

# Assistance Chatbot Exercise - Retrieval-Augmented Generation

## Table of Contents

Use Case	2
<b>Exercise Instructions</b>	<b>2</b>
Strategies to Employ	3
Implementation Details	6
Sample Repositories to Reference for Reusable Code	8

## Use Case

An assistant generative AI agent to provide a community of users with guidance, precedents, and training for their day-to-day roles. The assistant will answer questions related to its own training or procedural documentation.

## Exercise Instructions

1. Your goal is to build the back end of a successful chatbot Assistant using generative AI. You do not need to build the front end, although if it's easier to demonstrate, feel free to build it. Alternatively, you can show the success of your agent in your IDE of choice.
2. You will be provided with full guidance on how to complete this exercise within the upcoming sections. The intent is to help provide you with enough instruction or link-referenced options to make decisions and learn how to proceed.
3. You will use Microsoft Azure as your cloud service provider. Please create a free account to build your solution.
4. Please use the langchain framework and modules to build your chatbot.
5. Helpful links in the Implementation Details section provide you educational material to help you with the exercise.
6. The Sample Repositories to Reference can alternatively provide a jump start as well, although if utilized you'll need to determine how / where to customize the repo code to employ the strategies defined.
7. You can alternatively utilize code in the sample repositories to augment something that you personally build through the other reference links.
8. You may select the content that you will source for this exercise; please consider the target audience for the content covered by your generative AI assistant / agent, and explain why you chose the audience and related content.
9. Discuss strategies and implementation alternatives that you believe could have also worked, or even worked better for creating a helpful assistant.
10. Please ask questions as often as needed.
11. Feel free to use commercially available LLMs to help you close gaps in code. It's now expected that you utilize LLMs to improve your efficiency.

## Strategies to Employ

Implementing a Retrieval Augmented Generation (RAG) approach to optimize search and retrieval must consider the characteristics of both:

- Source Content, and
- User Queries

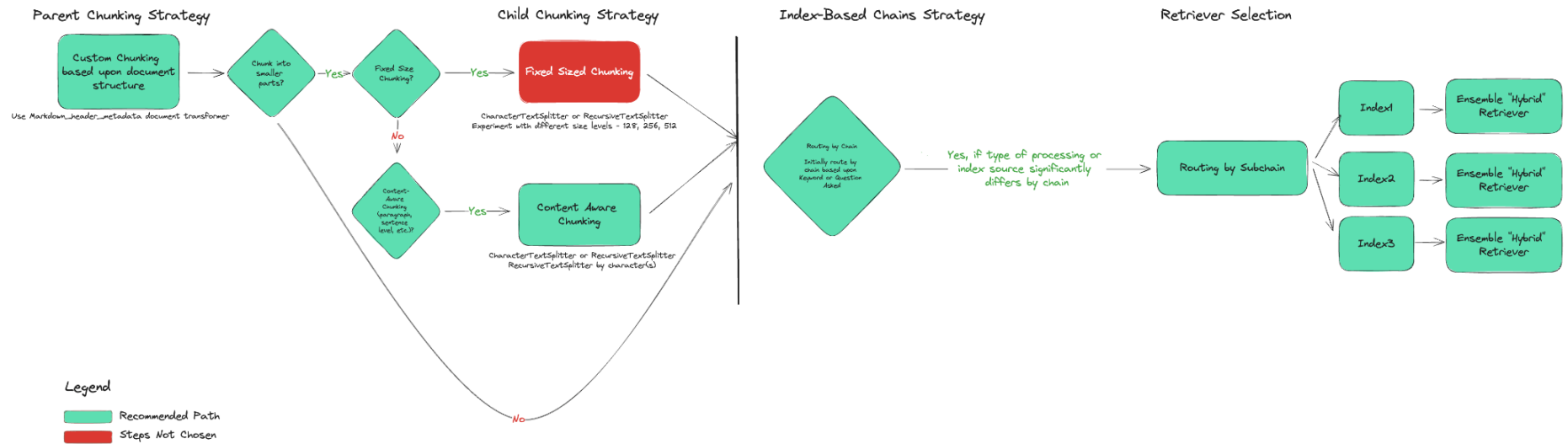
to determine the proper chunking, indexing, and retrieval strategies.

Below, we examine the specifics for each related to the above Use Case.

Category	Characteristics	Strategies to Employ
<b>Source Content</b>	<b>The Source Content</b> possesses a hierarchical structure, such as as: <ul style="list-style-type: none"><li>• Category</li><li>• Subcategory</li><li>• Chapter</li></ul>	<b>Chunk documents</b> aligned with the structure - either at the <b>section level</b> or the <b>paragraph level</b> .  <b>Generate a list of searchable keywords</b> representing each chunk.  <b>Attach metadata to each chunk</b> to represent <b>searchable keyword combinations</b> related to the chunk, along with a Title field and a link to the chapter ( <i>i.e. parent document</i> ).

<b>User Queries</b>	<p>More than half of the expected queries will require <b>content that spans multiple Chapters, and may also span multiple Subcategories.</b></p>	<p><b>Index documents</b> to employ a mix of indices representing keyword combinations (<i>representing queried "topics"</i>) for co-locating otherwise disparate chapter content to support accurate and performant cross-chapter retrieval through index pruning.</p> <p><b>Implement the Ensemble Hybrid Retrieval strategy</b> to drive initially filtered results with BM25-based lexical search, then semantic search.</p> <p><b>Note:</b> utilizing the <b>hybrid search strategy will auto-invoke Reciprocal Rank Fusion (RRF)</b> to run parallel queries for the multiple keywords associated with the query, decreasing retrieval latency.</p> <p><b>Optionally, also test search semantic re-ranking</b> to alter the resulting "top k" retrieved results for the model to formulate the answer.</p>
---------------------	---	--

