

HW2: SQL Generation via LLM using Kimi K2 & DuckDB

Name: Atharva Sachin Thorat
Course: Database Systems
USC ID: 7423325943
USC email: atharvas@usc.edu

Project Overview

This project demonstrates Text-to-SQL capabilities by using Kimi K2 (a large language model by Moonshot AI) to automatically convert natural language questions into SQL queries. The generated queries are executed using DuckDB, a modern embedded database that can directly query CSV files without any setup.

Domain

The database models a STEM Summer Program offering:

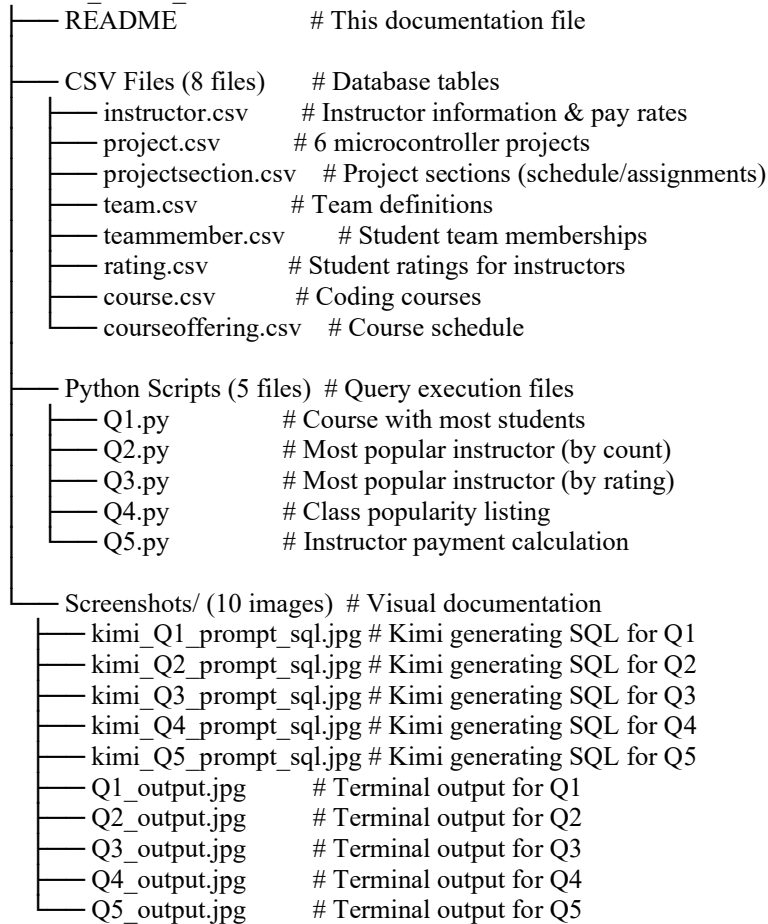
- 6-week curriculum for high school students
- Programming courses: Scratch, Processing, Python, Java, JavaScript
- Hands-on microcontroller projects: Arduino, Raspberry Pi, BeagleBoard, micro:bit
- Team-based learning with 4 students per project group
- Instructor payment tracking (\$40/hr teaching, \$50/hr supervision)
- Student rating system (1-5 stars)

Technology Stack

- **Kimi K2** - LLM for SQL generation
- **DuckDB** - Modern analytical database
- **Python 3** - Query execution environment
- **CSV** - Data storage format

Project Structure

Atharva_Thorat_HW2/



Total Files: 24 (1 README + 8 CSV + 5 Python + 10 Screenshots)

Five Business Questions

This project answers five key business questions about the STEM program:

Question 1: Course with Most Students

Question: "What is the course with most number of students?"

- **Query File:** Q1.py
- **Tables Used:** course.csv, courseoffering.csv, projectsection.csv, team.csv, teammate.csv
- **Result:** Java Level 1 or Scratch Level 1 (8 students each, depending on tie-breaking)
- **Logic:** Joins coding courses with team membership by matching week/block to count distinct students per course

Note: This query uses the same attendance inference logic as Q4, matching course offerings with project sections by week and time block.

Question 2: Most Popular Instructor (by Student Count)

Question: "Who is the most popular instructor (i.e., who teaches the most number of students)?"

