# RAG + AZURE AI

**1.** Set up Storage Account

**2.** Store Table in Azure Table Storage



**3.** Azure AI Search ⟶ Load Document



🔷 **Azure AI Search + Embedding Flow**
1. **Embed docs** → use `text-embedding-3-small/large`.
2. **Create index** → define fields: `id`, `content`, `vector`, `metadata`.
3. **Store** → push text + vectors into the index.
4. **Query** → embed user query → run **vector similarity search**.
   - Metrics supported:
     - ✅ **Cosine similarity** (most common, default)
     - ✅ **Dot product**
     - ✅ **Euclidean (L2) distance**
5. **Return** → top results → feed into RAG.

Go to Azure AI Search — Import Data

(BMW Data)

Choose RAG / Keyword



Vectorize your text

Column to vectorize: Key → column to be indexed

Model deployment: text-embedding-ada-tj → Embedding Model

**Imp:** So here I will be assigning

Text Embedding Model & Index Type

**RAG** ···
ai-search-tj

---

- ✓ Connect to your data
- ✓ Vectorize your text
- ● Advanced settings
- ○ Review and create

**Advanced ranking and relevancy**

Semantic ranker uses deep neural networks to provide relevant results and answers based on semantics, not just lexical analysis. Learn more ⧉

☑ Enable semantic ranker

**Index fields**

Shows a preview of the index fields and allo...

✏️ Preview and edit

(Schedule indexing)

Schedule

| | |
|---|---|
| Once | |
| 5 minutes | |
| 10 minutes | |
| 30 minutes | |
| **Hourly** | |
| Daily | |
| Custom | |

Hourly ▾

---

④ Set env Variables for AI Search

— AI Search Acc. ai-seach-tj

— Index Name rag-tj

— AI Search API Key. * * *

---

Home > ai-search-tj

**≣✓ ai-search-tj | Indexes** ☆ ···
Search service

🔍 Search

⊕ **Add index** ▾   ↻ Refresh   🗑 Delete

🔍 Filter by name...

| Name ▾ | Document count ▾ |
|---|---|
| rag-tj | 50 |
| tj-index | 5 |

- 🌐 Overview
- 📄 Activity log
- 👥 Access control (IAM)
- 🏷 Tags
- 🔧 Diagnose and solve problems
- 🔗 Resource visualizer
- ∨ Search management
  - (≣✓ Indexes)

---

**Code Example**

```
# These variables configure the search service and index for retrieving documents

# Set the Azure AI Search service name
os.environ["AZURE_AI_SEARCH_SERVICE_NAME"] = "ai-search-tj"

# Set the Azure AI Search index name to query
os.environ["AZURE_AI_SEARCH_INDEX_NAME"] = "rag-tj"

# Set the Azure AI Search API key for authentication
os.environ["AZURE_AI_SEARCH_API_KEY"] = "0CVzv2rS1feYv99m1BNkvTvDKiRICZCILzcn
```

} very imp for Retriever

---

1. **Azure AI Retriever**

```
# Step 1: Initialize the AzureAI Search Retriever
# This retrieves relevant documents based on the user query from the Azure Search index
retriever = AzureAISearchRetriever(
    content_key="Answer",   # The key for the content field in the search results change it accordingly as per your data
    top_k=1,                # Number of top results to retrieve
    index_name="rag-tj"     # Name of the Azure Search index to query
)
```

2. **Prompt**

```python
# Step 2: Define the prompt template for the language model
# This sets up how the context and question will be formatted for the model
prompt = ChatPromptTemplate.from_template(
    """Answer the question based only on the context provided.
Context: {context}  # Placeholder for the context from the retriever
Question: {question}  # Placeholder for the user question"""
)
```

3. **LLM**

```python
llm = AzureChatOpenAI(
    azure_deployment="my-first-gpt",    # The name of your deployed model in Azure
    api_version="2025-01-01-preview",
    azure_endpoint=AZURE_API_CLIENT,
    api_key=AZURE_API_KEY
)
```

4. **Chain** ⟶ Context, Prompt, LLM

```python
# This chain will process the retrieved context and the user question
chain = (
    {"context": retriever , "question": RunnablePassthrough()}  # Set context using the retriever and format it
    | prompt                                                     # Pass the formatted context and question to the prompt
    | llm                                                        # Generate a response using the language model
    | StrOutputParser()                                         # Parse the output to a string format
)

# Step 5: Infinite loop for user input
while True:
    user_question = input("Please enter your question (or type 'end' to exit): ")
    if user_question.lower() == "end":
        print("Exiting the loop. Goodbye!")
        break

    # Pass input as a dict matching the chain keys
    response = chain.invoke({"question": user_question})
    print("Response:", response)
```

