# a QA System for UCLA Statistics Seminars

Fanchao Lin
STATS 418 Final Project

# Data Scraping and Processing



**Seminar formats vary across years, making pattern matching difficult.**

**As a workaround, I extracted the full titles and content blocks directly.**

**Then, I used the OpenAI API to extract seminar topics, speaker info, and summaries line by line.**

June 2, 2025

# The Model

**VECTOR DATABASE:**

- The processed **title, abstract, and speaker bio** of each seminar are combined into a single chunk (450 words on average, around 600 tokens).

- Each chunk is embedded using the **text-embedding-3-small** model.

- The resulting embeddings are stored in a vector database using **ChromaDB**.

**RAG RETRIEVAL + Q&A:**

- The user's question is embedded using the same model.
- An approximate nearest neighbor (ANN) search is performed in ChromaDB based on cosine similarity, retrieving the **top 3** most relevant chunks.
- The retrieved chunks are combined with the user question and passed to **GPT-4 Turbo** to generate a final response.

June 2, 2025

# Deployment

- **Frontend**: Deployed R Shiny application on shinyapps.io

- **Backend API**: Flask-based REST API containerized with Docker

- **Cloud Hosting**: Flask API deployed on Google Cloud Run

- **Container Registry**: Docker image pushed to Docker Hub for reuse

- **Integration**: Shiny frontend communicates with Flask API via HTTPS (JSON POST requests)

June 2, 2025

# Demo

https://lfc787878.shinyapps.io/UCLA_seminar_QA/

## UCLA Statistics Seminar Q&A

**Enter your question:**

What seminars are about Bayesian statistics? summa

Ask

Example questions:

- What seminars are about Bayesian statistics? summarize their content

The seminars about Bayesian statistics

In the first seminar, the speaker disc

In the second seminar, the speaker tal

June 2, 2025

# Demo

**Q:** What seminars are about Bayesian statistics? summarize their content

**A:** The seminars about Bayesian statistics are "Bayesian Inference in Nonparanormal Graphical Models" and "Scaling and Generalizing Bayesian Inference."

In the first seminar, the speaker discusses Bayesian solutions for nonparanormal graphical models, which generalize Gaussian models for continuous variables. The talk focuses on two Bayesian solutions, one with priors on transformation functions and the other with a rank-based likelihood function, and studies their numerical performance through simulation and real data application.

In the second seminar, the speaker talks about recent research on addressing limitations in Bayesian statistics, including handling massive data sets and developing generic algorithms for approximating the posterior. The speaker also discusses how approximate posterior inference algorithms have revolutionized the field, making Bayesian statistics a usable and general-purpose language for data analysis.