

## Step-by-Step Plan for the “Veterans First” Claims AI Model

### Step 1: Define Objectives and Scope

Begin by clearly defining the project’s goals and boundaries. The objective is to create a “Veterans First” AI agent that helps navigate VA benefit claims and appeals, cutting through bureaucratic red tape. It should assist with filing claims, understanding decisions, and appealing denials, with a focus on U.S. Army veterans (active duty). Key capabilities will include: answering questions about veteran benefit laws and procedures, guiding users through claim filing steps, summarizing complex documents (e.g. decisions or regulations), and generating helpful text (like appeal letters or explanations). The agent should function in a chat interface, able to handle multi-turn Q&A, and use tools or external knowledge when needed. Defining these requirements guides what data and features the model will need.

### Step 2: Identify and Gather Comprehensive Data

Next, determine what data is required and how to collect it. This model must be grounded in the full range of laws, regulations, and documentation relevant to veterans’ benefits. Important data sources include:

**U.S. Code and VA Regulations:** Gather the laws governing veterans’ benefits (primarily Title 38 of the U.S. Code and Title 38 of the Code of Federal Regulations). These lay out eligibility rules, claim adjudication standards, rating schedules, and appeal processes. For example, 38 CFR Part 3 covers adjudication rules for compensation and pensions, Part 4 contains the disability rating schedule, and Parts 19–20 cover the Board of Veterans’ Appeals procedures. The eCFR website or official VA regulations publications can be used to obtain the entire text of relevant CFR sections. Ensure you have the latest versions since regulations are updated periodically.

**VA Policy and Procedure Manuals:** Acquire the VA’s internal manuals that describe how claims are processed. The primary one is the M21-1 Adjudication Procedures Manual, which is a comprehensive guide for VA staff on developing and adjudicating claims. This manual covers the end-to-end claims process (from application through decision) and is essential for understanding internal VA procedures. Similarly, the Higher-Level Review (HLR) guidelines (potentially in a newer M21-5 manual or policy memos) should be gathered to cover the appeals modernization process. These guides might be available through VA’s online portal (KnowVA) or via FOIA requests. If public access is restricted, look for secondary sources (e.g. veterans’ service organizations training materials or summaries) that detail these procedures.

**VA Claim Filing Guides and SOPs:** Collect any standard operating procedures, forms, and instructions for filing claims and appeals. This includes VA’s official how-to guides (for example, VA’s website instructions on “How to File a Disability Claim” and how to request an HLR), form instructions (like those for VA Form 21-526EZ for claims or VA Form 20-0996 for HLR), and checklists for fully developed claims. Veterans Service Organizations (VSOs) often publish service officer guides – for instance, the DAV or VFW guides – which summarize claim filing steps and evidence requirements in plain language. These can be invaluable data sources to train the AI to give step-by-step guidance.

**Legal Opinions and Case Decisions:** Include a broad range of decisions/opinions to teach the AI

how laws are applied in real cases. Two key sources are:

Board of Veterans' Appeals (BVA) Decisions: These are appellate decisions where Veterans Law Judges apply laws and regulations to individual cases. They illustrate how rules are interpreted and provide precedents on various claim scenarios. BVA decisions are public; the VA maintains a search tool and dataset of BVA rulings. For instance, the VA's data portal provides access to BVA decisions (with personal details redacted) searchable by keywords or docket numbers. You should plan to gather a substantial sample of BVA decisions, especially those relevant to Army veterans' claims (e.g. decisions involving Army service records, common claim issues, etc.). These can be collected via the VA's FOIA Reading Room or the [index.va.gov](#) search interface.

Discharge Review Board (DRB) Decisions (Army): Since the focus is on Army veterans, obtain decisions from the Army Discharge Review Board. The DRB reviews applications from former service members seeking to upgrade or correct their discharge status, which can affect benefit eligibility. The Department of Defense has an electronic reading room where DRB and Board for Correction of Military Records decisions are posted for each branch. For the Army, thousands of DRB decisional documents from 1998 to present are available (the Army Review Boards Agency site lists yearly folders with cases). These decisions will show how discharge issues are addressed and what "liberal consideration" policies or criteria are applied. Download a representative set of Army DRB rulings (ensuring PII is redacted) for the model's knowledge base.

