

AI Question Answering Agent - Complete Explanation

Let me break down each implementation in simple terms so you understand exactly what's happening!

I Part 1: Simple QA Agent (The Basics)

The Code Structure

```
from groq import Groq

class SimpleQAAgent:
    def __init__(self, api_key):
        self.client = Groq(api_key=api_key)
```

What's happening?

- from groq import Groq - We're importing the Groq library (it's like importing a tool)
- class SimpleQAAgent - We're creating a "blueprint" for our agent (like a recipe)
- def __init__ - This runs when you create the agent (initialization)
- self.client = Groq(api_key=api_key) - We're connecting to Groq's AI service

Why this way?

- Using a class keeps everything organized
- self.client stores our connection so we can reuse it
- We only connect once, then use it many times (efficient!)

The Answer Function

```
def answer(self, question, context=""):
    if context:
        prompt = f"""Based on the context below, answer the question.

Context: {context}
Question: {question}
Answer:"""

    else:
        prompt = question
```

What's happening?

- def answer - Function that takes a question and optional context
- if context: - Checks if you provided extra information (context)
- prompt = f"""\n..."""\n - Creates the instruction for the AI
- If you give context, it tells AI to answer based on that
- If no context, just sends the raw question

Why this way?

- Sometimes you want AI to answer from your own documents (context)
- Sometimes you want general knowledge answers (no context)
- This function handles both cases!

Calling the AI

```
response = self.client.chat.completions.create(  
model="llama-3.3-70b-versatile",  
messages=[{"role": "user", "content": prompt}],  
temperature=0.3,  
max_tokens=500  
)  
return response.choices[0].message.content
```

What's happening?

- `self.client.chat.completions.create()` - Send request to Groq's AI
- `model="llama-3.3-70b-versatile"` - Which AI model to use (free, fast, smart)
- `messages=[...]` - The conversation history (we're sending one message)
- `temperature=0.3` - How creative the AI should be (0.3 = more focused/accurate)
- `max_tokens=500` - Maximum length of answer (500 words roughly)
- `return response.choices[0].message.content` - Extract just the text answer

Why these settings?

- **Temperature 0.3** - Low = consistent, factual answers (good for Q&A)
- **Max tokens 500** - Prevents super long answers (saves API credits)
- **llama-3.3-70b** - Best free model on Groq (fast + accurate)

¶ Part 2: RAG QA Agent (The Smart One)

What is RAG?

RAG = Retrieval Augmented Generation

Think of it like this:

1. **Retrieval** - Find relevant info from your documents (like searching)
2. **Augmented** - Add that info to the question
3. **Generation** - AI generates answer using both question + found info

Why RAG?

- AI doesn't know about YOUR specific documents
- RAG lets AI answer from YOUR data
- More accurate, no hallucinations!

Initialization

```
class RAGQAAgent:  
    def __init__(self, groq_api_key):  
        self.embedder = SentenceTransformer('all-MiniLM-L6-v2')  
        self.llm = Groq(api_key=groq_api_key)  
        self.index = None  
        self.documents = []
```

What's happening?

- SentenceTransformer('all-MiniLM-L6-v2') - A model that converts text to numbers
- self.llm = Groq(...) - Connection to AI for answering
- self.index = None - Will store our searchable database
- self.documents = [] - Will store actual document texts

Why these components?

- **Embedder** - Turns text into vectors (arrays of numbers) so we can compare similarity
- **LLM** - The actual AI that generates answers
- **Index** - Fast search through thousands of documents (FAISS library)
- **Documents list** - Keep originals to show as sources

Adding Documents

```
def add_documents(self, documents):  
    self.documents.extend(documents)  
    embeddings = self.embedder.encode(documents)
```

What's happening?

- self.documents.extend(documents) - Save the original texts
- self.embedder.encode(documents) - Convert each document to a vector

Example:

"Python is a language" → [0.2, 0.8, 0.1, 0.5, ...] (384 numbers)
"Machine learning is AI" → [0.1, 0.3, 0.9, 0.2, ...] (384 numbers)

Why vectors?

- Computers can't directly compare "similarity" of text
- But they CAN compare similarity of number arrays!
- Similar meanings = similar vectors

Creating the Search Index

```
if self.index is None:  
    dimension = embeddings.shape[1] # 384 for all-MiniLM-L6-v2  
    self.index = faiss.IndexFlatL2(dimension)  
  
    self.index.add(embeddings.astype('float32'))
```

What's happening?

