Here's **LinkedIn post** you can use to showcase your achievement of building a **RAG-based chatbot with feedback learning** from a book (like Machine Learning textbook):

# **✓** LinkedIn Post: Showcase of RAG Chatbot Project

💋 Just Built My Own AI Chatbot Using RAG Architecture! 🔄 📘

Hey LinkedIn fam!

I'm super excited to share that I've recently completed a hands-on **AI project** where I built a **Retrieval-Augmented Generation (RAG) Chatbot** from scratch — using a full book on **Machine Learning** as the knowledge base!

#### **What It Does:**

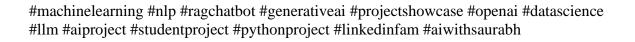
- ✓ Takes your question as input
- Retrieves relevant content chunks from the uploaded book
- ✓ Uses **Phi-2** (**LLM**) to generate human-like, context-aware answers
- ✓ Collects user feedback (Yes/No) to learn from responses paving the way for **Reinforcement Learning!**

## **Tools & Technologies Used:**

- Phi-2 LLM
- ChromaDB + SentenceTransformer for vector-based retrieval
- X Python, PyPDF2, Transformers, NLP
- Feedback collection loop for future RLHF tuning

# **Skills Applied:**

- Natural Language Processing (NLP)
- Vector Search & Embeddings
- Retrieval-Augmented Generation (RAG)
- LLM Prompt Engineering
- Feedback-Driven Learning
- End-to-End Project Deployment
- **○** Want to try something similar or collaborate on AI/ML projects? Let's connect!
- Feel free to drop me a message or leave your thoughts in the comments!



Here's a complete explanation of all the steps used to build a RAG-based (Retrieval Augmented Generation) chatbot, This will help you confidently answer when you're asked how you built or worked on a chatbot using LLMs, embeddings, or feedback loops.



# RAG-Based Chatbot: Step-by-Step with Explanation

### Step 1: Install Required Libraries

## 📌 Purpose:

Install Python libraries for PDF extraction, text embedding, vector search, and LLM usage.

## **Explanation:**

- PyPDF2: To extract text from PDF files.
- sentence-transformers: To generate vector embeddings of text.
- chromadb: A vector database for storing and querying embeddings.
- transformers: Load large language models like Phi-2.

## Step 2: Upload and Extract Text from PDF

## 📌 Purpose:

Get all the raw content from the book (PDF) to prepare for question answering.

## **Explanation:**

- PDFs are read page by page using PyPDF2.
- The extracted text will be later divided into chunks for retrieval.

#### • Step 3: Clean the Text (NLP Preprocessing)

## **≯** Purpose:

Make the text clean, standardized, and easier to embed and analyze.

#### **Explanation:**

- Converts to lowercase.
- Removes special characters.
- Removes stopwords (common words like *the*, *is*, *and*).
- Lemmatizes words (e.g.,  $running \rightarrow run$ ). This step ensures meaningful vector representation and reduces noise.

#### • Step 4: Generate Embeddings & Store in Vector DB (ChromaDB)

## 📌 Purpose:

Convert text chunks into numerical vectors and store them for retrieval.

#### **Explanation:**

- Text is split into chunks (e.g., 500 characters).
- Each chunk is embedded using SentenceTransformer.
- The embeddings are stored in ChromaDB with metadata (text content). Later, during a query, we'll use this database to find the most relevant context chunks.

#### Step 5: Load the Large Language Model (Phi-2)

## **★** Purpose:

Use a pre-trained LLM to generate natural language answers.

#### **Explanation:**

- Phi-2 is a small, powerful open-source model trained for reasoning and answering.
- transformers library loads both the model and its tokenizer. This model will take the **retrieved context** and **user query** and generate an answer.

#### Step 6: Define a Retrieval Function

## **≯** Purpose:

Search the vector database for relevant information based on the user's question.

## **Explanation:**

- The question is embedded using the same embedding model.
- ChromaDB is queried using this embedding to retrieve top-k similar chunks. This provides **grounded context** for the LLM to use (important in RAG!).

#### Step 7: RAG Chatbot Loop + Feedback System

## **≯** Purpose:

Interact with the user, answer their question, and get feedback on the answer.