Court Opinions and Other References (optional): If feasible, include higher-level court rulings (like CAVC – Court of Appeals for Veterans Claims – decisions) or notable GAO/IG reports on VA claims processing. These can add context on common errors or required standards. Additionally, any statistical reports or training materials on claims (for example, VA annual benefits reports, or the Veterans Benefits Administration's Program Guides like PG-21-2 which compiled policies) could be useful to ensure coverage of edge cases and legacy rules. However, prioritize the core sources above to keep the dataset focused.

For each data source, use official repositories where possible (e.g. eCFR for regulations, VA's websites for manuals, FOIA reading rooms for decisions). Plan the logistics of data gathering – for instance: downloading PDFs or text of regulations, writing scripts to scrape decisions from the BVA and DRB reading rooms, or using available bulk data downloads. Ensure all data is stored in text format for easy processing. This step may take significant time (several days) due to the volume of material – but it's crucial to "locate it all" so the AI has a complete knowledge base of the veterans' benefits domain.

### Step 3: Data Processing and Preparation

Once the raw data is collected, the next step is to prepare and organize it for modeling. This involves cleaning the data and structuring it in a way that an AI can learn from or retrieve from effectively:

Text Cleaning: Convert all documents into clean text. For PDFs or scanned files, use OCR if needed to extract text. Remove headers, footers, page numbers, and any unrelated content (e.g. web navigation text from scraped pages) so that only the meaningful content remains.

Redact or remove any residual personally identifiable information that might be present in case decisions (the DRB/BVA documents should already have PII redacted). Ensure uniform encoding and fix any textual errors from OCR. The goal is to have a corpus of plain text files for each source (laws, manuals, decisions, etc.).

Segmentation and Chunking: Break down large documents into logical sections or “chunks.” Long texts (such as an entire CFR part or a 30-page manual chapter) should be split into smaller units like sections, paragraphs, or subtopics. This is important both for training (so that fine-tuning examples can be reasonably sized) and for retrieval (so that the model can fetch specific relevant pieces). For example, you might split the M21-1 manual by its sections on claims, evidence, appeals, etc., and split lengthy BVA decisions by their issue discussions. Each chunk should be small enough (perhaps a few hundred words) to be used as context later. Maintain references or metadata for each chunk (like source name and section) so you know where any given text came from.

Annotation or Structuring (if needed): Consider organizing the data with labels or categories. For instance, tag each text chunk by its type ("Regulation", "Manual", "BVA\_decision", "DRB\_decision", etc.) and perhaps keywords (e.g. "PTSD claim", "Agent Orange", "discharge upgrade") if you plan to use metadata in retrieval. This can help if you want to filter or prioritize certain sources in responses. While not strictly necessary, a light taxonomy can improve the system’s clarity.

Creating Q&A or Instruction Examples: Since the plan is to fine-tune the model, you will need a training dataset of prompt-response examples. Leverage the collected texts to craft training pairs. For example, from a section of the CFR, you can create a question like “What are the criteria for service connection for PTSD?” and use the CFR text (simplified) as the answer. From a BVA decision, you might form a question about why a claim was denied and provide the decision’s reasoning as the answer. You can also include procedural questions (e.g., “How does a veteran file a Higher-Level Review?”) with answers drawn from the manuals or guides. Aim to create a diverse set of examples covering definitions, procedures, and reasoning. If possible, also include a few interactive examples (like a short dialogue) to fine-tune the AI’s conversational ability – for instance, a user question followed by the AI guiding them through an answer step-by-step. Each example should be formatted in the style you want the model to produce (preferably instructive and accurate).

Format Data for Training: Organize the fine-tuning data in the required JSONL format for OpenAI. Each entry will typically have a "prompt" (which could be a user question or an instruction) and a "completion" (the ideal answer, including any formatting or phrasing you expect the model to use). Ensure the completions are written in a helpful, factual tone (almost like an accredited claims agent’s answer). If your fine-tuning strategy is to teach the model tool usage or to always cite sources, include those patterns in the examples. For instance, you might have some completions that say “According to 38 CFR §3.X, …” to habituate the model to referencing regulations. Finally, split the dataset into training and validation sets so you can later evaluate performance on held-out examples.

