ask contains the answer and reference data (e.g., in JSON format).



The UI then displays this to the user in a readable form (rendering citations appropriately, perhaps as footnotes or hyperlinks).

## Data Storage and Schema

The primary data storage in this system is the **Knowledge Base**, which includes two parts: the **document store** and the **vector index (and also metadata)**.

- **Document Store:** This stores the actual content of all source documents and their metadata. The schema for each document could include fields such as:
  - doc\_id (unique identifier for the document)
  - title or source\_name (human-readable name of the document/source)
  - text\_content (the full text or text excerpt of the document)
  - embedding (vector representation of the document's content, for retrieval)
  - metadata (any additional info like date, author, tags, or permissions)
- The documents can be stored in a database or simply as files on disk with an index mapping. In our design, a lightweight approach is to store them in a Qdrant database, since the primary access pattern will be via the vector index.
- **Vector Index:** For efficient similarity search, all document embeddings are stored in a vector index structure. At a high level, the vector index can be thought of as a table with rows corresponding to documents and columns representing dimensions of the embedding. A separate index (like an ANN – Approximate Nearest Neighbor index) is built on this data to allow rapid retrieval of nearest vectors to a query vector.

During the **knowledge ingestion phase**, we will populate this database. All relevant documents (say PDFs, webpages, text files) are processed by an **Embedding Service**. This service computes a fixed-size vector for each document's content using a transformer model or embeddings API. Each vector is then stored in the index alongside the doc\_id. We will also store the mapping of doc\_id to full document content in the document store.

In addition to the knowledge base, there may be other storage needs (our stretch goal for the future):

- **Conversation History** (optional): If the chatbot supports multi-turn conversations, we might store the recent dialogue history for each user session (for context tracking). This could be kept in memory or a fast in-memory store like Redis. The schema might be as simple as a list of past Q&A pairs for each session.
- **Logging Database:** For analytics and debugging, we may log each user query and the system's response, including which documents were retrieved. This could be a simple log file or a database table with fields (timestamp, user\_id, query, retrieved\_doc\_ids, answer).

This helps in later evaluating system performance and errors.

### High-Level Schema Summary:

- **Documents**(doc\_id, title, text\_content, metadata...) – stored in the document store.
- **Embeddings Index**(doc\_id, embedding\_vector) – stored in vector index structure (the actual schema depends on the vector DB implementation; typically you load vectors and query by vector, not manual SQL).
- **ConversationHistory**(session\_id, turn\_index, user\_query, bot\_reply) – optional table if needed for multi-turn context.
- **Logs**(log\_id, timestamp, user\_id, query, retrieved\_docs, response) – optional, for development analysis.

The system does not have a traditional relational schema with many interrelated tables, since most of the heavy lifting is in the vector index. However, we assume the document corpus is static or updated periodically. If documents are updated, we would need to recompute their embeddings and update the index to keep the knowledge base in sync (this could be done offline or during off-peak hours as part of maintenance).

### Architectural Assumptions

Our architecture is based on several key assumptions:

- **Quality and Volume of Data:** We assume that the knowledge base contains sufficient information to answer the majority of user queries in our domain. The success of RAG is predicated on having relevant data to retrieve. If the data is sparse or irrelevant, even the best retrieval algorithm will not yield useful context for the generator. We also assume the size of the document corpus is manageable (few hundreds documents, not millions initially) so that indexing and search remain efficient on our hardware.
- **Embedding Model Efficacy:** We assume that we can obtain good quality semantic embeddings for our documents and queries. These embeddings need to capture the meaning of questions and documents such that similar content is near in vector space. We rely on the assumption that the chosen embedding model is well-suited to our domain's language; otherwise, retrieval performance may suffer.
- **LLM Context Window:** The chosen LLM has a limited context window (e.g., a few thousand tokens). We assume that by retrieving a handful of top documents (5-10 documents), the combined length of those documents plus the query will fit within the model's input limit. This influences how many documents we retrieve and how much of each document we can include. It's assumed we won't exceed these limits.
- **Real-Time Performance:** We assume the system can meet near real-time response

requirements (on the order of a few seconds per query). This means our design assumes adequate computational resources (CPU/GPU) to run the retrieval and the LLM inference quickly enough. We anticipate using a hosted LLM API or a reasonably sized model locally to achieve responses likely within a couple of seconds. If the model or indexing is significantly slower, the architecture might need adjustment.

- **Statelessness of Queries:** For the initial implementation, we assume each query-response is handled independently (unless we explicitly implement multi-turn memory). This simplifies the architecture because the Retriever can treat each user query in isolation. We do not assume long-term memory or user-specific learning in the basic design (aside from context provided explicitly in a conversation session).
- **Security and Privacy:** We assume that the data in the knowledge base is allowed to be used with the LLM (we received that permission from our client). If using an external LLM API, we assume no highly sensitive information is present (or if it is, that the API use is compliant with privacy requirements).
- **Failure Modes:** We assume that network calls (if any, e.g., to an LLM service) are reliable most of the time. We do plan for some basic timeout and error-handling (for instance, if the LLM API fails or times out, the system should catch that and possibly return a fallback message like “Sorry, I’m having trouble generating an answer.”). Similarly, if the Retriever finds nothing, we assume the system will handle it gracefully by informing the user it cannot find relevant info.

