

CSC3170 Project: Database for Models and Datasets (Draft)

- Notation: different parts to complete.
 - [s] secure
 - [d] database schema
 - [f] frontend
 - [i] data insight
 - [?] other todo list

1. Introduction and motivation

1.1. Introduction

- Our project is a database for machine learning models and datasets.
 - **Basic database operations:** It allows users to browse the information about the models and datasets, upload and download models and datasets.
 - **Schema:** Apart from the basic schemas such as dataset, model, user, we also included schemas that are especially helpful for machine learning developers, such as tables describing the modular structures of different architecture of models (CNN, RNN, Transformer).
 - **GUI:** A beautifully designed graphic user interface is implemented, where users and administrators can perform multiple types of operations.
 - **LLM:** An LLM agent is implemented, to translate user's natural language query into SQL language. User can also customize their query by selecting different tables and different fields.
 - **Security:** Methods are implemented to protect data security,

1.2. Motivation

- We are motivated by [huggingface](#), one of the most influential platform in the AI community that facilitates the sharing and collaboration of machine learning models and datasets.

1.3. How to run our code

- **Step 1-3 has to be done ONLY when running it at the first time; if it's not the first time, you can skip 1-3, and also can skip 4 if you don't need to initialize the database.**
- [q] *update this part after startup.py is finished.*

1. Install dependencies according to `requirement.txt` [?]
2. Create an `.env` file at the root directory of the project, and add the following lines to it (repace `$your_api_key` and `$your_base_url` with your own values):

```
# -----database-----  
DB_USERNAME=root  
DB_PASSWORD=123  
DB_HOST=0.0.0.0  
DB_PORT=3306  
TARGET_DB=openmodelhub  
# -----agent-----  
API_KEY=$your_api_key  
BASE_URL=$your_base_url
```

3. Test connection by running `database/db_connection_check.py`.

4. Initialize the database with the records stored in `database/records/demo.json`, by running:

```
database/load_data.py
```

- then you'll be asked to choose a .json file stored in `database/records` to initialize it; just choose `demo.json`.

5. Run the GUI:

```
streamlit run frontend/app.py
```

6. Login as common user or admin

- Login to admin with **username: admin, password: admin**.
- After logging in as admin, you can see the list of all users in the page `user management`. Note that some users are admin, too, as indicated on the page.
- Every user's password is admin.
- You can register your own user, too.

2. Design and implementation

2.0. Project Structure

- our project is composed of the following components:
 1. Database.
 2. Data.
 3. Frontend.
 4. Agent.
 5. Security.
 6. Data analysis.

2.1. Database

- [d][THE WHOLE PART needs fact-checking!! whether my description is accurate?]

Schema Design

- Our database follows the relational model and the 4th normal form.
- Our schema are as follows:
- [d] [please insert a markdown format table here to show the schema. can be generated from our slides.]
- [?] llm optimized design

Implementation

- In `database/database_schema.py`, schemas are represented by python classes.
- In `database/database_interface.py`, we have encapsulated interfaces to perform SQL operations safely. Therefore, in other programs where we have to execute SQL, we can call an encapsulated functions instead of executing the SQL operations directly.

2.2. Data

Initialization

- We created a set of records to initialize our database; although more records can be inserted to or deleted from the database during use. It is stored in `database/records/demo.json`, and can be run by `database/load_data.py`, as indicated previously.
- The records consist of:
 1. 12 affiliations;
 2. 28 users from these affiliations;
 3. 100 datasets;
 4. 92 models.
- The models' names, corresponding architecture, media type, train method (fine-tuned or pre-trained) are real; the dataset's names and media types are real, because they are copied from models and datasets that are actually posted to [huggingface](#). However, some other attributes, such as parameter number and authors, are made up.

Upload and Download

- `database/load_data.py` can initialize the database by inserting records stored in json formats, containing entities among `affiliation`, `user`, `dataset`, `model`.
- [f] [should explain how to download and what programs are responsible.]

2.3. Frontend

login/regsiter

1. common user login
2. common user register and login

- 3. admin login: has some pages that common users don't have.
 - username: admin; password: admin.
- The pages visible to a common user / an admin is different.

type

user

admin

sidebar

Open Model Hub

Welcome, Lewis!

Log Out

Navigation Menu

- Home
- Model Repository
- Datasets

Open Model Hub

Welcome, admin!

Log Out

Navigation Menu

- Home
- Model Repository
- Datasets
- User Management
- data insight

Home page

- can export and download ata.

Model/Dataset Repository page

- the following screenshots are from the model page; but the dataset page is very similar.

LLM assisted search, with specifying the entity in the drop-down box

upload model

click "view details", and 2 tables representing the detailed information of that model will be displayed. paging are implemented for improved user experiment.