By the end of this step, you should have: (1) a cleaned corpus of source texts organized for retrieval, and (2) a prepared fine-tuning dataset of Q&A/instruction examples. This ensures you’re ready for both retrieval-augmented answering and supervised fine-tuning of the model.

## Step 4: Build a Knowledge Base with Embeddings

To enable the model to use the vast information you gathered, create a searchable knowledge base using embeddings and a vector store. This is the backbone of a Retrieval-Augmented Generation approach. Here's how to do it:

**Choosing an Embedding Model:** Use a semantic text embedding model (OpenAI offers `text-embedding-ada-002`, which is well-suited for semantic search on textual data). This model will convert each chunk of text into a high-dimensional numeric vector that captures its meaning. By using embeddings, the AI can later find conceptually relevant information even if the user's question is worded differently than the source text.

**Embedding Text Chunks:** Take each chunk of your documents (from Step 3) and feed it into the embedding model to get its vector representation. Store these vectors along with an identifier or reference to the chunk (so you can retrieve the actual text later). For example, each chunk from a BVA decision would get an embedding and retain a link to that snippet of text. This process will produce a vector database of all your knowledge. You can use open-source vector databases (like FAISS, Pinecone, Weaviate, etc.) or even OpenAI's own hosted solution if using their platform – the principle is the same: you'll end up with a set of embeddings ready for similarity search.

**Indexing and Similarity Search:** Ensure the vector database supports fast nearest-neighbor search. This allows the system to find which stored chunks are most semantically similar to a new query. For instance, if a user asks about "benefits for Agent Orange exposure", the system should quickly retrieve chunks about Agent Orange presumptive conditions from the CFR or VA manual, even if the question didn't use the exact same phrasing. The semantic search will match on meaning, not just keywords.

**Testing the Knowledge Base:** Before integrating it with the AI model, do some dry runs on the embedding search. Manually query the vector store with a few sample questions and verify that the top results include the relevant sections of text. For example, test with questions like "What is a fully developed claim?" or "How do I appeal a denied PTSD claim?" and see if the right chunks (e.g. the manual's section on Fully Developed Claims, or a BVA decision on PTSD) come up. This will validate that your chunking and embedding are effective. If some queries fail, you might refine the chunk sizes or add keywords to certain chunks to improve retrievability.

At the end of this step, you have a semantic search knowledge base covering laws, procedures, and decisions. This will later be used to fetch supporting context so the AI's answers stay accurate and grounded in the official documentation.

## Step 5: Model Selection and Fine-Tuning

With data in hand, select an appropriate base model and fine-tune it for your task. The user suggested using "03" – which likely refers to OpenAI's GPT-3 model (e.g. `text-davinci-003`). This is a strong choice as it's a powerful language model and can be fine-tuned. Key steps in this phase:

**Choose a Base Model:** For OpenAI, the davinci series (especially text-davinci-003) is the most capable GPT-3 model, with high capacity and understanding. GPT-3.5 turbo or GPT-4 are even more advanced, but at the time of writing, fine-tuning is (or will soon be) available for GPT-3.5. You can proceed with text-davinci-003 for fine-tuning, given it's explicitly mentioned. This model already understands instructions well, which is advantageous for our Q&A/chat use case.

**Set Up Fine-Tuning Process:** Use OpenAI's fine-tuning API or the Playground interface to upload your prepared training dataset (from Step 3). The dataset should be in the proper JSONL format. You will follow OpenAI's fine-tuning procedure: \*“Collect a dataset of examples, upload it, then create a fine-tuning job”\*. In practice, this means running a command or API call that starts the fine-tuning, specifying the base model (davinci) and your training file. The OpenAI platform will handle the training run on their servers.