These assumptions were made to scope the project feasibility. For example, by assuming a static knowledge base, we avoid implementing complex streaming updates or real-time crawling of data. By assuming a suitable embedding model exists, we focus effort on integration rather than NLP research. Each assumption, of course, could be revisited if they prove false – e.g., if response times are too slow, we might need to choose a smaller model or invest in optimization.

## Alternative Architectural Choices (Pros & Cons)

In designing a RAG-based chatbot, we considered alternative architectures and approaches. Below are some key alternatives we evaluated and the trade-offs associated with each:

- **Fine-Tuned LLM (No Retrieval):** One alternative to RAG is to fine-tune a large language model on our project’s knowledge data so that it can answer questions directly without an external knowledge base at runtime. The advantage of fine-tuning is that the model potentially integrates the knowledge internally, leading to fast responses (just the model inference) and a simpler architecture (no need for a retriever or external index at runtime). However, this approach has significant downsides: retraining a model is time-consuming and expensive, especially for large LLMs. Moreover, a model that’s

fine-tuned once will become stale as data changes – it would need continuous re-training to stay up-to-date. We also risk the model “hallucinating” or giving false answers if it didn’t perfectly learn the facts, and it might be unclear on what basis it gives an answer (no straightforward way to cite sources). In short, while fine-tuning can make the system self-contained, it sacrifices flexibility and maintainability (the “knowledge update” problem and lack of transparency). Given our constraints, we chose RAG over fine-tuning so we can easily update the knowledge base and have the model always reference the latest information.

- **Pure Retrieval + FAQ (No Generation):** Another option was to build a system that doesn’t generate novel sentences but instead retrieves and returns existing text (like a smarter search engine). For example, upon a query, simply return the paragraph from the knowledge base that likely contains the answer, perhaps with some highlighting. The benefit here is avoiding hallucination entirely – the system only shows exact text from trusted sources. This is also simpler to implement; it’s essentially an information retrieval system with maybe some basic ranking. We considered this approach, but the downside is the user experience: raw paragraphs may be too verbose or not directly answer the question. Users often prefer a concise, synthesized answer rather than having to read through a document excerpt. Additionally, a generative model can phrase the answer in a conversational or user-friendly manner and can combine information from multiple sources if needed. Pure retrieval would struggle if the answer needs synthesis of multiple documents. Thus, while more predictable, the pure-retrieval approach was deemed less helpful for our end users’ needs.
- **Using Traditional Keyword Search vs. Semantic Vector Search:** Within the retrieval component, an architectural choice was whether to use a dense vector approach (semantic search) or a traditional sparse keyword search (like an inverted index using TF-IDF/BM25). A keyword-based system would be simpler in some respects and very transparent – it finds documents with matching keywords. This can be effective for certain queries, and is easier to debug (you can see which words matched). However, keyword search may fail when queries are phrased differently from the text or rely on synonyms and context. Semantic vector search was chosen because it excels at finding conceptually relevant information even if wording differs. The con of semantic search is that it requires computing embeddings and can be heavier computationally; plus, it may retrieve something that is topically similar but not actually containing the factual answer (which then the LLM might run with incorrectly). For our use case, we decided the improved recall of semantic search outweighs those concerns, especially since we will combine it with careful prompt design to mitigate errors. We will, however, keep the option to fall back to keyword search for debugging and verify results.
- **Microservices vs. Monolith:** We also considered how to structure the deployment of the system. One approach is a monolithic application where all components (UI, retrieval,

generation) run within one service or process. The other extreme is to have microservices (e.g., a dedicated embedding service, a separate vector server, a separate generation service). A microservice architecture could improve scalability – for instance, the vector database could be scaled independently or replaced by a managed service. It also adds flexibility in technology choices for each component. The downside is increased complexity in development and deployment: more services to manage, network calls between them and potential points of failure at service boundaries. Given our project's size and timeline, we lean towards a simpler monolith or a minimally tiered design (a front-end service and a back-end service at most). This will reduce overhead and make development/testing easier.

- **Choice of LLM Integration:** Architecturally, using an external API vs. an open-source local model is another consideration. An external API (like OpenAI) is convenient and offers powerful models, but it introduces external dependency, potential latency, and cost per request. A local model (running on our own hardware) gives us more control and possibly lower variable costs, but might require GPU infrastructure and might not be as capable as the cutting-edge API models. We decided initially to use a hosted API for the LLM to ensure we get high-quality generation out of the box. The pros of our choice: quick to implement and likely best quality answers; cons: need to budget for API usage and ensure no sensitive data goes to third-party servers. This trade-off was acceptable for the first version of our product.

Each alternative was carefully weighed. Ultimately, the chosen architecture (RAG with semantic retrieval and generative answer synthesis) was deemed the best fit for our requirements of accuracy, transparency, and maintainability. The alternatives provide useful perspective and remain possible future improvements (for instance, once the system is stable, we *could* experiment with fine-tuning the model on chat transcripts to further improve it, effectively combining RAG with some model tuning). For now, our design maximizes immediate utility and flexibility, while minimizing development risk.

