# YOUTUBE Q&A CHATBOT

AI Technology Final Project Report
Developed by Group 10
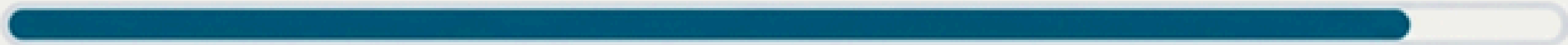
# The search for a single answer can feel like an entire lecture.

**?** **? ? ?** **?** **?**

2:34:51

People need key information from long educational videos without watching for hours. The goal was to build an intelligent agent to solve this.

# 3BROWN1BLUE

Problem Solving

3b1b

•live

# SHELDON COOPER
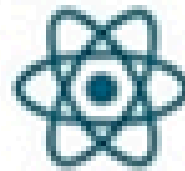
# Meet the Agent: Instant, Accurate, and... Peculiar.

Whisper

Pinecone

LangChain

GPT-4o-mini

Lovable

## Dr. Cooper

What is the core principle of RAG?

*Bazinga!* An intriguing question. The foundational principle of Retrieval-Augmented Generation is elementary: one fortifies a Large Language Model with external, verified knowledge *before* it presumes to formulate an answer. This prevents it from succumbing to the... let's call it 'creative confabulation' inherent to its nature. It is, in essence, intellectual scaffolding. Simple, yet effective.

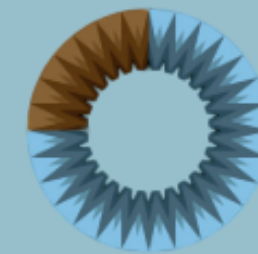# TECHNICAL SUMMARY

Architecture based on a **RAG pipeline** :

1. Video transcription
2. Chunking and embedding
3. Storing embeddings
4. Retrieving top-3 relevant chunks
5. Feeding context + user question
6. Maintaining conversation memory

⇨ Allows the agent to give answers grounded in the videos, avoiding hallucinations

# DATASET

**How did we do it ?**

- full whisper transcription
- Cleaned the text
- Split transcripts into chunk size of 1000 tokens and chuck overlaps of 200 tokens
- Added metadeta :
  - Video title
  - URL
  - Timestamp
  - chunk index
- split our entire set of video transcriptions into 207 chunks

Dataset is scalable : Adding more videos → re-run embedding script → Pinecone updates the index automatically

3Blue1Brown

# Q&A AGENT
## RAG + Persona + Memory

Chatbot powered by LangChain's ConversationalRetrievalChain

**How it works:**
1. The user asks a question.
2. Retrieval of the most relevant transcript chunks from Pinecone.
3. GPT model generates an answer strictly based on the chunks.
   a. Answer is from the persona (aka has a personality) → temperature of 0.7
4. The agent has buffer memory so it can understand follow-up questions

# CHALLENGES FACED & SOLUTIONS

**Challenge 1 — Tool/Environment Issues**
- Challenge: Lovable and code behaved differently on each computer.
- Solution: So we had to change the code for it to work on each computer separately

**Challenge 2 — Whisper Transcription Errors**
- Challenge: Inconsistent, messy transcripts.
- Solution: Added preprocessing to clean and normalize text.

**Challenge 3 — LLM Hallucinations**
- Challenge: Model answered with general knowledge instead of context.
- Solution: Stricter prompts, minimum 3 retrieved chunks, and "say you don't know" rule.
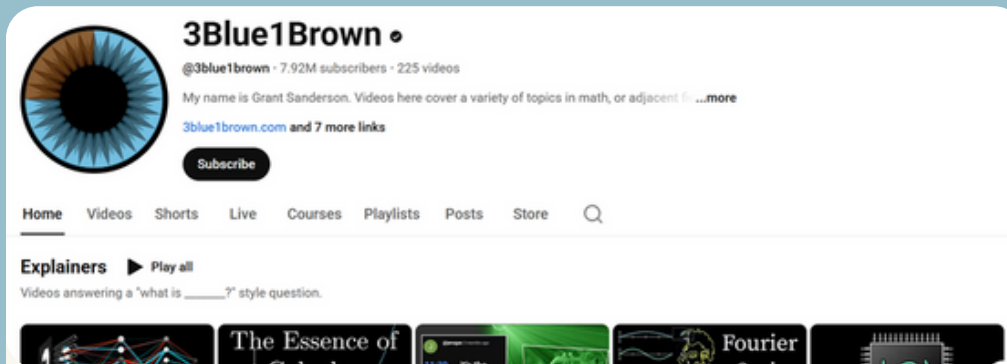
**Challenge 4 — Persona Drift**
- Challenge: Model drifted from Sheldon's tone.
- Solution: Lowered temperature to 0.7 and reinforced persona rules.