Model Repository

Search Models

Model Type

Model ID

Model Name

Search

Model Type

Model ID

Model Name

Search

Model Type

Model ID

Model Name

Search

Model Type

Model ID

Model Name

Search

Model Repository

Search Models

Model Type

Model ID

Model Name

Search

Model Type

Model ID

Model Name

Search

Model Type

Model ID

Model Name

Search

Model Type

Model ID

Model Name

Search

Open Model Hub

Welcome, admin!

Log Out

Navigation Menu

- Home
- Model Repository
- Datasets
- User Management
- data insight

SimpleRNN-TextGen-5

View Details

Model Details - SimpleRNN-TextGen-5

Basic Information	
Model ID	SimpleRNN-TextGen-5
Architecture Type	RNN
Parameter Count	10,000,000
Training Type	Supervised

Related Information	
Type	generator
Supported Tasks	Text Generation
Authors	OpenAI Research
Related Datasets	gpt2_text_1.1

(Admin Privilege) User Management / Data Insight

4. (Admin Privilege) User Management

5. (Admin Privilege) Data Insights

data insight, page 2

illustration of the analysis on the data in the database.

create/edit user

User Management

Search Users

Enter natural language query

Select Type

id

Select Field

Enter Query Value

user_id

Search

+ Add New User

Username*

Rain

Password*

☒ Admin Privileges

Organization

The Chinese University of Hong Kong, Shenzhen

Create User

User ID	Username	Organization	Admin Status
1	Bradley	The Chinese University of Hong Kong, Shenzhen	<input checked="" type="checkbox"/>
2	Resaland	The Chinese University of Hong Kong, Shenzhen	<input checked="" type="checkbox"/>
3	Alice	Tsinghua University	<input checked="" type="checkbox"/>
4	Bob	Tsinghua University	<input checked="" type="checkbox"/>
5	Charlie	Peking University	<input checked="" type="checkbox"/>
6	David	Peking University	<input checked="" type="checkbox"/>
7	Eve	Fudan University	<input checked="" type="checkbox"/>
8	Frank	Fudan University	<input checked="" type="checkbox"/>
9	Grace	Shanghai Jiaotong University	<input checked="" type="checkbox"/>
10	Hannah	Shanghai Jiaotong University	<input checked="" type="checkbox"/>

Admin features

Select Field

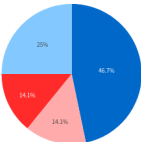
Enter Query Value

user_id

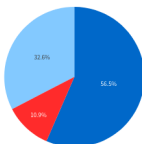
Search

Model

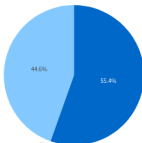
number & percentage of each media_type of model



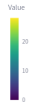
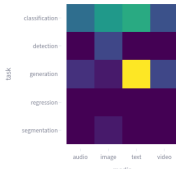
number & percentage of each arch_name of model



number & percentage of each trainname of model



media_task_relation



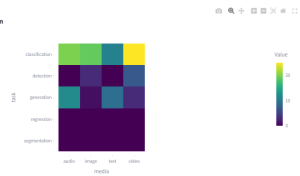
Also illustration.

param_num summary:

value	
max	1700000000
min	3500000
mean	138770860.2629
std	362491249.1381

dataset

media_task_relation

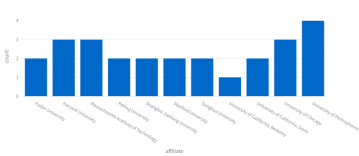


param_num summary:

value	
max	300000000
min	5
mean	11075168.45
std	38033822.1054

user

user in different affiliations



2.4. Agent

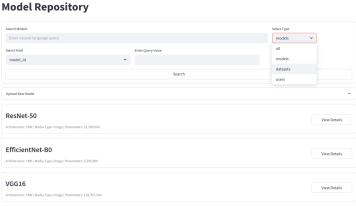
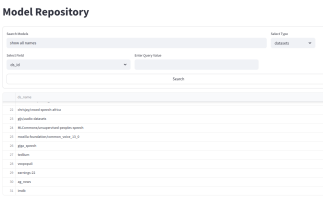
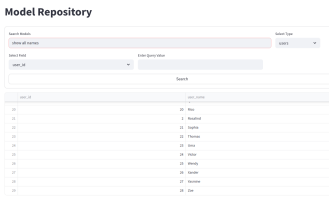
Implementation

- We incorporated gpt-4o as an LLM agent that translates user's natural language input into SQL queries.
- **Input/Output**
 - Input includes a natural language query, and a integer specifying the type of entity user is asking for. A corresponding string will be appended to the natural language query. This integer is by default 0, indicating no specific constraints.
 - Output: a dictionary, consisting of:
 - an error code indicating whether a grammatically correct sql is generated;
 - a SQL query generated
 - the result of the SQL query