## Data Ingestion

- **Sources and Formats:** The chatbot will ingest knowledge from two primary sources: **text-based documents** and **audio files** (specifically WEBM audio, though other audio formats can be supported with minor adjustments). text files may include reports, manuals, or any text-based resources. Audio files might be recordings of meetings, lectures, or user-provided voice notes that contain information we want the bot to leverage.
- **Text documents Processing:** To extract text from text-based documents, we will use Python-based processing libraries. Possible choices include PyPDF2,

pdfplumber, or PyMuPDF (fitz) – these allow reading content programmatically. The ingestion pipeline will iterate through each document, extracting textual content page by page. During extraction, we will also consider basic cleaning: removing headers/footers, OCR-ing scanned documents if needed (though that adds complexity), and perhaps segmenting the text by sections or headings if the file is structured (to retain some logical grouping).

- **Text Splitting:** Raw text from a document can be very large (many thousands of words). It's not efficient or practical to embed extremely long texts as a single chunk (there are token limits for embedding models, and semantically, we want finer granularity). Therefore, after extracting text, we will split it into chunks of appropriate size (*~500 tokens per chunk*). We will ensure chunks don't cut off in the middle of a sentence for coherence. Each chunk will then be treated as a separate "document" for the purposes of embedding and retrieval. We will keep metadata such as the source document name, page number, or section title along with each chunk (metadata will be stored in the vector database as payload). This metadata is useful for providing context (like citing the source) in answers.
- **Audio Transcription:** For audio (WEBM) files, we utilize the **OpenAI Whisper API** to transcribe speech to text. Whisper is an advanced speech recognition model by OpenAI that can handle multiple languages and accents, producing text transcripts with high accuracy. Using the Whisper API via OpenAI is straightforward; the audio file is sent to the API and a text transcript is returned. We will likely use OpenAI's Python client library for this. For example, the code to transcribe an audio file looks like:

'''

```
import openai

audio_file = open("meeting_recording.webm", "rb")

transcript = openai.Audio.transcribe("whisper-1", audio_file)
```

'''

- **Storing Raw Text vs. Direct Embedding:** We have a design choice – whether to embed the text immediately upon ingestion or store the raw text and embed on demand. Our plan is to perform embedding as part of the ingestion pipeline so that we populate the vector database right away. We will not store the raw text in a

separate database; instead, the vector will hold the text in payload. However, for backup or reprocessing, we might keep a simple storage (even as JSON or in filesystem) of the extracted raw texts and transcripts. This can help if we need to re-embed (for example, if we switch embedding models or want to tweak chunk size without re-reading PDFs).

In summary, the data ingestion stage **converts heterogeneous data (text documents, audio) into text chunks** ready for vectorization. By the end of this stage, we have a collection of textual chunks, each associated with some metadata (source identifier, position in source, etc.). These chunks will feed into the embedding generation process.

### Embedding Generation

- Once we have textual chunks from the documents and transcripts, the next step is to convert each text chunk into a numeric vector representation (embedding). We use OpenAI's **text-embedding-ada-002** model for this purpose. This model is part of OpenAI's embedding API offering and is well-suited for semantic search tasks due to its balance of high dimensionality and cost-efficiency.
- **Why text-embedding-ada-002?** This model produces 1536-dimensional vector embeddings for any given text input. These embeddings capture the semantic meaning of the text such that texts with similar content will have vectors that are close in the vector space. The model can handle quite large input text (up to around 8191 tokens, which is several pages of text), although we aim to keep chunks much smaller than that for relevance and performance. Another reason to choose ada-002 is its cost-effectiveness: as of OpenAI's pricing, it's one of the most affordable embedding models. It is also widely used and robust for a variety of content, making it a safe default choice.
- **Process:** For each text chunk obtained in the ingestion stage, we will call the OpenAI embedding API. This can be done via the OpenAI Python library.
- **Integration with Qdrant:** We will utilize it to generate embeddings easily, and then use the Qdrant vector store utility to store them.
- **Quality considerations:** The semantic quality of embeddings is crucial – if they are good, then relevant pieces of text will be retrieved for a query. text-embedding-ada-002 is a general-purpose model, so it should capture semantics of most content well. If our documents have very domain-specific jargon, we might consider domain adaptation in future (or using a different model), but initial tests

with ada-002 on a variety of content have shown it to be quite effective. We should ensure to normalize or handle any unusual characters in text (like large code blocks, if any, or tables) since those might not embed meaningfully. But given our use-case (mainly narrative text in PDFs and spoken words), we are in a good position.

The output of this stage is that **each text chunk now has a corresponding embedding vector**. These (vector, text, metadata) tuples will be inserted into the vector database. By doing the embedding generation at ingestion time, we front-load the computation and ensure query time is fast (the heavy lifting of embedding the entire corpus is done only once initially, not on every query).

### **Vector Storage & Retrieval (Qdrant)**

We employ **Qdrant** as the vector database to store all the embeddings and enable fast similarity search. Qdrant is an open-source vector search engine designed for high-dimensional data and provides efficient nearest-neighbor search for our embeddings.

- **Storing Vectors in Qdrant:** Each embedded chunk (vector) will be stored as a point in Qdrant. A point consists of the vector itself and an optional payload (which can be any JSON-like data). We will use the payload to store metadata and the original text chunk. For example, a payload might look like: {"text": "...chunk content...", "source": "DocumentX.pdf", "page": 5}. We also assign each point an ID (could be a simple incrementing ID or a UUID). All points will be stored in a collection within Qdrant. We might name the collection something like "documents" or a project-specific name. When creating the collection, we specify the vector size (1536) and the distance measure (we can use cosine or dot product for similarity since OpenAI embeddings are typically normalized; cosine similarity is a common choice).
- **Processing Retrieved Data:** The retrieved chunks' text will be pulled out to be given to the language model. We may also use the metadata (for instance, to tell the model the source or to later display the source to the user). If some of these chunks are very similar or redundant, we might consider filtering or merging, but initially the plan is to just take the top-k as separate context pieces.
- **Relevance Feedback (Potential Extension–stretch goal):** While not in the initial scope, we note that Qdrant supports filtering (e.g., we could search within certain documents or metadata criteria if needed). We could also implement a reranking



step where we use the language model to pick the most relevant of the retrieved chunks (if  $k$  is large). For now, we rely on the vector similarity's ranking which is usually strong.