**Fine-Tune the Model:** Monitor the fine-tuning job as it trains on your data. With a substantial dataset (potentially hundreds or a few thousand examples covering various questions and scenarios), the fine-tuning might take some time. OpenAI may provide metrics like training loss; keep an eye on those to ensure the model is learning properly. You might do an initial epoch and evaluate, then possibly continue for a few epochs if needed. Avoid over-fitting – since your model should generalize to unseen questions, it's not necessary to force the model to memorize entire documents (the RAG approach will supply documents at runtime). The fine-tuning should primarily teach the model how to use the information: e.g. to follow the style of a helpful VA claims expert, to structure answers in steps if needed, to cite regs by number, and so on.

**Validate the Fine-Tuned Model:** Once fine-tuning is complete, test the model on your validation set or some sample questions. Does it follow instructions and give relevant answers? For example, ask it “Explain the steps to file a VA disability claim” or “Under what conditions can tinnitus be rated above 10%?” and see if it responds correctly (using the kind of content from your data). Since the fine-tuned model may have some knowledge baked in from training examples, it should perform better in the veterans domain than a generic model. Note any weaknesses: if it still hallucinates information or misses details, you may need to add more examples to the training set or adjust your approach (sometimes reducing the complexity of answers and relying more on retrieval can help).

**Iterate if Necessary:** Fine-tuning may be an iterative process. You might realize you need additional training data (e.g., you forgot to include examples about a certain benefit or a certain type of appeal). You can update the dataset and fine-tune further or start a new fine-tune. Remember that fine-tuning doesn't update in real-time with new regulations or cases – that's why we incorporate retrieval – but it does help the model align with the domain and desired style. For now, proceed with the best model you have fine-tuned on the data available.

By the end of this step, you will have a custom fine-tuned GPT-3 model (let's call it “VeteransFirstGPT-3”) that is tailored to veteran claims Q&A. This model can be invoked via the OpenAI API or Playground, and will serve as the brain of the AI agent.

## Step 6: Implement Retrieval-Augmented Question Answering (RAG)

Now, combine the knowledge base and the fine-tuned model into a single pipeline that can answer user questions effectively. This is the Retrieval-Augmented Generation (RAG) approach,

which means the model will retrieve relevant info from your data at runtime to craft its answers. The pipeline will work as follows:

Diagram: Basic retrieval-augmented generation workflow. The LLM (language model) first issues a query to the data source (knowledge base) to fetch relevant context, then consumes the user's prompt plus that context to generate a response.

1. User Query Handling: When a user asks a question or makes a request (for example: "I was denied benefits for sleep apnea, how can I appeal?"), the system will take this input and preprocess it if necessary (for instance, strip out any irrelevant text or split multiple questions). The query is then passed to the next stage.
2. Semantic Search in Knowledge Base: The system generates an embedding for the user's query (using the same embedding model from Step 4). It then performs a similarity search in the vector database of documents to find the most relevant pieces of text. For the sleep apnea example, the top results might include: a chunk from 38 CFR about service connection for sleep apnea, a section of the VA manual on evidence requirements for sleep apnea, and perhaps a BVA decision where a sleep apnea appeal was granted or denied. Retrieve the top N chunks (commonly 3–5) that are most relevant.
3. Context Assembly: Take the retrieved text chunks and assemble them into a context passage. You might format it as a brief citation list or as raw text context. For instance, you could create a prompt section that says: "Relevant Reference: [Excerpt from M21-1 about sleep apnea claims] ... [Excerpt from BVA decision] ... User's Question: [the user's query]". The exact formatting can be refined – the key is to provide the language model with the info it needs. Keep the combined context within model token limits (for GPT-3, a few thousand tokens max). If the retrieved text is very long, you may truncate or summarize it before inclusion.
4. LLM Response Generation: Send the composed prompt (user question + retrieved context) to the fine-tuned model (VeteransFirstGPT-3) to generate an answer. Because the model has both the question and authoritative context, it can formulate an accurate, specific answer. It will "augment" its knowledge with the provided snippets rather than relying on memory alone. For example, the model's answer might explain how the veteran can appeal the sleep apnea denial by filing a supplemental claim or HLR, referencing the criteria from the CFR that were provided. The fine-tuning will help it phrase the answer helpfully (e.g., "According to VA's guidelines, you should..."). Ideally, the model might even quote or cite from the context (especially if we included examples of that during fine-tuning).
5. Return Answer to User: The AI's answer is then returned to the user in the chat interface. It should read as a coherent, informative response. At this point, you might also include any source citations if desired (for instance, if the agent is meant to show which CFR section or BVA case it used, you could have the model include that in the answer).