- **System prompt**
 - **Schema:** In the system prompt, we describe our database, the integrity constraints, and other information required.
 - **Synonyms:** In practice, we find it necessary to add some synonyms to help agent understand user's needs in this context. For example, if user asks for a **language model**, user is referring to **models where media_type includes 'text'**
 - **Instance type:** The constraints on the type of entity user's asking for is also indicated in the system prompt.
- **2-stage error-detection leveraging agent's self-correction:**
 - After the SQL is created, it will be executed to check its grammatical correctness, instead of directly returning the SQL.
 - **If incorrect, agent will perform another attempt to generate SQL, based on the previous failure.** However, if it fails again, no more attempts will be made.

Demonstration: using LLM assisted search in the GUI

- Feature: generate query of a specific entity.

entity type	model	dataset	user
query: show all names			

- other queries

Query	Find all transformer models	(same as previous)	top 10 users with the most published datasets

QueryFind all transformer models

(same as previous)

top 10 users with the most published datasets

More complicated queries can also be executed.

Result can be represented in a table.

Result

Model Repository

model_id	model_name	dataset_size	model_type	arch_name	description	status
0	AI-001	1000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
1	AI-002	1000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
2	AI-003	8000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
3	AI-004	11000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
4	AI-005	10000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
5	AI-006	11000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
6	AI-007	1700000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
7	AI-008	20000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
8	AI-009	10000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000
9	AI-010	8000000	TEXT	TRANSFORMER	Transformer PECTRAL	6500000000

Can also view the corresponding SQL query.

Query Details

```

{
  "Natural Language Query": "Find all transformer models",
  "Generated SQL": "SELECT * FROM model WHERE arch_name = 'Transformer'",
  "Error Code": 0,
  "Has Results": true,
  "Error Message": null,
  "Query Results": [
    {
      "model_id": 0,
      "model_name": "AI-001",
      "dataset_size": 1000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 1,
      "model_name": "AI-002",
      "dataset_size": 1000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 2,
      "model_name": "AI-003",
      "dataset_size": 8000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 3,
      "model_name": "AI-004",
      "dataset_size": 11000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 4,
      "model_name": "AI-005",
      "dataset_size": 10000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 5,
      "model_name": "AI-006",
      "dataset_size": 11000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 6,
      "model_name": "AI-007",
      "dataset_size": 1700000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 7,
      "model_name": "AI-008",
      "dataset_size": 20000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 8,
      "model_name": "AI-009",
      "dataset_size": 10000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    },
    {
      "model_id": 9,
      "model_name": "AI-010",
      "dataset_size": 8000000,
      "model_type": "TEXT",
      "arch_name": "TRANSFORMER",
      "description": "Transformer PECTRAL",
      "status": "6500000000"
    }
  ]
}
```

Model Repository

Search Models: top 10 users with the most published datasets. Select Type: all. Select Field: user_id. Enter Query Value: user_id. Search

user_name	dataset_count
0 Grace	15
1 Haruto	15
2 Fange	14
3 Hay	13
4 Eve	13
5 Wendy	12
6 Jack	12
7 Bob	12
8 Yasmine	11
9 Rosalind	11

Query Details

```

{
  "Natural Language Query": "top 10 users with the most published datasets",
  "Generated SQL": "SELECT user_name, COUNT(dataset_id) AS dataset_count FROM user u JOIN dataset ds ON u.user_id = ds.user_id GROUP BY u.user_id ORDER BY dataset_count DESC LIMIT 10",
  "Error Code": 0
}
```

2.5. Security

- [s]

2.6. Data Insight

- [i]

3. Conclusion and self-evaluation

3.1. Conclusion

- We has completed task [?] indicated in the project guideline.
- [?] mention detailed implementation here.

3.2. Self-Evaluation

- Work division is as follows: (members' names follows alphabetical order)

Yimeng Teng

- Implemented the entire agent part. Generated test cases to evaluate and refine it.
- Collaborated with Linyong Gan to generate demo.json, which contains sufficient amounts of records for initializing the database.
- Collaborated with Wentao Lin in implementing a data loader that load json files and insert records to the database. Designed the first version and help completed the final version.
- Participated in the formulation of the database schema (but not the implementation).

4. References

- <https://huggingface.co/>
- Feistel, H. (1973). Cryptography and computer privacy. Scientific american, 228(5), 15-23.
- Rivest, R. L., Shamir, A., & Adleman, L. (1978). A method for obtaining digital signatures and public-key cryptosystems. Communications of the ACM, 21(2), 120-126.

5. Appendices

[?] what to include