**Why Qdrant?** Qdrant is chosen for its performance and developer-friendly features:

- It provides a straightforward REST and gRPC API, with official clients in Python (which we use) and other languages. This makes integration easier.
- Qdrant can be self-hosted or used as a managed service (Qdrant Cloud). For our project, during development we might run it via Docker locally, and for deployment we could either containerize it alongside our app or use the cloud service with an API key.
- Features like filtering, payload support, and even full-text indexing of payload (if needed) give us flexibility for advanced usage.
- It's open source and has a strong community, meaning we can get support if we run into issues and we're not locked into a proprietary system.

## Language Model Response Generation

With relevant document excerpts retrieved, the next stage is generating a final answer to the user's question. We use **OpenAI's GPT-4** model (accessed via API) to produce the response. GPT-4 is a state-of-the-art large language model known for its advanced reasoning, understanding of context, and fluent generation capabilities, making it well-suited to synthesize answers from provided information.

- **Prompt Construction:** The core challenge in this step is constructing an effective prompt for GPT-4 that includes the retrieved context and the user's query, and instructs the model to use them. A typical prompt format might be:

'''

System: You are an intelligent assistant helping answer questions based on provided documents. Use the given context to answer the question. If the answer is not in the context, say you don't know. Do not fabricate information.

User: Question: "<user's question>"

Assistant: (The assistant is given the following context to help answer the question)

Context:

1. "<text of retrieved chunk 1>"

2. "<text of retrieved chunk 2>"

... (maybe multiple chunks listed)

Using only the above context, answer the user's question in a helpful and concise manner.

'''

We might put the context into the system or user message. Another strategy is to prepend each chunk with something like [Document excerpt] and then the text. The exact prompt will be refined through testing. The key instruction we include is that the model should only use the provided context for answering and not rely on outside knowledge, to minimize any hallucination and ensure the answer is grounded in our data. We also tell it to say if it doesn't know or if the info wasn't found, as a fallback.

Because GPT-4 can handle a large prompt (the 8K version can handle about 8,000 tokens and there is also a 32K token version), we have some room to include several chunks of context. If each chunk is, say, 500 tokens and we include 5 chunks, that's ~2500 tokens of context, which is easily handled. The user query might be a few dozen tokens, and the answer could be a couple of hundred tokens, all within comfortable limits. We do have to remain mindful of the token limit to not exceed it, especially if we one day feed a bigger context or have a multi-turn conversation where previous Q&A are included. But initial single-turn usage is safe.

- **Calling GPT-4 API:** We will use OpenAI's ChatCompletion API to get GPT-4's answer.
- **Post-processing the Answer:** We will likely do minimal post-processing. At most, we might trim excessive detail if needed or ensure the answer is polite and formatted. GPT-4 is generally good at following the prompt formatting instructions, so if we ask for a concise answer, it usually complies. If the answer is supposed to include references, we might prompt it to include something like "[Source: Document X]" if possible. Alternatively, since our system can know which chunks were used, we might append a list of sources after the answer in the

UI level (not generated by GPT). This is an implementation detail: either have the model cite or have the system cite. Often, a safer approach is the system adding citations to avoid any chance the model mismatches them. For now, the focus is on the answer content itself; handling citation display can be an extension.

- **Multi-turn Conversations (stretch goal):** Our initial goal is to answer single questions with provided context. However, to make it a chatbot in the true sense, we will allow the user to have a conversation. This means the user might ask a follow-up question like “What about the second topic it mentioned?” which relies on the previous answer or context. To handle this, we can use a conversation history or memory.

### Why GPT-4?

- GPT-4 is chosen because of its superior capabilities in understanding context and generating coherent answers. Compared to its predecessor GPT-3.5 (text-davinci-003 or gpt-3.5-turbo), GPT-4 generally yields more accurate and nuanced responses, especially important in a setting where it needs to strictly use provided information. GPT-4 is less prone to make up facts when properly instructed, and can handle more complex queries. The trade-off is cost and speed: GPT-4 is slower and more expensive per token. However, since our use case often deals with detailed internal knowledge, accuracy is paramount. A single incorrect answer could undermine user trust in the system. If needed, we can use GPT-3.5 for faster responses in non-critical scenarios or as a fallback, but our primary target is quality, hence GPT-4.
- Additionally, GPT-4’s larger context window (and availability of a 32k context version) means it can handle a lot of information if needed – useful if a user asks a very broad question that pulls in multiple documents worth of info (though we’d likely summarize rather than feed extremely large context due to cost).

**Integration of Results:** Once GPT-4 generates the answer, that answer is sent back to the user via the interface. We will integrate this with the UI (Streamlit) such that after the user submits a question, the system does: embed query -> search Qdrant -> assemble prompt -> call GPT-4 -> show answer. This flow will be wrapped into a function or chain for clarity.