This RAG loop ensures the AI stays grounded in the actual data you collected, providing up-to-date and specific answers. It mitigates the risk of the model "hallucinating" false info, because it leans on the reference text for facts. As OpenAI's docs note, \*“RAG improves a model’s responses by injecting external context into its prompt at runtime”\* – your system will do exactly that, retrieving knowledge on the fly.

**Handling Follow-up Questions:** In a chat setting, users might ask follow-ups (e.g., “What does that mean for my case?”). You should design the system to carry context forward. This can be done by appending a summary of prior Q&A or relevant details into the next prompt. The fine-tuned model should be able to handle some dialogue and recall (you can use conversation history truncation strategies within token limits).

By implementing this retrieval-augmented QA pipeline, your AI agent becomes a “veterans claims investigator” that dynamically researches its own knowledge base to answer questions accurately. This forms the core functionality for addressing user queries about claims and appeals.

#### Step 7: Enable Summarization and Document Generation Features

Beyond Q&A, the agent should assist with summarizing documents and generating custom text (like letters or explanations). These capabilities use the same model but possibly require slight tweaks in how prompts are constructed:

**Summarization of Long Texts:** Veterans might have lengthy decision letters or medical opinions they want summarized. To enable this, you can feed the text of a document into the model (possibly in chunks if it’s very long) and prompt the model to summarize it. For example: “Summarize the following BVA decision in plain language for the veteran:” followed by the decision text. The fine-tuned model, being domain-aware, should highlight key points (like what was decided and why). If the document is extremely long (exceeding token limits), implement a strategy: break the text into sections, have the model summarize each section, and then perhaps summarize the summaries. This iterative approach yields a concise final summary. You might incorporate this into your pipeline as a special mode (user could upload a PDF which gets OCR’d and then summarized, for instance). Ensure the summary is accurate by comparing it against the source for a few test cases.

**Paragraph and Letter Generation:** The AI can help draft documents such as appeal letters, claim statements, or evidence summaries. This is essentially a guided generation task. You’ll prompt the model with relevant context (via retrieval) plus an instruction like: “Draft a letter to the VA explaining the veteran’s disagreement with the decision and citing relevant evidence/laws.” The model would then produce a well-structured paragraph or multi-paragraph letter. Fine-tuning can assist here if you included example letter formats in training. For instance, the model could produce: “Dear Sir or Madam, I am writing to request a Higher-Level Review of my claim...” and proceed to cite facts (“According to 38 CFR...”) in support of the veteran’s case. This fully developed claim letter generation can save a lot of time for veterans or representatives. You may want the model to use a formal, polite tone and include necessary details (identification info, specifics of the dispute, etc.), so adjust your prompt or fine-tuning examples accordingly.

**Multi-Turn Guidance:** The agent might sometimes need to walk a user through a process step-by-step. Rather than one-shot answers, it could produce a step-by-step list or engage in a back-and-forth. For example, if a user doesn’t know what evidence to gather for a PTSD claim, the agent could respond with a series of questions or prompts (“Do you have a PTSD diagnosis from a doctor?”) to gather info and then give tailored advice. To handle this, incorporate

conditional logic or plan for multiple turns in your conversation design. The model itself can output questions, but you'll need to feed the user's answers back in context on subsequent prompts.

Implementing summarization and generation uses the same underlying model; it's mostly about crafting the prompt appropriately and possibly using the retrieval mechanism for any needed facts. Thanks to fine-tuning, the model should be adept at these tasks (since it has seen domain-specific examples). Always test the outputs – for instance, check that a generated letter is factually consistent with the file or scenario given, and that summaries don't omit critical details. Adjust prompt instructions as needed (e.g. "Focus on the outcome and reasons given by the Board" for summarizing a BVA decision).