### **Explanation:**

- A loop asks the user for a question.
- If the user types byy or stop, the loop ends.
- Otherwise:
  - o The question is embedded.
  - o Top relevant text chunks are retrieved.
  - o These are combined with the user query in a prompt.
  - o The LLM generates a response using the context.
  - o The user can provide feedback (yes/no), which is saved for improvement.

#### Step 8: Save Feedback for Future Learning

## **≯** Purpose:

Store user feedback for analysis or model fine-tuning.

#### **Explanation:**

- Each interaction (question, context, answer, feedback) is saved in a JSON file.
- This can help:
  - o Analyze performance.
  - o Fine-tune future models.
  - o Create supervised datasets for reinforcement learning or RAG optimization.

# ✓ Interview Points You Can Emphasize

Topic	🦞 What You Can Say	
RAG Concept	Combines retrieval of knowledge (using vector DB) and generation of answers (using LLMs).	
Embedding	Used SentenceTransformer to convert text and questions into numerical vectors.	
Vector Store	Used ChromaDB to store and search relevant chunks of the book.	
Prompt Engineering	Carefully structured prompt to provide retrieved context to the model.	
Feedback Loop	User feedback collected to improve accuracy and understand model performance.	
LLM Used	Used Phi-2 model (lightweight but effective) to generate natural language responses.	
Use Case	Built a book-based chatbot that can answer questions on topics like machine learning, AI, etc.	
Improvement	Could integrate reinforcement learning from user feedback to adapt responses over time.	

Q = why use this vector embedding model why not other model

Great question! • Let's break it down:

# Why use a vector embedding model (like SentenceTransformer) in a RAG-based chatbot?

# **★** What are vector embeddings?

Vector embeddings convert text (words/sentences) into **numerical vectors**—think of it like putting meanings into numbers. These vectors capture **semantic meaning**, so similar meanings are **closer together** in the vector space.

Reason	🖳 Explanation	
Semantic Search	Embeddings help find relevant chunks that are <i>semantically</i> (not just word-for-word) related to the user's question.	
<b>♦</b> Efficient & Fast	SentenceTransformer models like all-MinilM-L6-v2 are fast and light—perfect for real-time retrieval.	
► Better than Traditional TF-IDF	Unlike TF-IDF or Bag of Words, embeddings capture <b>context and word meaning</b> , which is critical for understanding questions.	
Compatible with ChromaDB	These embeddings integrate well with vector databases like ChromaDB for quick similarity search.	
<b>%</b> Pre-trained and Fine-tuned	SentenceTransformers are already trained on tasks like Semantic Textual Similarity, making them ready for use.	
Plug-and-play	Easy to use with a single .encode() line. No need to build or fine-tune deep models from scratch.	

# **◯** Why not use a different model?

Here's a comparison:

<b>Model Type</b>	Why Not Ideal		
★ TF-IDF / CountVectorizer	Only uses word frequency; no understanding of <b>meaning</b> or <b>context</b> .		
➤ Word2Vec / GloVe	Only word-level embeddings, not sentence-level. Hard to use for full queries.		
X BERT-base directly	BERT isn't optimized for sentence similarity unless fine-tuned; SentenceTransformer wraps and fine-tunes it.		
X Custom-trained embeddings	Training from scratch requires data, time, compute, and expertise—overkill for most RAG applications.		

# **✓** Summary:

We use vector embedding models like SentenceTransformer because they provide fast, meaningful, and context-aware numerical representations of text—perfect for finding relevant content in RAG-based systems.

### Q = why use chroma DB database why not other

Awesome follow-up question! Let's break down **why we use ChromaDB** in a **RAG-based chatbot**, especially in interviews where you need crisp and clear reasoning.



# Why use ChromaDB in a RAG-based chatbot?

ChromaDB is a vector database—it stores and retrieves high-dimensional embeddings (vector representations of text). In RAG (Retrieval-Augmented Generation), once we convert documents into vectors using models like SentenceTransformer, we need to store them in a way that allows fast similarity search. That's exactly what ChromaDB does!