In summary, the language model stage takes the relevant bits of knowledge and **crafts a human-friendly answer**. It is the stage where “information retrieval” turns into “answer

generation.” By leveraging GPT-4 with a carefully designed prompt, we ensure the final output is not just a raw excerpt of text, but a synthesized, coherent answer addressing the user’s query directly, supported by the content in our knowledge base.

## Component Details and Interaction

In this section, we detail the internal design of each software component identified in the architecture and how they interact at a technical level. We also illustrate the sequence of operations using UML-style descriptions.

**User Interface / Frontend:** This will be a simple web application that allows the user to input queries and see responses. In design terms it will have a front-end hosting (Streamlit-based/Amazon Web Services (AWS) hosting) that will call to a backend API. We will design the UI to be minimalistic: a chat window where each user query and chatbot answer (with citations) appear as a dialog bubble. The UI will call the backend’s /ask endpoint and handle the JSON response. The UI design will also include a “loading” indicator while an answer is being generated, to let the user know the system is working.

**Backend Orchestrator (Controller):** This is the central piece of code that receives requests from the UI and orchestrates the call to Retriever and Generator. Conceptually, it corresponds to the Query Processor and overall flow management. In code, this might be an `ask_question(query)` function that ties everything together. It will:

1. Take the incoming query string.
2. If multi-turn context is supported, retrieve past conversation context for this user and prepend or otherwise incorporate it.
3. Pass the query to the Retriever and get back documents.
4. Pass the query + documents to the Generator to get an answer.
5. Pass the answer and sources to the Response Composer to format the final answer.
6. Return the final answer to the caller (UI).

This controller will handle exceptions (e.g., if the Retriever fails or returns nothing, it may invoke a fallback or return a polite “I don’t know” message).

**Retriever (and Embedding Index):** The Retriever in implementation will use a framework. We plan to use an embedding model to vectorize the query. For example, we might use `sentence_transformers` in Python: `query_vec = model.encode(query)`. Then use FAISS (if in-memory) or an API call to a vector to search for similar vectors: `results =`

`index.search(query_vec, top_k)`. The result is a list of document IDs with highest similarity. Then we fetch those documents from the document store.

- We'll implement the retriever logic in a class, say `DocumentRetriever`, with methods like `get_relevant_docs(query_text) -> List[Document]`. This method encapsulates the steps: embed the query, perform similarity search, and return `Document` objects.
- We will have a data class or structure for `Document` that includes at least id, content, and possibly metadata (like title). This allows passing around documents easily.
- The embedding model itself could be loaded once at service startup and kept in memory to embed queries quickly. Similarly, the FAISS index can be loaded at startup from a file and kept ready for queries.
- **UML Class Structure (simplified):** We might have classes like `KnowledgeBase` (with subcomponent `VectorIndex` and `DocumentStore`), and a `Retriever` class that uses them. For example:
  - `KnowledgeBase`: has methods `search_embeddings(query_vec, k)` and `get_doc(doc_id)`.
  - `Retriever`: has a reference to `KnowledgeBase` and an embedder model. It calls `embedder.encode(query)` and then `kb.search_embeddings(query_vec, k)` internally, then fetches each doc via `kb.get_doc(id)`.
  - `Document`: a simple data class with attributes `id`, `text`, `title`.
- By structuring it this way, if we switch out FAISS for another search backend, we only modify `KnowledgeBase`'s implementation. The `Retriever` code remains mostly the same.

**Generator (LLM Integration):** The Generator will be designed to interface with whichever LLM we choose. This will be done through an API call to GPT 4o. We will wrap this in a class or function such as `LLMGenerator.generate_answer(query, docs) -> string`. The steps inside will include:

Constructing the prompt: e.g., we might use a template like:

"You are an expert chatbot. Using the following information from our knowledge base, answer the user's question.\n\n

Context:\n{doc1}\n{doc2}\n...\n\n

Question: {user\_query}\n

Answer:"

- Where `{doc1}`... are the text of retrieved documents (possibly trimmed to fit).
- Sending this prompt to the LLM. Using OpenAI API, this means calling the

- `openai.Completion.create` or chat completion with the prompt.
- Receiving the generated output text from the model.

The Generator design will also handle token limits – if the retrieved content is too large, it might truncate or select key sentences. We might implement a simple strategy: each retrieved document is truncated to a certain number of characters or sentences before inclusion. Another strategy is to summarize documents first if they are long, but that might be out of scope; so, truncation is the likely approach.

The class `LLMGenerator` may also include some configuration like which model/engine to use, API keys, etc., to avoid scattering those details in the codebase. By encapsulating generation, we also ease testing.

**Response Composer:** This part of design takes the raw answer and the document references to produce the final formatted answer. We have decided to use a bracketed citation style (e.g., “[Source: Document Title]” or numeric references) in the text. Implementation-wise, after generation, we know which documents were retrieved (we have their IDs and titles). We can map those to a reference notation. A simple approach: use numeric indices for each doc in the answer (like [1], [2]) and then append a reference list in the answer with those numbers mapping to the doc titles or URLs. This keeps the answer text clean yet traceable.

For instance, if let’s say 3 docs were used, the final answer could end with something like:

Sources:

1. Title of Document A
2. Title of Document B
3. Title of Document C

We will ensure the composer does not modify the core answer content, just augments it with reference notations.

