

CHATDB FOR LARGE-SCALE ENTERTAINMENT DATASETS

Team Details

Team members background and skills

- Quynh Tran: Background in Business,
 - My name is Quynh Tran and my undergraduate is Business Administration. I have experience in managing and making business decisions based on statistics. I just transferred to Data Science and am doing my master's in person at USC.
- Kevin Bui:
 - My name is Kevin Bui, and I am currently a DEN Student majoring in Computer Science Data Science for my Master of Science. I work full-time as an Application Support Specialist at Center for Creative Leadership. This is a fully remote position, and I work heavily within SQL Server to maintain live data in many different applications and API's. This involves performing data fixes, creating/executing stored procedures, and creating reports for data analysis. I studied at Azusa Pacific University with a Bachelor of Science in Computer Science. I am comfortable with Python, Java, SQL, CSS, PHP, React, and Javascript.
- Riya Berry:
 - My name is Riya Berry, and I am currently an in-person Masters in Applied Data Science student at USC. I completed my BA in Data Science at University of California, Berkeley, and then worked for 2 years as a Data Analyst at a market research consulting firm called Protagonist. I have significant experience in building and developing ML systems, NLP, statistical modeling, and data mining in Python/R/SQL. I am also currently conducting Computer Science Research for the Norman Lear Center at USC Annenberg. As a researcher, I am developing a web-interface platform that allows users to search for key terms within a large-scale database of film and television scripts. Through this role, I have gained increased exposure to front and back-end development, AWS EC2 and ElasticBeanstalk, and distributed database management.

Project requirements

In this project, our task is to create a NLI (natural language interface) that interfaces with at least 3 distinct SQL and NoSQL database systems. We will design an NLI that performs CRUD operations using natural language inputs. In other words, users can write out requests to the NLI in natural/plain English – such as “how many rows does this dataset contain?” and other questions to explore the kinds of data that exist in the database. Users can also query our database to find out what tables exist, apply functions such as *match* and *sort*, and modify data in the database – i.e., “delete all data from the year 2022” or “update John’s location to California.”

Our code will leverage Large Language Model (LLM) APIs to translate user’s natural language requests into the appropriate SQL/NoSQL query. Depending on the user’s query, our NLI should also be able to create, update, and delete data. As a three-person group, we will also be using full-stack development to create a web platform for users to interact with ChatDB. This platform will allow users to ask plain language questions about the data, process requests in the back-end, and translate the results into output that is easily understandable to non-technical users. In other words, this platform will serve as an intelligent “ChatGPT” website users can visit for modifying/exploring/and querying our database (ChatDB).

For this project, we must make a proposal, submit progress reports, and demo our project to the rest of the class. The majority of our work lies in the implementation, which must be completed before our in-class demo.

Planned Implementation

1. Development Tools & Environment

- **Version Control:** GitHub (for code sharing)
- **Database Types:**
 - **SQL:** MySQL
 - **NoSQL:** options
 - MongoDB
 - Amazon DynamoDB (integrates with EC2, S3)
 - Cassandra (Apache)
 - *Consideration:* Firebase (with REST API integration)

2. Large Language Model (LLM) API Integration

- **Primary Choice:** Mistral API (via HuggingFace)
- **Decision Needed:** Determine if other LLM APIs might be a better fit through exploration

3. Dataset Selection & Preparation

- **Industry:** Entertainment/Media
- **Potential Large-Scale Datasets for Database:**

- Kaggle (for large, easy-to-access datasets)
- Movie Reviews
- Concerts (Scraped from LiveNation, TicketMaster)
- [Google Dataset Search](#) (for CSVs, Excel Sheets)
- **Tasks:**
 - Scraping and cleaning the dataset

4. NLI Development

- **Goal:** Develop a web app for the NLI interface
- **Task:**
 1. Figure out which type of LLM API we're using
 2. Front end for inputting responses and outputting results (need to have a beta NLI interface)
 3. Exploration of database
 - First get NLI working on SQL database, then the NoSQL database

Team Responsibilities:

- Rachel: project & timeline management
 - Primary role in developing and finalizing our written reports, as well as creating presentation for live demo
 - SQL server management
- Riya:
 - Researching/setting up LLM API to manage natural language inputs
 - Distributed data management
 - SQL server management
- Kevin:
 - NoSQL server management
 - Front/back end development

Timeline

Task #	Task Name	Assigned Team Member(s)- Team Responsibility	Start Date	End Date	Status	Notes - Planned Implementation
1	Project Kickoff & Planning	Everyone	02/06	02/12	C... ▾	<input checked="" type="checkbox"/> Initial brainstorming, project scope finalized <input checked="" type="checkbox"/> Research plan on implementing ChatDB <input checked="" type="checkbox"/> Identify database source <input type="checkbox"/> https://www.kaggle.com/datasets/saiivvvaam/most-watched-movies-and-tv-shows <input checked="" type="checkbox"/> Set up Git <input checked="" type="checkbox"/> Set up all the tools we need and do research how we are going to connect and use them <input checked="" type="checkbox"/> LLM
2	Submitting Proposal	Ri...	F...	📅 Date	C... ▾	100 POINTS
2	Dataset Preparation	Riya		2/12	C... ▾	<input checked="" type="checkbox"/> Data cleaning using Pandas, export to CSV & JSON Use pandas for cleaning, export to CSV <input checked="" type="checkbox"/> Find a way to make data structured for SQL and NoSQL during cleaning <input checked="" type="checkbox"/> Convert to JSON for NoSQL
3	SQL Database Implementation	Riya, Rachel		2/17	C... ▾	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> MySQL server set up (Riya and Rachel)
4	NoSQL Database Implementation	Kevin		2/17	C... ▾	<input checked="" type="checkbox"/> Set up MongoDB/DynamoDB/Cassandra <input checked="" type="checkbox"/> Mysql and postgresql
5	Natural Language Processing (LLM API)	Riya		2/24	C... ▾	<input checked="" type="checkbox"/> Integrate Mistral API via Hugging Face to NLP