- dimension = embeddings.shape[1] - Get vector size (384 numbers)
- faiss.IndexFlatL2(dimension) - Create FAISS search index
- .add(embeddings) - Add all document vectors to index

Why FAISS?

- Super fast search (even with millions of documents)
- Used by Facebook, Google in production
- Free and works locally (no API needed)
- L2 means "Euclidean distance" (measures vector similarity)

Retrieving Relevant Documents

```
def retrieve(self, query, k=3):
    query_embedding = self.embedder.encode([query])
    distances, indices = self.index.search(
        query_embedding.astype('float32'), k
    )
    return [self.documents[i] for i in indices[0]]
```

What's happening step by step:

1. Convert question to vector

```
query_embedding = self.embedder.encode([query])
```

"Who created Python?" → [0.3, 0.7, 0.2, ...]

2. Search for similar vectors

```
distances, indices = self.index.search(query_embedding, k=3)
```

Returns: indices = [5, 2, 8] (positions of closest documents)

distances = [0.2, 0.4, 0.6] (how close they are)

3. Get actual documents

```
return [self.documents[i] for i in indices]
```

Returns the 3 most relevant document texts

Why k=3?

- 3 documents usually give enough context
- More = slower and might confuse AI with irrelevant info
- Less = might miss important information

Answering with RAG

```
def answer(self, question, return_sources=True):
    # Step 1: Find relevant documents
    relevant_docs = self.retrieve(question, k=3)

    # Step 2: Combine them into context
    context = "\n\n".join(relevant_docs)

    # Step 3: Create prompt
    prompt = f"""Based on the following context, answer the question.

Context:
{context}

Question: {question}

Answer:"""

    # Step 4: Get AI answer
    response = self.llm.chat.completions.create(...)
```

Context:
{context}

Question: {question}

Answer:"""

```
# Step 4: Get AI answer
response = self.llm.chat.completions.create(...)
```

The Complete Flow:

User asks: "When was Python created?"

↓

1. Convert question to vector: [0.3, 0.7, 0.2, ...]
↓
2. Search FAISS index for similar document vectors
↓
3. Find top 3 matches:
 - "Python was created by Guido van Rossum in 1991..."
 - "Python emphasizes code readability..."
 - "Python is high-level language..."
↓

4. Combine into context string
- ↓
5. Send to AI: "Based on context [above docs], answer: When was Python created?"
- ↓
6. AI reads context and generates: "Python was created in 1991 by Guido van Rossum."
- ↓
7. Return answer + sources

Why this works:

- AI only sees relevant documents (not all 1000s)
 - AI answers from YOUR data, not general knowledge
 - Fast because FAISS pre-computed embeddings
 - Accurate because context is provided
-

Part 3: Loading Documents from Files

Loading PDFs

```
def load_pdf(pdf_path):
    with open(pdf_path, 'rb') as file:
        reader = PyPDF2.PdfReader(file)
        text = ""
        for page in reader.pages:
            text += page.extract_text() + "\n"
    return text
```

What's happening?

- `open(pdf_path, 'rb')` - Open PDF file in binary read mode
- `PyPDF2.PdfReader(file)` - Create PDF reader object
- `for page in reader.pages` - Loop through each page
- `page.extract_text()` - Pull out text from that page
- Combine all pages into one big string

Why this way?

- PDFs store text in complex format
- PyPDF2 handles all the complexity
- We get clean text output

Chunking Text

```
def chunk_text(text, chunk_size=500, overlap=50):
    words = text.split()
    chunks = []
```

```
for i in range(0, len(words), chunk_size - overlap):
    chunk = ' '.join(words[i:i + chunk_size])
    chunks.append(chunk)
```

```
return chunks
```

What's happening?

- `text.split()` - Split text into individual words
- `range(0, len(words), chunk_size - overlap)` - Create windows
- `words[i:i + chunk_size]` - Take 500 words at a time
- `chunk_size - overlap` - Step by 450 (500-50) to create overlap

Visual Example:

Document: "word1 word2 word3 word4 word5 word6 word7 word8..."

Chunk 1: words 1-500 (overlap region)

Chunk 2: words 451-950 (overlap region)

Chunk 3: words 901-1400

Why overlap?

- Prevents cutting sentences/ideas in half
- If answer spans chunk boundary, overlap captures it
- 50-word overlap = ~2-3 sentences

Why 500 words per chunk?

- Small enough for fast search
- Large enough for complete context
- Fits well in AI's context window
- Good balance for embedding model

Part 4: LangChain Version (Framework Approach)

What is LangChain?