**Control Flow (Sequence):** To clarify the interactions, here is the sequence of operations when a user asks a question:

1. **User** enters a query in the UI (e.g., “What are the benefits of Human Desing?”) and hits send.
2. **UI** calls the backend API `ask_question(query)`.
3. **Backend Orchestrator** (Query Processor) receives the request. It creates a new session context if not already or attaches a session ID.

4. Then we call the **Retriever** component: `get_relevant_docs(query)`.
  - The Retriever uses the **Embedder** to convert the text query into a vector.
  - It then queries the **Vector Index** in the **Knowledge Base** with this vector to find, say, the top 5 most similar documents. The Knowledge Base returns their IDs and similarity scores.
  - For each returned ID, the Retriever fetches the document content from the **Document Store** (which could be within Knowledge Base as well).
  - The Retriever returns a list of Document objects (e.g., DocA, DocB, DocC).
5. We then call the **Generator**: `generate_answer(user_query, [DocA, DocB, DocC])`.
  - The Generator prepares the prompt using the text of DocA, DocB, DocC (their content might be concatenated or listed as context) along with the question.
  - It sends the prompt to the LLM (via API).
  - The LLM processes the prompt and returns a completion – e.g., a paragraph of answer text. The Generator captures this text.
6. We now call the **Response Composer**: `compose_answer(answer_text, [DocA, DocB, DocC])`.
  - The Composer assigns reference numbers or markers to DocA, DocB, DocC (for example, based on their titles or a predetermined ordering).
  - It inserts the reference markers into the `answer_text` at appropriate points.
  - It produces a final answer string with citations and maybe a “Sources” list.
7. It then returns this final answer string (or structured data with answer and source list) to the **UI** via the API response.
8. **UI** receives the answer and displays it to the user. The answer appears as the chatbot’s reply, including any citations or footnotes for transparency.
9. The interaction for that query ends. If this were a multi-turn scenario, we would also store the question and answer in the **Conversation History** for that session before finishing, so that context can be used for the next turn.

This sequence mirrors the typical RAG workflow of Query -> Retrieve -> Read/Generate -> Respond. Each component in our design corresponds to a step in this workflow, ensuring a clear separation of concerns.

## Diagram Illustrations

**Class Diagram (Simplified):** The system can be depicted with the following main classes and relationships:

- ChatbotController – orchestrates the process. It has associations with Retriever, Generator, and ResponseComposer.

- Retriever – uses KnowledgeBase. It calls methods on KnowledgeBase to search and fetch documents. It also has an Embedder (which could be a separate class or external API for encoding text).
- KnowledgeBase – contains VectorIndex and DocumentStore. In the diagram, KnowledgeBase has a composition relation with these, indicating it manages their lifecycle. VectorIndex provides search(vector, k) functionality. DocumentStore provides get(doc\_id) functionality.
- Generator – may contain a reference to an LLMClient or API wrapper. It has a method like generate(query, docs) and internally uses the LLMClient.
- ResponseComposer – a utility class with method compose(answer, docs).
- Document – a data structure (could be a class or just a dict) representing knowledge base entries. KnowledgeBase.get returns one or more Document objects.

### Relationships:

- ChatbotController -> Retriever (calls retrieve(query)).
- ChatbotController -> Generator (calls generate\_answer(query, docs)).
- ChatbotController -> ResponseComposer (calls compose\_answer(ans, docs)).
- Retriever -> KnowledgeBase (calls search\_index(query\_vec) and get\_document(id)).
- Generator -> LLMClient (calls submit\_prompt(prompt)).
- KnowledgeBase -> VectorIndex and DocumentStore (calls internal methods on them; this could be shown as composition).
- Both Retriever and KnowledgeBase might use an Embedder component.

This design ensures that each class has a single responsibility: Retriever only worries about finding documents, KnowledgeBase about storing and retrieving data, Generator about LLM communication, etc. It also implies we could replace the internals of one without affecting others.

### Sequence Diagram (Query handling):

- User -> UI: types query.
- UI -> Controller: ask\_question(query).
- Controller -> Retriever: get\_relevant\_docs(query).
- Retriever -> Embedder: encode(query) [return: query\_vec].
- Retriever -> KnowledgeBase: search\_index(query\_vec, top\_k) [return: doc\_ids].
- Retriever -> KnowledgeBase: get\_documents(doc\_ids) [return: docs].
- Retriever -> Controller: returns list of Document objects.
- Controller -> Generator: generate\_answer(query, docs).
- Generator -> LLMClient/LLM: prompt [LLM processes and returns answer\_text].
- Generator -> Controller: returns answer\_text.



- Controller -> ResponseComposer: `compose_answer(answer_text, docs)`.
- ResponseComposer -> Controller: returns `final_answer`.
- Controller -> UI: returns `final_answer`.
- UI -> User: displays answer.

Throughout this, if an error occurs (say LLM fails), exceptions would be returned up the chain.