# **Key Reasons to Use ChromaDB:**

Feature	Explanation	
Blazing fast retrieval	ChromaDB is optimized for <b>real-time nearest neighbor search</b> using Approximate Nearest Neighbor (ANN) techniques.	
In-memory & persistent modes	You can run it in RAM for fast prototyping or store data permanently on disk.	
<b>*</b> Easy to integrate	Seamless integration with Python, LangChain, and embedding models like SentenceTransformer.	
Schema-free and document-rich	Allows storing <b>metadata</b> along with embeddings (like page number, source, tags, etc.)	
No server setup required	Can run <b>locally</b> , no need for external hosting or Docker unless scaling.	
<b>X</b> Open-source and lightweight	Great for small to medium projects or solo devs (like your RAG chatbot)!	

# Why not use other databases?

<b>X</b> Alternative	<b>Why not ideal in this use case</b>	
FAISS (by	Super fast but more low-level; harder to use, no persistent storage out of	
Facebook)	the box.	
Pinecone	Great managed cloud solution, but <b>requires API key and internet</b> . Not ideal for offline/local projects.	
Weaviate /	Heavy, need setup with Docker/Kubernetes. Overkill for simple or loca	
Milvus	RAG bots.	
Elasticsearch	Powerful but made for keyword search, not optimized for deep vector similarity unless extra plugins are added.	

# **©** Summary for Interview:

"ChromaDB is an open-source vector database that's easy to use, lightweight, and optimized for semantic search. It integrates well with Python-based ML pipelines and is perfect for

storing and retrieving text embeddings in real time. Compared to heavier tools like FAISS or Pinecone, ChromaDB provides a great balance between speed, flexibility, and ease of use for local or medium-scale RAG applications."

### Q = how to search top K similar chunks which technique used,

Great interview-style question! Let's break down how we **search top K similar chunks** in a RAG (Retrieval-Augmented Generation) pipeline and what **techniques** are used behind the scenes.

# How to Search Top K Similar Chunks?

We search the **top K similar chunks** using **vector similarity search**. Here's how it works:

# **Step-by-Step:**

- 1. Text  $\rightarrow$  Embedding (vector)
  - Convert user query and document chunks into vector form using an embedding model (e.g., SentenceTransformer).
- 2. Query Vector vs. Chunk Vectors
  - Compare the query vector to the stored vectors in the vector database (like ChromaDB).
- 3. Compute Similarity Score
  - o Use a similarity/distance metric (e.g., cosine similarity).
- 4. Return Top K Results
  - o Fetch the **K** most similar chunks (highest similarity or lowest distance).

# **Mhich Technique is Used?**

- **✓** 1. Similarity Metric (Core of Top-K Search)
  - The most common metrics:
    - Cosine Similarity ( Most commonly used in NLP)
    - o Euclidean Distance
    - o Dot Product

ChromaDB by default uses Cosine Similarity.

## **2.** ANN (Approximate Nearest Neighbors)

- To speed up the search when the dataset is large.
- Instead of scanning all vectors (which is slow), it uses efficient indexing structures like:
  - **HNSW (Hierarchical Navigable Small World)**
  - IVF (Inverted File Index)
  - LSH (Locality-Sensitive Hashing)

These are **indexing techniques** that let you search billions of vectors **fast**.



# In Code (ChromaDB)

```
python
CopyEdit
results = collection.query(
     query_embeddings=[query_vector],
n_results=3 # Top K = 3
```

This line asks ChromaDB to return the top 3 most similar chunks to your query vector using cosine similarity.



## **Short Interview Answer:**

"To retrieve the top K similar chunks in a RAG pipeline, we use vector similarity search. The user query is embedded into a vector, and we compare it to pre-stored document vectors using similarity metrics like cosine similarity. Libraries like ChromaDB perform this search efficiently using Approximate Nearest Neighbor (ANN) techniques such as HNSW, ensuring fast and accurate retrieval."

Q = which is the LLM model u used and why not use other LLM model

Awesome follow-up! Let's break this down clearly to help you answer interview questions confidently -



## Which LLM Model Did We Use?

We used:



✓ Phi-2 by Microsoft

It's a small, fast, and efficient LLM designed for reasoning and instruction-following tasks.