LangChain = Pre-built components for AI apps

Instead of writing everything yourself:

- ✓ Use pre-built chains
- ✓ Use pre-built retrievers
- ✓ Use pre-built memory
- ✓ Add tools easily

Trade-off:

- **Pros:** Faster development, less code, many features
- **Cons:** More dependencies, slightly slower, less control

The Setup

```
self.llm = ChatGroq(  
    groq_api_key=groq_api_key,  
    model_name="llama-3.3-70b-versatile",  
    temperature=0.3  
)  
  
self.embeddings = HuggingFaceEmbeddings(  
    model_name="all-MiniLM-L6-v2"  
)
```

What's different from manual RAG?

- LangChain wraps Groq in ChatGroq class
- Wraps embeddings in HuggingFaceEmbeddings class
- Same models, just LangChain's interface

Creating Knowledge Base

```
text_splitter = RecursiveCharacterTextSplitter(  
    chunk_size=500,  
    chunk_overlap=50  
)  
texts = text_splitter.create_documents(documents)
```

What's happening?

- RecursiveCharacterTextSplitter - Smart chunker from LangChain
- Splits on paragraphs first, then sentences, then words
- More intelligent than our simple word-based chunker

Why "Recursive"?

Try to split by: \n\n (paragraphs)
↓ If chunk still too big
Try to split by: . (sentences)
↓ If chunk still too big
Try to split by: spaces (words)

Creating QA Chain

```
self.qa_chain = RetrievalQA.from_chain_type(  
    llm=self.llm,  
    chain_type="stuff",  
    retriever=self.vectorstore.as_retriever(search_kwargs={"k": 3}),  
    return_source_documents=True  
)
```

What's happening?

- RetrievalQA - Pre-built QA chain from LangChain
- chain_type="stuff" - Put all retrieved docs into one prompt
- retriever=... - How to search (same FAISS under the hood)

- `search_kwargs={"k": 3}` - Return 3 documents
- `return_source_documents=True` - Give us sources too

Chain Types Explained:

Type	How it works	Best for
stuff	Puts all docs in one prompt	Small docs (what we use)
map_reduce	Summarize each doc, then combine	Large docs
refine	Answer from doc 1, refine with doc 2, etc.	Very detailed answers
map_rerank	Score each doc's answer, pick best	Multiple possible answers

Why "stuff"?

- Simplest and fastest
- Works when 3 docs fit in context window
- No extra API calls needed

Using the Chain

```
result = self.qa_chain.invoke({"query": question})

return {
    'answer': result['result'],
    'sources': [doc.page_content for doc in result['source_documents']]
}
```

What happens internally:

1. qa_chain receives question
↓
2. Retriever searches FAISS index
↓
3. Gets top 3 documents
↓
4. Creates prompt automatically
↓
5. Calls LLM with prompt
↓
6. Returns formatted result

You write: 1 line

LangChain does: All the RAG logic we wrote manually!

II Comparison: Manual vs LangChain

Manual RAG Agent

Pros:

- ✓ Full control over every step
- ✓ Less dependencies (lighter)
- ✓ Easier to debug (you wrote it)
- ✓ Learn how RAG actually works
- ✓ Faster execution (no framework overhead)

Cons:

- ✗ More code to write
- ✗ Need to handle edge cases yourself
- ✗ No built-in memory/tools

Best for:

- Learning AI/ML concepts
- Production apps needing speed
- Custom RAG logic
- Portfolio projects (shows deeper understanding)

LangChain Version

Pros:

- ✓ Much less code
- ✓ Many built-in features (memory, tools, agents)
- ✓ Easy to add new capabilities
- ✓ Large community support

Cons:

- ✗ More dependencies
- ✗ Harder to debug (framework abstractions)
- ✗ Slightly slower
- ✗ Less control

Best for:

- Rapid prototyping
 - Complex multi-agent systems
 - When you need memory/tools/chains
 - Enterprise projects with many features
-

¶ Key Concepts Summary

1. Embeddings (Vector Representations)

Simple explanation: Converting text to numbers so computers can understand similarity

"I love Python" → [0.2, 0.8, 0.1, 0.5, ...]

"Python is great" → [0.3, 0.7, 0.2, 0.4, ...] ← Similar vectors!

"I love cooking" → [0.9, 0.1, 0.8, 0.2, ...] ← Different vector

Why 384 dimensions?