## Notable Design Decisions

Some specific design decisions included in the above that are worth highlighting:

- We chose to encapsulate the vector search inside a `KnowledgeBase` class. This hides the complexity of FAISS or any index library from the rest of the code. Other modules simply ask `KnowledgeBase` for results.
- We decided to treat the retrieval and generation as stateless pure functions in the system. This makes testing easier (given an input, we expect a certain output). The state (like conversation memory) is kept separate (in history storage) and can be injected into the query when needed, rather than being baked into those components.
- For the LLM prompt, an important design element is prompt engineering. Although part of content, it's essentially configuration. We will document the prompt template and keep it configurable (in a config file or easily editable string to use in testing) so we can tweak it without changing code.
- The design also accounts for possible multi-turn extension (stretch goal): While initial implementation may handle one query at a time, the structure allows adding a mechanism to include previous Q&A as context. We would adjust the Query Processor to prepend the last user question and answer (or a summary of the conversation so far) to the new query before embedding or including it in the prompt for generation.
- We emphasize modularity for testing: Each component can be tested in isolation. E.g., we can unit test the Retriever by injecting a fake knowledge base with known documents and verifying that `get_relevant_docs()` returns the expected ones for a sample query. Similarly, we can stub the Generator (have it return a canned answer) to test the end-to-end flow without needing an actual LLM in the loop. This design is key to a robust testing strategy.

## Security Measures

We outline the measures for authentication, access control, and data security below:

- **Authentication:** We will enforce that only authorized users can access the chatbot. For a simpler implementation (like a demo via Streamlit), we will implement a password-based access: e.g., the app will require an access key that we distribute to

intended users. Additionally, the back-end will authenticate requests – ensuring that even if someone tries to query the API directly, they need a valid token or credentials.

- **Stretch Goal (beyond the scope of this class) Access Control:** Beyond just login, we will consider **authorization levels**. For example, if we have documents with different confidentiality levels, not every logged-in user should query all data. Qdrant can support tags or we could use separate collections for different data sensitivity. We can implement a filter on queries: the system will check the user's role/permissions and adjust the Qdrant query filter accordingly. For instance, if a certain document is marked "Managers only", the chatbot should not retrieve those chunks for a regular client's queries, but only for the consultant.
- **API Keys Protection:** We will be using API keys for OpenAI and possibly for Qdrant Cloud. These keys will be kept out of the source code and out of the client-side. We'll use environment variables or a secure config file on the server to store them. This will prevent accidental exposure of keys in code repositories.

## Statement of AI/External Sources Use

- The proposed system is based on the client's needs and requirements, my prior experiences, and my general knowledge gained from literature and industry practices in the fields of information retrieval, Retrieval-Augmented Generation (RAG), and chatbot systems. While I cannot cite any particular sources (since I did not utilize any specific papers, necessarily), I want to emphasize that many statements made throughout this document and project proposal are derived from my experiences, built upon prior explorations, research, and familiarity with the literature, foundational principles, and technologies used in the domain mentioned above.
- ChatGPT was used for research during the project proposal stage, which is also often referenced here, and was also consulted to finalize decisions regarding the system architecture and RAG pipeline. Additionally, ChatGPT served as a discussion partner and consultant, helping me analyze and map my prior experiences with RAG (including my summer internship and independent projects). It was instrumental in brainstorming which methodologies I had previously applied should be used in this case and which ones should be reconsidered, and, once I finalized the ideas, helped polish them and work together on best formulation of them.
- ChatGPT was also used to assist with formatting the writing and aggregating information across various files and documents, consolidating decisions that had been made, and refining the project proposal and midterm presentation.

- ChatGPT also played a role in debugging, brainstorming, structuring syntax, and outlining pseudo-code data processing flows.
- ChatGPT was directly utilized for structuring text, grammar, formatting, and spell-checking. Largely because **due to my dyslexia**, writing, spelling, and formatting written materials can be particularly challenging for me, and ChatGPT helped ensure clarity and correctness in my written work.
- To summarize everything stated above: the ideas, architecture, data pipeline, testing methodologies, requirements, and specifications used in this project still represent my original work. As mentioned earlier, the project is largely based on my past experiences and knowledge, as well as discussions with the client and consultations with research literature. However, ChatGPT was used as a **supplementary** tool and assistant in various aspects of the project. It did not generate the main corpus of work for me but rather aided in refining, structuring, and improving the quality of my output.

## Reflections

**Here are some things that i would like to note about my experience with this project:**

### **1. Embracing Comprehensive Responsibility**

This project marked the first time I was solely in charge of coordinating, designing, and managing a full-stack product of such significant scope. In prior experiences, especially with AI-driven or collaborative endeavors, I was often just one member of a team, able to lean on others when it came to implementing certain functionalities or ensuring the project stayed on track. Stepping into this role where I was almost entirely responsible for every decision, large and small, felt daunting at first. Yet, it was also deeply empowering. The sheer volume of decision points, from architectural choices to integration details, forced me to operate in a more methodical and strategic manner than ever before. I found myself carefully mapping out each stage, identifying dependencies between components, and setting clear goals for every sprint. This higher-level oversight fundamentally changed my perspective on what it means to “lead” or to “own” a project from end to end.