	Integration)					
6	Backend Development	Kevin		3/5	C... ▾	<input checked="" type="checkbox"/> Develop API endpoints for SQL/NoSQL queries <input checked="" type="checkbox"/> Set up MySQL server, test queries <input checked="" type="checkbox"/> Integrate Hugging face open ai 3.5 (rachel) <input checked="" type="checkbox"/> Firebase setup
7	Mid-term Progress report	All	M...	Mar ...	C... ▾	<input checked="" type="checkbox"/> Write the report - rachel is writing - verify some information with Riya and Kevin <input checked="" type="checkbox"/> Tested out the modification request with Langchain -> Worked Task needed completed 03/07/2025 <input checked="" type="checkbox"/> Complete these parts in the report <ul style="list-style-type: none"> <input checked="" type="checkbox"/> What has been implemented so far? What are the next steps? <input checked="" type="checkbox"/> Describe your approach to converting natural language into database queries. <input checked="" type="checkbox"/> Ask kevin for his part of NoSQLI to confirm <input checked="" type="checkbox"/> Copy the timeline to the mid-progress report <input type="checkbox"/> Riya Berry <input checked="" type="checkbox"/> Duplicates in the database for tv shows (each season) <ul style="list-style-type: none"> <input checked="" type="checkbox"/> Choose entry with the highest watch time <input checked="" type="checkbox"/> Join in data about tv show episode and season count (TMDB) <ul style="list-style-type: none"> <input checked="" type="checkbox"/> Have to join at least two tables <input checked="" type="checkbox"/> Check that "titles" match when joining these tables <input type="checkbox"/>
8	Tools Setup Check+ Checking NLI			3/12	In... ▾	Task Needed completed before March 12th 2025 <input type="checkbox"/> Set up ChatDB with Langchain <ul style="list-style-type: none"> <input type="checkbox"/> Team Members need to make sure How to run it on local machine: <ul style="list-style-type: none"> <input type="checkbox"/> Langchain, OpenAI,

	(Necessity of having ec2)					<ul style="list-style-type: none"> <input type="checkbox"/> streamlite and Langchain to run it <input type="checkbox"/> All run in one file <input type="checkbox"/> Kevin Bui Could you upload the code to Google Collab for test run for Riya and Rachel in our loca machine <input type="checkbox"/> Build UI for query submission & visualization <input type="checkbox"/> Kevin - turn data to sql <input type="checkbox"/> Kevin idea: all run the demos, by uploading the code to run - run it locally <input type="checkbox"/> Integrate our specific database <input type="checkbox"/> Everyone: <ul style="list-style-type: none"> <input type="checkbox"/> Push everything to github (cleaned, nosql, setting up) <input type="checkbox"/> Riya Berry <ul style="list-style-type: none"> <input type="checkbox"/> Clean duplicate tv shows (titles should be unique) <ul style="list-style-type: none"> <input type="checkbox"/> Create a JSON file where the key is title <input type="checkbox"/> Clean new dataset about tv show episode / season count
9	Testing and Debugging			03/19	N... ▾	<input type="checkbox"/> Unit testing for query handling
10	Testing and Debugging			3/26	N... ▾	<input type="checkbox"/> Report structure, challenges & solutions
11	Testing and Debugging	👤 P...	📅 D.	Apr ...	N... ▾	
12	Testing and Debugging	👤 P...	📅 D.	Apr ...	N... ▾	
14	Testing and Debugging	👤 P...	📅 D.	Apr ...	N... ▾	<input type="checkbox"/> Prepare for the Demo

14	DEMO DAY	👤 P...	📅 D.	Apr ...	N... ▾	<input type="checkbox"/> Demo (inclass, 4/21 and 4/23, 10 points): <input type="checkbox"/> Give a live demo of your project. All project members should be present during the demo, presenting his/her contribution. If you're absent, we will assume that you have not contributed to the implementation of the project
	FINAL TOUCH_UP	👤 P...	📅 D.	Apr ...	N... ▾	<input type="checkbox"/> We should touch up on Final Report one more time before the MAY 09
15	FINAL REPORT	👤 P...	M...	📅 Date	N... ▾	<input type="checkbox"/> the final report should be comprehensive, details your design and implementation, and your learning experiences. <input type="checkbox"/> Implementation (due on your demo time, 60 points): note your project should be fully implemented before the demo. You should include in your final report a link to Google drive where you will upload your project codebase and documentations. Make sure you give access to your project folder.
		👤 P...	📅 D.	📅 Date	N... ▾	