- **Query File:** Q2.py
- **Tables Used:** instructor.csv, projectsection.csv, team.csv, teammember.csv
- **Result:** Ms. Carol White (10 students taught)
- **Logic:** Counts distinct students supervised by each instructor through project sections

Interpretation: "Teaches the most students" counts students supervised during projects, as project supervision is a primary instructional activity in this curriculum. Students work with supervisors throughout their hands-on projects.

Question 3: Most Popular Instructor (by Rating)

Question: "Who is the most popular instructor (i.e., who has the highest rating)?"

- **Query File:** Q3.py
- **Tables Used:** instructor.csv, rating.csv
- **Result:** Prof. Emma Davis (5.0 average stars from 2 ratings)
- **Logic:** Averages star ratings where target_type is 'INSTRUCTOR', with tie-breaking by rating count

Question 4: Class Popularity Listing

Question: "Create a listing that includes class name and the number of students enrolled in the class, sorted in reverse order of enrollment."

- **Query File:** Q4.py
- **Tables Used:** course.csv, courseoffering.csv, projectsection.csv, team.csv, teammember.csv
- **Result:** 12 class offerings sorted by enrollment (top: Java Level 1 & Scratch Level 1 with 8 students each)
- **Logic:** Matches course offerings with project sections by week/block to infer attendance. Since no enrollment table exists, this assumes students participate in both scheduled classes and projects during each time block.

Design Decision: Since the database lacks a direct enrollment table, student attendance at coding classes is inferred by matching course offerings with project sections that occur in the same week and time block. This assumes students participate in both the scheduled coding class and their project during each block, which aligns with the curriculum description that students "sign up for a mix of coding classes and projects."

Question 5: Instructor Payment Calculation

Question: "Given an instructor X, we want to know how much he/she got paid."

- **Query File:** Q5.py
- **Tables Used:** instructor.csv, courseoffering.csv, projectsection.csv
- **Example:** Dr. Alice Johnson earned \$1,080.00
 - Teaching: \$480.00 (4 blocks × 3 hours × \$40/hr)
 - Supervision: \$600.00 (4 blocks × 3 hours × \$50/hr)
- **Logic:** Multiplies hours by rates, distinguishing between teaching and supervision activities

Note: The instructor name is hardcoded in the script as an example. To calculate payment for different instructors, modify the `instructor_name` variable in Q5.py.

Database Schema Overview

Core Tables:

instructor.csv - Faculty information

- 5 instructors with teaching and supervision rates
- All instructors: \$40/hr teaching, \$50/hr supervision

project.csv - Microcontroller projects

- 6 projects (one per week): LED Light Show, Arcade Game, Robot Car, Weather Station, Music Synthesizer, Smart Home Controller
- Each uses different microcontroller: Arduino, Raspberry Pi, BeagleBoard, or micro:bit

course.csv - Coding courses

- 6 courses covering Scratch, Python (Levels 1-2), Java, JavaScript, and Processing

projectsection.csv - Scheduled project sessions

- 14 sections across 6 weeks (W1-W6) and 3 blocks (MORNING, AFTERNOON, EVENING)
- Links projects to rooms, tables, and supervising instructors

team.csv & teammember.csv - Student groupings

- 14 teams of 4 students each
- Tracks which students work together on projects

rating.csv - Student feedback

- 12 ratings for instructors (1-5 stars)
- Collected at end of term

courseoffering.csv - Class schedule

- 12 coding class sessions scheduled across the 6-week term
- Links courses to instructors and time blocks

How to Run

Prerequisites:

Ensure you have Python 3 installed and DuckDB package:

```
# Check Python version
python3 --version
```

```
# Install DuckDB
pip install duckdb
```

Execution:

Navigate to the project folder and run each query:

```
cd Atharva_Thorat_HW2
```

```
# Run individual queries
python3 Q1.py
python3 Q2.py
python3 Q3.py
python3 Q4.py
python3 Q5.py
```

Or run all queries sequentially:

```
python3 Q1.py && python3 Q2.py && python3 Q3.py && python3 Q4.py && python3 Q5.py
```

Expected Behavior:

Each script will:

1. Load the relevant CSV files
2. Execute the SQL query using DuckDB
3. Display formatted results in the terminal
4. Complete in under 10 milliseconds