### **2. Reconciling AI Expertise with Real-World Implementation**

Personal Reflections on My RAG Implementation Journey

Looking back at my work implementing the FutureFrame AI-driven Human Design consulting assistant, I'm struck by how significantly my initial assumptions were challenged. My academic and professional background in artificial intelligence led me to believe that the AI components

(model architecture, fine-tuning, and prompt engineering) would represent the most substantial technical challenges of the project. This assumption was fundamentally incorrect; the AI aspects of the system turned out to be surprisingly manageable, largely due to the mature ecosystem of libraries, APIs, and frameworks now available. The OpenAI embeddings API greatly simplified vector generation. Even the LLM integration through GPT-4 was straightforward with well-documented endpoints and clear response handling. And what proved genuinely challenging was constructing the supporting infrastructure around these AI components. The data ingestion pipeline became unexpectedly complex, specifically through dealing with inconsistent source materials ranging from text documents to webinar recordings required robust parsing, cleaning, and chunking strategies. The retrieval component required significant tuning beyond just implementing vector search. I discovered that naive cosine similarity alone wasn't sufficient; we needed to implement hybrid retrieval combining semantic search with keyword matching to handle edge cases. Additionally, chunking strategies dramatically affected retrieval quality; naive character-based splitting often broke conceptual units, while more sophisticated paragraph-based approaches with overlap created redundancy issues.

On the deployment side, ensuring consistent performance under varying load proved more challenging than anticipated. I encountered issues with token rate limits, timeout handling, and graceful degradation that weren't covered in any AI course I'd taken. We also needed to implement robust logging and monitoring to track system behaviors in production, especially to identify cases where the model produced incorrect information despite high retrieval confidence scores. Testing was another unexpectedly complex area. Creating evaluation frameworks to assess retrieval accuracy and answer quality objectively required careful design. We developed test suites with known question-answer pairs from the Human Design domain and implemented quantitative metrics for context relevance that went far beyond the simplified examples from academic papers.

Combined, this experience fundamentally shifted my perspective on AI system development. I now understand that AI expertise alone is insufficient for creating production-ready applications. The real challenge isn't implementing the core algorithms—it's building reliable, efficient systems that connect these algorithms with real-world data and user needs. The engineering skills needed to handle error states, optimize resource usage, ensure security, and maintain system reliability are equally critical to the specialized AI knowledge.

Looking forward, this realization informs how I'll approach future projects. I'll allocate more resources upfront to data pipeline engineering, retrieval optimization, and integration infrastructure. I'll design more comprehensive testing frameworks that evaluate system behavior under realistic conditions.

#### **4. The Underestimated Complexity of Tooling and Configuration**

I've always been enthusiastic about diving into new libraries or frameworks, but this project tested how deeply I understood the trade-offs between them. The entire process of deciding which hosting platform to use, starting with Amazon Web Services, illustrates how critical it is to evaluate not just feature sets but also ease of integration and maintainability under tight deadlines. Similarly, the choice of libraries for retrieving, embedding, or vector indexing had ramifications on performance, deployment costs, and even how easily I could debug certain issues. In the past, these “tooling” or “configuration” questions may have felt secondary compared to AI modeling, but they turned out to be pivotal. In fact, picking the right (or wrong) tool sometimes had a more pronounced impact on productivity and the final product's stability than any modeling tweak could.

Also, perhaps the most surprising element of this entire journey was just how many errors or incompatibilities surfaced during seemingly routine tasks, particularly during deployment to a live web environment. I often found that the code which ran smoothly on my local machine would behave quite differently once deployed, forcing me to become much more meticulous about both debugging and logging. There was a point where I felt that 20–30% of my total debugging happened during the final hosting phase, when I was convinced my codebase was already “finished.” So much so, this scenario taught me the importance of adopting continuous integration and continuous deployment (CI/CD) mindsets, even in smaller projects, and how hosting can serve as an invaluable final validation step that uncovers issues missed by local tests.

### **Lessons Learned**

#### **1. Managing LLM Variance and “Randomness”**

One of the standout takeaways from this project was confronting the inherent unpredictability of large language models, even when using advanced techniques like retrieval-augmented generation (RAG). Despite calibrating temperature settings, carefully constructing prompts, or providing curated contexts, I found out that there remain occasions when the LLM behaves erratically. The sheer randomness in responses (especially when a query sits at the edge of the model's understanding), I think, can frustrate both developers and end-users. It became clear that while domain-specific data and well-tuned parameters help reduce spurious outputs, they don't fully eliminate the underlying “probabilistic DNA” of these models. Hence, I counseled for myself that building robust user interfaces and fallback mechanisms (such as disclaimers or suggestion prompts) is crucial for mitigating users' potential confusion and maintaining trust.

#### **2. The Value of Team Collaboration**

Another key lesson was how working largely in isolation, though beneficial for honing my independent problem-solving skills, made me crave the synergy that comes from a diverse

team... In past collaborations of mine, there was always someone else to bounce ideas off of, someone with a different area of expertise who could propose an alternative library or catch a bug I overlooked. Going forward, I realize that, while solo projects are great crash courses in self-reliance, I definitely value the collective learning and shared accountability that a team brings.

### **3. Rapid Onboarding to Web Hosting and Deployment**

Finally, I learned just how steep the learning curve can be when it comes to web deployment if you've never done it before. From domain configurations and SSL certificates to ensuring that your backend and frontend can handle real-world traffic, the challenges were plentiful. In my case, initial attempts at hosting via AWS gave way to solutions using Streamlit, partly because Streamlit integrated smoothly with the Python-based environment and custom AI code I had already written. The process revealed that deployment isn't just about "going live" but is itself a powerful form of testing (as I mentioned before). Each deployment iteration revealed new edge cases: missing dependencies, version mismatches, or front-end rendering quirks... (hopefully now I have a complete set of requirements in my txt:) By the end, I realized that web hosting is not just a one-time step but a continuous process of refinement and there is definitely much more for me to learn about it.