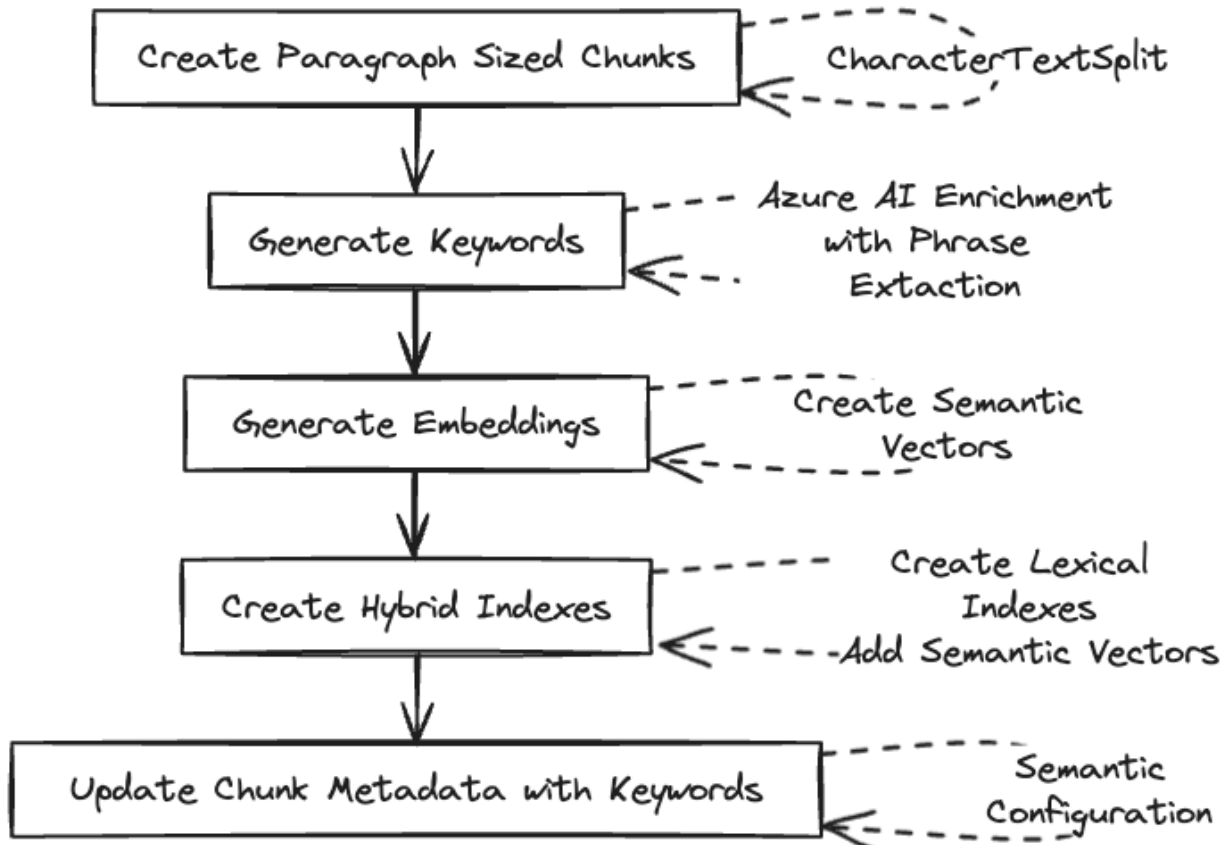
Given the above characteristics spanning both source content and queries, we employ an implementation pattern that invokes the following strategies reflected in **green**.



## Implementation Details

The following flow summarizes the sequencing for chunking and indexing, with references to implementation concepts:

### Chunking & Index Generation Flow



**We can further break down the sequenced concepts into more concrete implementation steps, with reference to Microsoft Azure and Langchain 'how to' documentation step details, as applicable:**

Step	Instruction	Extended Notes and 'How to' Reference Links
Create Paragraph Sized Chunks	<p>Split text into chunks/passages with a targeted paragraph size.</p> <p>Utilize "\n\n" as the character to split content should avoid splitting a paragraph that includes bullets.</p>	<p>Uniform chunks facilitate fuller vector representation for semantic relevance to specific queries.</p> <p><b>Reference Link</b>  <a href="#">Langchain Character Text Splitter</a></p>
Generate Keywords	Analyze text and identify important topics, entities, ideas for tags	<p>Could use various NLP techniques like topic modeling, named entity recognition, concept extraction. Useful for classification, metadata generation, search optimization.</p> <p><b>Reference Links</b>  <a href="#">AI Enrichment to Generate Metadata (using Skillsets)</a>   <a href="#">Cognitive Skill - Key Phrase Extraction</a></p> <p><b>-OR -</b>   <a href="#">Docugami for Metadata Extraction</a></p>
Generate Embeddings	Map chunks into numerical vector space	<p>Represents semantic context. Allows similarity calculations for search matching.</p> <p><b>Reference Link</b>  <a href="#">Azure AI Search LangChain vector code sample</a></p>

Create Hybrid Indexes	Create a lexical index, and then add semantic vector fields to it.	<b>Reference Links</b> <a href="#">Create an Index in Azure AI Search</a> <a href="#">Add vector fields to a search index</a>
Update Chunk Metadata with Keywords	Store keywords from text with paragraph chunks as metadata.	Associate metadata with chunks to drive filtering and ranking.  <b>Reference Links</b> <a href="#">Enable Semantic Ranking</a> <a href="#">Configure Semantic Ranking</a>  <a href="#">Azure AI Search LangChain vector code sample</a>

## Sample Repositories to Reference for Reusable Code

[Azure-search-openai-demo](#)  
[chat-with-your-data-solution-accelerator](#)