- More dimensions = more nuanced understanding
- all-MiniLM-L6-v2 uses 384
- Larger models use 768 or 1024
- Trade-off: accuracy vs speed

2. Vector Search (FAISS)

Simple explanation: Finding similar vectors quickly

How it works:

Add documents

```
index.add([
    [0.2, 0.8, ...], # Doc 1
    [0.3, 0.7, ...], # Doc 2
    [0.9, 0.1, ...] # Doc 3
])
```

Search

```
query = [0.25, 0.75, ...]
results = index.search(query, k=2)
```

Returns: Doc 1 and Doc 2 (most similar)

Why L2 distance?

- Measures straight-line distance between vectors
- Close vectors = similar meaning
- Fast to compute

3. Temperature

Simple explanation: How "creative" or "random" the AI is

Temperature 0.0: "Python was created in 1991."

(Same answer every time, factual)

Temperature 0.5: "Python was created in 1991 by Guido van Rossum."

(Slightly varied but still accurate)

Temperature 1.0: "Python, the elegant language, emerged in 1991..."

(More creative, sometimes less accurate)

For Q&A, use 0.2-0.4:

- Want consistent, factual answers
- Don't want creative embellishments
- Minimize hallucinations

4. Context Window

Simple explanation: How much text AI can "see" at once

Llama-3.3-70b: ~8000 tokens (~6000 words)

If your prompt is:

- System message: 100 tokens
- 3 retrieved docs: 1500 tokens
- Question: 20 tokens
- Total: 1620 tokens used

Remaining for answer: 6380 tokens

Why this matters:

- If docs too long, they get cut off
- That's why we chunk to 500 words
- Always leave room for answer

¶ Practical Tips

For Your Projects

1. **Start with Simple QA** - Understand the basics
2. **Move to RAG** - Most real-world applications
3. **Add LangChain** - When you need advanced features

For Colab/Kaggle

These work on free tier:

- ✓ sentence-transformers (runs on free GPU)
- ✓ FAISS (CPU is fine, super fast)
- ✓ Groq API (14,400 free requests/day)

Download embedder once, reuse:

```
embedder = SentenceTransformer('all-MiniLM-L6-v2')
embedder.save('./saved_model')
```

**Later: embedder =
SentenceTransformer('./saved_model')**

For Freelancing

RAG QA is perfect for:

- Customer support chatbots (answer from docs)
- Internal company knowledge base
- Legal/medical document Q&A
- Educational course assistants
- Research paper Q&A systems

Clients love it because:

- Answers from THEIR data (not generic AI)
- Can verify sources
- No hallucinations
- Easy to update (add new docs)

For Interviews

Know these concepts:

- "Explain RAG" → Retrieval + Generation
 - "Why embeddings?" → Semantic similarity search
 - "Why chunking?" → Fit in context window, better retrieval
 - "FAISS vs alternatives?" → Speed, production-ready
 - "Temperature setting?" → 0.3 for factual Q&A
-

¶ Next Steps for You

Based on your background (AI engineering student, building projects):

Week 1: Master the Basics

- [] Run Simple QA Agent
- [] Understand embeddings
- [] Try different temperatures
- [] Test with various questions

Week 2: Build RAG System

- [] Implement RAG QA Agent
- [] Load your own PDFs (study materials!)
- [] Experiment with chunk sizes
- [] Measure retrieval accuracy

Week 3: Production Ready

- [] Add error handling
- [] Deploy with FastAPI
- [] Add Streamlit UI
- [] Push to GitHub

Week 4: Advanced Features

- [] Try LangChain version
- [] Add conversation memory
- [] Add multiple document sources
- [] Deploy on free hosting (Render, Railway)

Portfolio Project Ideas

- 1. Personal Study Assistant**
 - Upload your course PDFs
 - Ask questions about material
 - Get instant answers with sources
- 2. Research Paper Q&A**
 - Load ML/AI papers
 - Query about methods, results
 - Great for literature review
- 3. Code Documentation Helper**
 - Load library docs
 - Ask "how to" questions
 - Instant code examples

All of these are freelance-ready! ¶

🔗 Resources to Learn More

Embeddings:

- Sentence Transformers documentation
- "How embeddings work" tutorials
- Try different models (compare accuracy vs speed)

FAISS:

- FAISS GitHub wiki
- "Vector databases explained"
- When to use FAISS vs Pinecone vs Weaviate

LLMs:

- Groq documentation
- Llama model papers
- Prompt engineering guides

RAG:

- LangChain RAG tutorial
- "Building production RAG systems"
- RAG evaluation metrics

Keep building and learning! 🔗