# FUTURE IMPROVEMENTS

- Better memory system
- Improved embedding models
- UI enhancements (avatars, layouts, personas)
- Multi-modal features (frames, diagrams, audio)
- Different languages

# CONCLUSION

## We built a YouTube Q&A chatbot that:



Transcribes and processes YouTube videos.
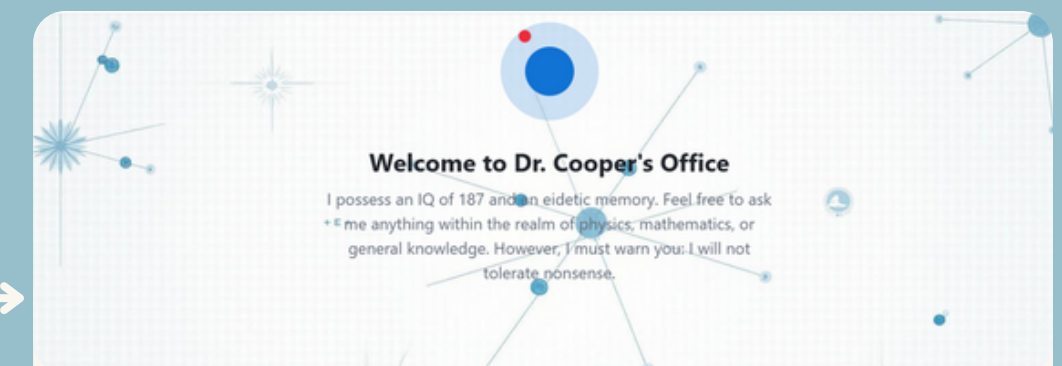


Stores and retrieves vectorized knowledge.



Uses RAG to answer user questions with contextual accuracy.



Maintains conversation history.
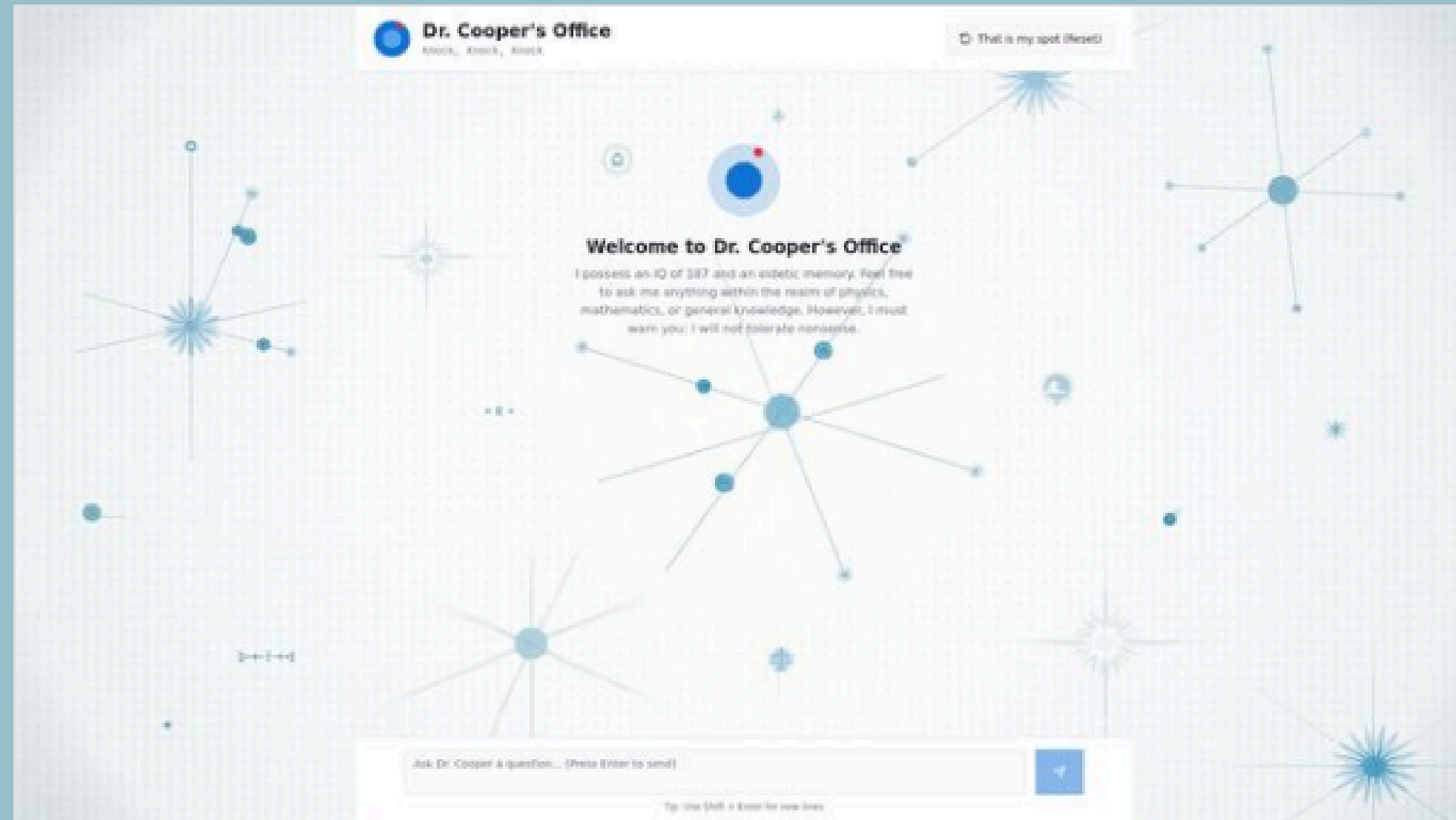


Features a fun and consistent persona.



Runs within an accessible web interface.

# LET'S TRY IT!



## Dr. Sheldon Cooper AI Chatbot

Have a conversation with Dr. Sheldon Cooper, theoretical physicist with an IQ of 187

Lovable