#### Step 8: Integrate Tool Use and Multimodal Capabilities

To make the AI agent truly robust, you should allow it to use external tools or handle multiple modalities when necessary. In practice, this means extending the agent beyond pure text responses:

**Knowledge Base Updates:** One “tool” already in use is the knowledge vector search. You can further empower the model by allowing it to trigger searches. For example, if the user asks about something not in the current data (say a very new policy update from 2025), the agent could perform a web search or query an updated database. Implementing this could involve a wrapper where the model's output is parsed for an intent to search (using a special token or format), and then an API call is made to fetch new info, which is then fed back into the model for a final answer. This way, the agent remains current without requiring constant re-training.

**OpenAI Function Calling or Agents:** OpenAI's APIs now support function calling, which allows the model to request using a function (like a calculator, a search, etc.) by outputting a JSON snippet. You could set up a function for knowledge base queries or one for simple computations (for example, back pay calculations given start dates and rates). The model, upon recognizing a query that needs that, would invoke the function and get results that you then feed into its context. This makes the agent a “tool-using model” – it can extend its capabilities through defined tools, rather than trying to do everything in the prompt. Plan out which tools are relevant: likely candidates are a search function (for policy lookup), a database query (if you have a database of claim statuses or stats), or perhaps an OCR tool if you want the model to handle images of letters.

**Multimodal Inputs:** If you anticipate users might upload images or PDFs (like a scan of a denial letter or a DD-214 discharge form), incorporate an OCR and parsing step. For example, you can use a vision model or an OCR library to turn an image into text, then feed that text into the AI for analysis. Likewise, if you wanted the agent to interpret military service records or medical PDFs, you could integrate those as input modalities. This essentially means the agent can accept not just text prompts but also file uploads. Google's Gemini AI, for instance, is expected to be multimodal, combining text and other inputs. In your implementation, you can simulate this by manually converting non-text to text behind the scenes and then proceeding as normal. The key is to design the interface such that a user can drop in a document and ask the AI about it (e.g. “Summarize this letter” or “Does this decision mention my PTSD diagnosis?”), and the system will handle the necessary conversions and prompt the model appropriately.

Example – Tool Integration: Suppose a veteran asks, “Has there been any change in law for burn pit exposures this year?” If your static knowledge base is updated only to last year, the agent might not have that info. With a tool integration, the agent could call a live web search API. Upon finding a news update or a VA press release, it can pull that text and then respond: “Yes, in January this year Congress expanded presumptions for burn pit exposure-related conditions. …” Another example: a user could ask “Calculate my disability combined rating if I have 30% and 20%.” The agent could use a simple function to do the VA math (which is not a straight sum but a formula) and then give the result, rather than guess. These enhancements make the agent more interactive and practical.

Implementing tool usage requires extra coding (often using an AI orchestration framework like LangChain or custom logic to intercept model outputs). During development, test that the model knows when to use a tool. You might include hints in the system prompt like, “If the query is about new information beyond your data or requires calculation, use the appropriate function.” Keep the user experience seamless – the user just asks, and behind the scenes the model might be doing multiple steps. This multi-step, multi-modal ability will set your “Veterans First” agent apart, making it a versatile assistant rather than a static QA bot.

#### Step 9: Testing and Evaluation

With all components in place (model, retrieval system, and tools), conduct thorough testing. This ensures the system works end-to-end and meets the needs of veterans and claims representatives:

Functional Testing: Try a wide range of sample queries covering every aspect of the claims process. For filing stage: “How do I file a fully developed claim for PTSD?” For mid-process: “What does it mean when VA sends a duty-to-assist letter?” For appeals: “What are my options after an HLR denial?” Also test edge cases: “Can I get benefits for a condition diagnosed after discharge?” and branch-specific queries: “I was in the Army and got an Other-than-Honorable discharge – can the VA help me?” Ensure the agent provides correct and helpful answers each time, using the retrieved knowledge (check that it isn’t making things up). If you find inaccuracies, trace whether the retrieval pulled wrong context or the model misinterpreted it. This might lead to refining your search query formation or adding training examples for those cases.