*Imp.* The order is imp in Chain

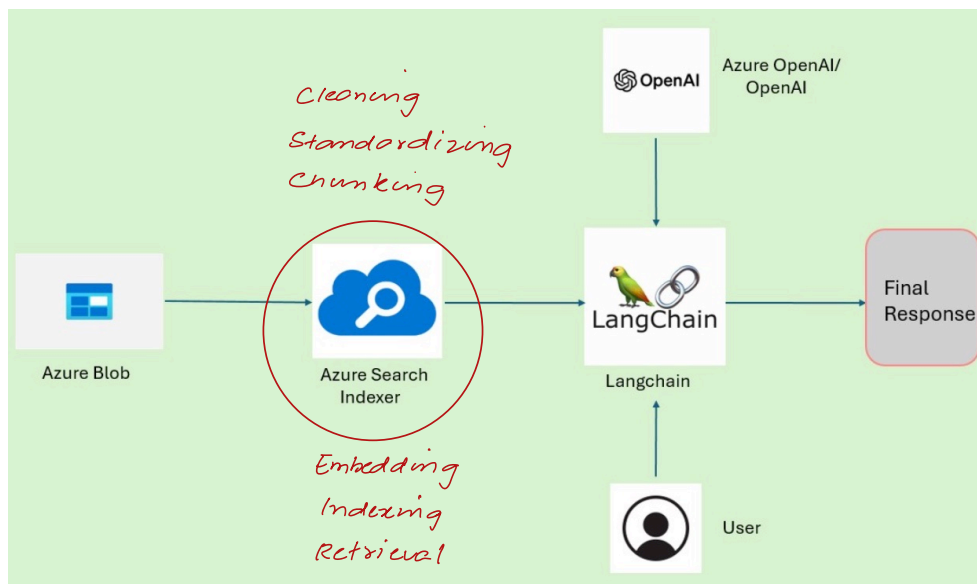✅ Rule of thumb:

Retriever → Prompt → LLM → Parser
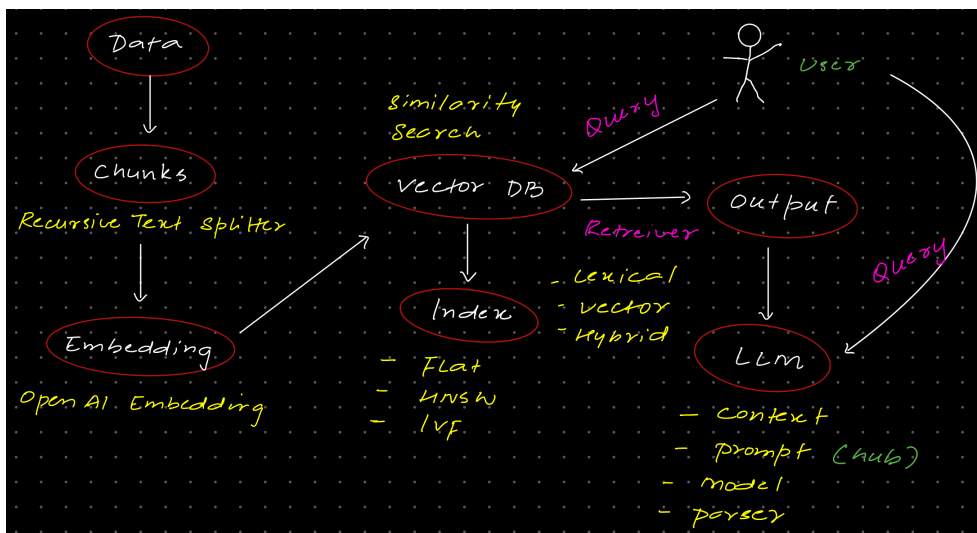
1      2      3      4

# Azure AI Search + Langchain

AI Search makes process really easy.



Cleaning
Standardizing
Chunking

Azure OpenAI/
OpenAI

Azure Blob

Azure Search
Indexer

Embedding
Indexing
Retrieval

LangChain
Langchain

Final
Response

User

# RAG + Langchain

Complex Process



Data

Chunks

Recursive Text Splitter

Embedding

Open AI. Embedding

Similarity
Search

Vector DB

Index

- Flat
- HNSW
- IVF

User

Query

Retriever

Output

- Lexical
- Vector
- Hybrid

LLM

Query

- Context
- Prompt (hub)
- Model
- Parser

Imp. Azure AI Search makes process easy.