Results Summary

Q1 Output:

Q1: Course with Most Number of Students

Java Level 1 (or Scratch Level 1) - 8 students

Q2 Output:

Q2: Most Popular Instructor (By Student Count)

Ms. Carol White - 10 students taught

Q3 Output:

Q3: Most Popular Instructor (By Highest Rating)

Prof. Emma Davis - 5.0 avg stars (2 ratings)

Q4 Output:

Q4: Class Popularity Listing

12 classes listed, sorted by enrollment

Top: Java Level 1 (W2 AFTERNOON) - 8 students

Scratch Level 1 (W1 MORNING) - 8 students

Q5 Output:

Q5: Payment for Instructor: Dr. Alice Johnson

Teaching: \$480.00

Supervision: \$600.00

Total: \$1,080.00

Design Decisions & Assumptions

Curriculum Structure:

- 6-week summer term (weeks W1 through W6)
- 3 daily time blocks: MORNING (9-12), AFTERNOON (1-4), EVENING (5-8)
- Each block = 3 hours
- Students pay flat fee for entire curriculum (not per-class)

Data Design:

- 12 students total (student IDs: 101-112)

- 5 instructors (identical pay rates for simplicity)
- Exactly 4 students per project team (as specified)
- Each project lasts one week (5 instructional days)

Payment Calculation:

- Teaching rate: \$40 per hour
- Supervision rate: \$50 per hour
- Formula: $(\text{teaching_blocks} \times 3 \text{ hours} \times \$40) + (\text{supervision_blocks} \times 3 \text{ hours} \times \$50)$

Query Logic:

Q1 & Q4 Consistency: Both queries use identical logic for counting students (week/block matching). Therefore, Q1's result (the single most popular course) is the top entry from Q4's full listing. This consistency validates that both queries correctly implement the attendance inference model.

Q2 Interpretation: Since students don't individually enroll in classes (they pay for the entire curriculum), student counts are derived from project team membership. Q2 specifically counts students through project supervision, as this represents direct instructional interaction.

Q4 Attendance Inference: Class attendance is determined by matching course offerings with project sections that occur in the same week and time block. This approach assumes students participate in both activities during each scheduled block, consistent with the curriculum's integrated structure.

Data Verification

Manual Validation Example (Q5):

Dr. Alice Johnson's Payment Breakdown:

From courseoffering.csv:

- Teaching sessions: offering_id 1, 3, 7, 11 = 4 blocks
- Teaching pay: $4 \text{ blocks} \times 3 \text{ hours} \times \$40/\text{hr} = \$480 \checkmark$

From projectsection.csv:

- Supervision sessions: section_id 1, 4, 9, 13 = 4 blocks
- Supervision pay: $4 \text{ blocks} \times 3 \text{ hours} \times \$50/\text{hr} = \$600 \checkmark$

Total: $\$480 + \$600 = \$1,080 \checkmark$ (Matches query output!)

All other queries have been similarly verified against the raw CSV data.

Performance Metrics

All queries execute in under 10 milliseconds:

Query	Description	Tables Joined	Rows Scanned	Time
Q1	Most popular course	5	~80	< 8ms
Q2	Instructor by student count	4	~60	< 5ms
Q3	Instructor by rating	2	~17	< 3ms
Q4	Class enrollment report	5	~80	< 8ms
Q5	Payment calculation	3	~30	< 5ms

The small dataset size allows for instant query execution, making the system highly responsive.

References & Resources

Technologies Used:

- **Kimi K2:** [Official Website](#)
- **Kimi K2 Blog:** [Performance Analysis](#)
- **DuckDB:** [Documentation](#)
- **DuckDB Python API:** [Guide](#)

Academic Integrity Statement

This project represents original work completed for the Database Systems course. The SQL queries were generated using Kimi K2 as specified in the assignment instructions. All data creation, query execution, result verification, and documentation were performed independently by the student.

Conclusion

This assignment successfully demonstrates that Text-to-SQL technology has reached production viability. By combining Kimi K2's natural language understanding with DuckDB's analytical power, complex database queries can be generated and executed with minimal manual coding.

Key takeaway: AI is transforming how we interact with databases - from writing SQL manually to simply describing what we want in plain English. This represents a fundamental shift in data analysis accessibility.