User Scenario Testing: Simulate actual user sessions. For example, walk through a scenario where a veteran chats: they describe their situation, the agent asks follow-up questions, they answer, and then the agent gives advice. Evaluate if the conversation remains coherent and contextually aware. The agent should remember details the user shared (within the same session) – if it forgets, you may need to pass conversation history into the prompt. Assess the helpfulness: does the agent not only give correct info but also empathize or clarify as a human advisor would? Adjust the tone via the prompts if needed (you can prepend a system message like “You are a friendly VA claims assistant…” to reinforce the style).

Performance and Limits: Measure how long responses take, especially when using retrieval and tools. The system should ideally respond in a few seconds per query. If it’s slow (perhaps due to long vector searches or large context), optimize by reducing chunk size or pre-indexing more

efficiently. Also check the token usage – answers should typically remain within a few hundred tokens. If the model tends to run on too long, consider setting a max tokens limit or instructing it to be concise for certain answers. Conversely, for something like letter drafting, you may allow longer outputs. Find the balance based on testing feedback.

**Accuracy Checks:** It's important to validate the content of answers. You or subject matter experts (like accredited claims agents) should review the AI's responses for correctness. Does it cite the right regulation? Does it give sound advice (e.g. not telling someone to do something harmful to their claim)? Because this domain affects real benefits, accuracy is paramount. If any systematic errors are found (say the AI consistently misinterprets a certain regulation), update the knowledge base or add a special rule/prompt to address it. In some cases, you might "hard-code" certain critical pieces of knowledge rather than rely on the model (for example, a list of presumptive conditions might be better as a reference table the model can retrieve exactly).

**User Feedback Loop:** If possible, incorporate feedback. For a deployed agent, allow testers or early users to flag incorrect answers or ask questions the AI couldn't handle. Use this to improve the model – either by fine-tuning on those missed Q&As or by expanding the knowledge base. For instance, if a new BVA precedent decision comes out changing how sleep apnea is evaluated and users ask about it, you'd add that to the data and possibly fine-tune the model on an explanation of that change.

Testing should be iterative. It's normal to go back and tweak earlier steps (e.g., adjust how retrieval works or add more training data) as you discover issues. Only once the agent consistently produces reliable, helpful responses across a variety of scenarios should you consider it ready for deployment.

#### Step 10: Deployment and Front-End Integration (Gemini Chat App)

Finally, deploy the system and integrate it into a user-facing chat application. The user mentioned using Gemini for the front-end – presumably referring to Google's Gemini AI/chat interface or a similarly named platform – to provide the chat experience. The deployment steps include:

**Set Up Backend Service:** Package your AI pipeline (the fine-tuned model + retrieval + tools) into a backend service. This could be a cloud function, a web server, or a container running your code. The service will expose an API endpoint for the chat app to communicate with. For instance, you might create an endpoint /chat where the front-end sends user messages and receives the AI's replies. Within this service, implement the logic we outlined (receive query -> do embedding search -> call OpenAI model -> return answer). Make sure to secure the service (use API keys or authentication, especially since it will access potentially sensitive info and your OpenAI credentials).

**Integrate with Gemini Front-End:** Gemini likely provides tools or SDKs to build chat interfaces (Google's Gemini might be a powerful model itself, but here we are using it for the interface presumably). Design a chat UI that is intuitive for users. It should have a text input box, and possibly options to upload documents or select the type of request (Q&A vs. summarization, etc.). The UI will send the user's messages to your backend and display the AI's responses in a conversational format. Since the target users may be veterans who are not tech-savvy, keep the

interface simple and mobile-friendly (many might use a phone to chat). Use clear prompts in the UI, like “Ask the Veteran Claims Assistant” or “Upload a VA letter for analysis.” Also, integrate any needed disclaimers (e.g. “This AI provides information but not legal advice” if relevant).

**Real-Time Conversation Features:** Ensure that the front-end handles conversation context. A user should be able to ask follow-ups without repeating everything. You might maintain a session state on the backend that stores the last few interactions (and include them in each new prompt) for context continuity. The front-end can handle user session management (maybe via cookies or session IDs). Also, provide the ability to reset the conversation if needed.