# 🥰 Why Did We Use Phi-2?

Here are the main reasons:

#### 1. Lightweight & Efficient

- It's smaller than GPT-style models (like LLaMA or GPT-3), making it:
  - o Faster to load
  - o Less memory-intensive
  - Works on a normal laptop or Colab (good for prototyping)

## 2. Good for Question Answering

Pretrained for **instruction-following tasks**, so it gives better answers when you pass in a prompt like:

```
makefile
CopyEdit
Context: ...
Question: ...
Answer:
```

## 3. Open Source and Easy to Use

- Available via Hugging Face
- Supports PyTorch, Transformers, Accelerate (all tools you already use in Colab)

# Why Not Use Other LLM Models?

Let's compare -

Model	Why Not Used Here?	Notes
<b>GPT-3.5 / GPT-4</b>	Requires <b>API key</b> + \$\$ <b>6</b>	Not local, costs money
LLaMA 2	Too large for 8GB RAM or Colab, needs GPU	Powerful, but heavy
Mistral / Mixtral	Large model, often needs <b>quantization</b> to run locally	Also great, but more setup
OpenChat / Vicuna	Finetuned models but larger in size	Needs >16GB RAM

# Short Interview Answer:

"I used the Phi-2 model from Microsoft because it's lightweight, optimized for instruction tasks, and easy to run on systems with lower hardware specs like Colab or an 8GB RAM laptop. It balances performance and accuracy well. Heavier models like GPT or LLaMA offer more capabilities but require more computational power or API access, which wasn't ideal for my setup."

Let me know if you want help switching to a bigger model like LLaMA 2 using quantization (like GGUF or GPTQ) on your hardware

Q = how to work the feedback loop and how to improve the answer quality/ accuracy

Great question! Feedback loops are **crucial** in Retrieval-Augmented Generation (RAG) systems to **improve answer quality over time**—especially if you're aiming for a more *adaptive*, *intelligent chatbot*. Here's a clear breakdown

# What is a Feedback Loop in RAG?

A feedback loop is when user feedback (positive or negative) is used to:

- Improve retrieval quality
- Enhance **LLM responses**
- Adapt to user preferences over time

# Types of Feedback You Can Use

#### 

# How Feedback Loop Works in RAG

## **Step-by-step Workflow:**

- 1. User asks a question
- 2. RAG retrieves chunks + LLM gives an answer ✓
- 3. User provides feedback
  - o "Good answer" 🗸
  - o "Incorrect answer" X
- 4. System logs query, context, feedback
- 5. That feedback is used to:
  - Update retrieval ranking
  - o Improve embedding similarity
  - o Or even fine-tune the LLM or reranker

# Techniques to Improve Accuracy Using Feedback

# **✓** 1. Reranking with Feedback

Use models like **BGE-Reranker**, **Cohere Reranker**, or train your own:

- User gives feedback  $\rightarrow$  label the chunk relevance  $\rightarrow$  fine-tune a reranker
- Helps rank better chunks in future retrievals

# 2. Query Expansion or Reformulation

Use feedback to **reformulate queries** (e.g., paraphrase or add keywords) to fetch better context next time.

## **✓** 3. LLM Fine-Tuning or Prompt Tuning

If you collect enough Q&A + feedback data:

• Fine-tune Phi-2 or another model to align better with your domain or user style.

## 4. Personalized Context Weighting

Adjust how you select or score documents:

- Give more weight to docs previously rated "helpful"
- Learn what kinds of chunks produce better answers

# Notional: Reinforcement Learning (Advanced)

Use Reinforcement Learning from Human Feedback (RLHF) or Reinforcement Learning with AI Feedback (RLAIF):

- Rewards model when it gets good feedback
- Penalizes when it performs poorly

Think: GPT models fine-tuned with RLHF use a similar loop.



# Interview Answer Template:

"In a RAG system, a feedback loop lets the model learn from user input. For instance, if users rate answers as helpful or not, we can use that to improve document retrieval, re-rank the chunks more accurately, or fine-tune the LLM. Over time, this loop helps the model better understand user intent and generate more accurate responses."