**Testing in the Live Environment:** Before full launch, test the entire system through the front-end as a user would. This can reveal integration bugs (for example, encoding issues that scramble quotes, or timeouts if the AI call is slow). Optimize the experience – maybe add typing indicators while the AI is formulating an answer, and handle errors gracefully (if the AI fails to get an answer, the front-end could show “I’m sorry, I didn’t get that. Can you rephrase?” rather than just breaking). Test file upload flows if you have them, ensuring the file is properly sent to the backend and parsed.

**Monitoring and Scaling:** Once deployed, monitor usage. Track metrics like number of questions answered, response time, and any errors. Given that the model calls OpenAI, keep an eye on your API usage and cost. If usage grows, you might need to scale the backend or optimize calls (e.g., maybe cache certain embedding search results for common queries, etc.). Also be prepared for updates: as laws change or new data becomes available (for instance, new BVA decisions each month), plan a schedule to update the knowledge base and possibly fine-tune the model periodically so it stays current.

Deploying via a chat app like Gemini means veterans and advocates will be able to interact with the AI agent seamlessly, as if texting with a knowledgeable assistant. The result is a multimodal, embedding-powered, retrieval-augmented chat agent that can answer questions, summarize documents, and help draft claim-related text – a powerful tool to cut through VA’s bureaucratic maze.

#### Step 11: Maintenance and Ongoing Improvement

(Building the model is not a one-and-done task—maintaining accuracy and usefulness over time is crucial):

**Regular Data Updates:** The landscape of veterans benefits can change (new legislation, policy updates, yearly regulation changes, etc.). Schedule regular updates to your knowledge repository. For example, check quarterly for any updates to 38 CFR or VA manuals and incorporate those. Continuously fetch new BVA decisions (the [data.gov](#) index updates monthly) and add them to the vector store so the AI learns from the latest precedents. Similarly, update the DRB decisions annually as new ones are published.

**Re-Fine-Tune if Needed:** If significant new knowledge needs to be baked into the model (say a new type of claim or a major process overhaul like the introduction of a new appeal lane), consider fine-tuning a new version of the model with additional examples. Alternately, leverage the retrieval system to handle it without re-training by ensuring the new info is retrievable and

perhaps weighting it higher in search results.

**Monitor Performance and Feedback:** Deploy analytics or feedback mechanisms. For instance, allow users to rate answers or flag if something was incorrect or unhelpful. Use this feedback to identify areas for improvement. You might discover, for example, that users often ask about a topic you didn't anticipate – you can then add content for that. Or if the agent frequently uses too much legal jargon, you might adjust the prompt to simplify language. Make it a continuous loop of improvement.

**Scaling to Other Branches or Use Cases:** Currently the focus is Army full-time veterans. In the future, you could extend the model to cover nuances for other branches (Navy, Air Force, etc. have their own DRBs and sometimes slightly different processes for records). The core laws are the same, but you might add branch-specific data if demand arises. The modular design of your system (with data and retrieval) will allow adding these easily – e.g., include Navy DRB decisions in the vector DB and perhaps fine-tune on any unique terms (like Navy rates or Marine Corps specific issues).

**Ethical and Legal Compliance:** Ensure the AI remains a tool for assistance and does not give detrimental advice. Always remind users (perhaps via the UI or occasional AI message) that it is providing information based on VA sources, and encourage them to consult a human expert for complex cases. Adhere to privacy – do not log any personal details users share beyond what is necessary, and secure any uploaded documents. Since the model deals with potentially sensitive personal and medical topics, maintaining confidentiality and ethical use is critical.

By following these steps – from data gathering to deployment and maintenance – you will have built a custom AI model for veterans' claims that leverages the full body of relevant knowledge (laws, rules, procedures, and precedent cases) and provides interactive, helpful assistance. This "Veterans First" model, powered by GPT-3 and enhanced with embeddings, RAG, and tool use, aims to empower veterans to navigate the VA benefits system more effectively than ever before. Good luck with your implementation, and thank you for creating a tool to help those who served!