

Project Final Report

Check our database on [Github](#)!

Project Final Report

- Introduction

- Planned Implementation

 - SQL Database

 - NoSQL Database

 - CLI

- Architecture Design

 - Overview

 - Key Components

 - Flow Diagram

- Implementation

 - Functionalities

 - Tech Stack

 - Implementation Screenshots

- Learning Outcomes

- Individual Contribution

- Conclusion

- Future Scope

Introduction

The AvA Database project integrates SQL and NoSQL databases, offering versatile data management capabilities. It supports various field types and customizable settings, including default values and file extensions. The system can import data from CSV and JSON files, automatically adapting to the appropriate database format. A custom logging module enhances monitoring and debugging. Detailed usage guidelines and setup instructions ensure easy operation. The project includes constants and definitions to standardize operations and parameters across the system.

A Flask web server handles database operations like insert, update, and delete (and more advanced function). The project context manages configurations and database connections. Advanced field processing techniques support complex queries and database operations. Table management includes creation, deletion, and metadata handling, ensuring data integrity. Metadata handling is critical for maintaining data structure in both SQL and NoSQL systems. The table manipulation module implements advanced functionalities like sorting, grouping, and joining, providing robust data processing capabilities. This project represents a comprehensive solution for efficient and complex data management across diverse database systems.

Planned Implementation

SQL Database

Data Model

- Each row in a table represents an entry.
- Tables are organized in a fixed schema that outlines the type and constraints of each column.
- Each entry will be assigned a PK (Primary Key) per table.
- Foreign Keys will be used to represent relationships between tables.

Storage

Data will be stored in tables, Each table will be a separate file or set of files managed by the SQL database engine. We don't manage these files directly.

Programming Language

SQL-based database will be implemented by Java(after our discussion, now we have a unified language for both SQL and NOSQL), since the language has better performance.

Supported Operations

SQL naturally supports JOIN operations, which allows for combining rows from two or more tables based on a related column.

The other operations mentioned

- projection
- filtering
- grouping
- aggregation
- ordering
- insertion
- updation
- deletion are also fully supported in SQL databases.

"JOIN" will be supported in SQL database, since there are multiple tables can be operated.

Big Data

- SQL database is able to store large volumes of data efficiently, but after the amount of data reaches to million, optimizations such as indexing, caching, and partitioning can be used(). (we now use chunk to store large volumes of data)
- SQL database allows transactions, locking, and other concurrency controls, to manage multiple simultaneous operations.(no more concurrency is considered)

Other Details

SQL database will have a local interface to issue commands and show query results.

NoSQL Database

Data Model

- Each entry will be presented as an object. Each object will be stored as a JSON document.
- Objects in the same class will be stored together as a collection.
- Every entry will be assigned with a unique ID.
- The entry will hold the IDs of other entries as JSON fields.

Storage

We suppose every single entry can be loaded into memory. Each JSON document will be stored in a distinct file. Objects in the same collection will be stored under the same folder.

Programming Language

NoSQL database will be implemented with Python.

Supported Operations

"JOIN" is not supported by our NoSQL database (No tables to join for NoSQL!)

Other operations are supported, including:

- projection
- filtering
- grouping
- aggregation
- ordering
- Insertion, updation, and deletion are all supported.

Big Data

Only one JSON document will be loaded into memory for each thread. The JSON object will be cleared or replaced before loading the next JSON document. Only IDs of different objects will be kept in memory at the same time.

Other Details

NoSQL Database will be served as an HTTP server and handle HTTP requests only.

CLI

It's designed as a user interface to communicate with SQL or NoSQL servers. Queries will be sent by HTTP protocol, while the query results will be parsed and presented in standard output.

The CLI will be implemented with Python.

Architecture Design

Overview

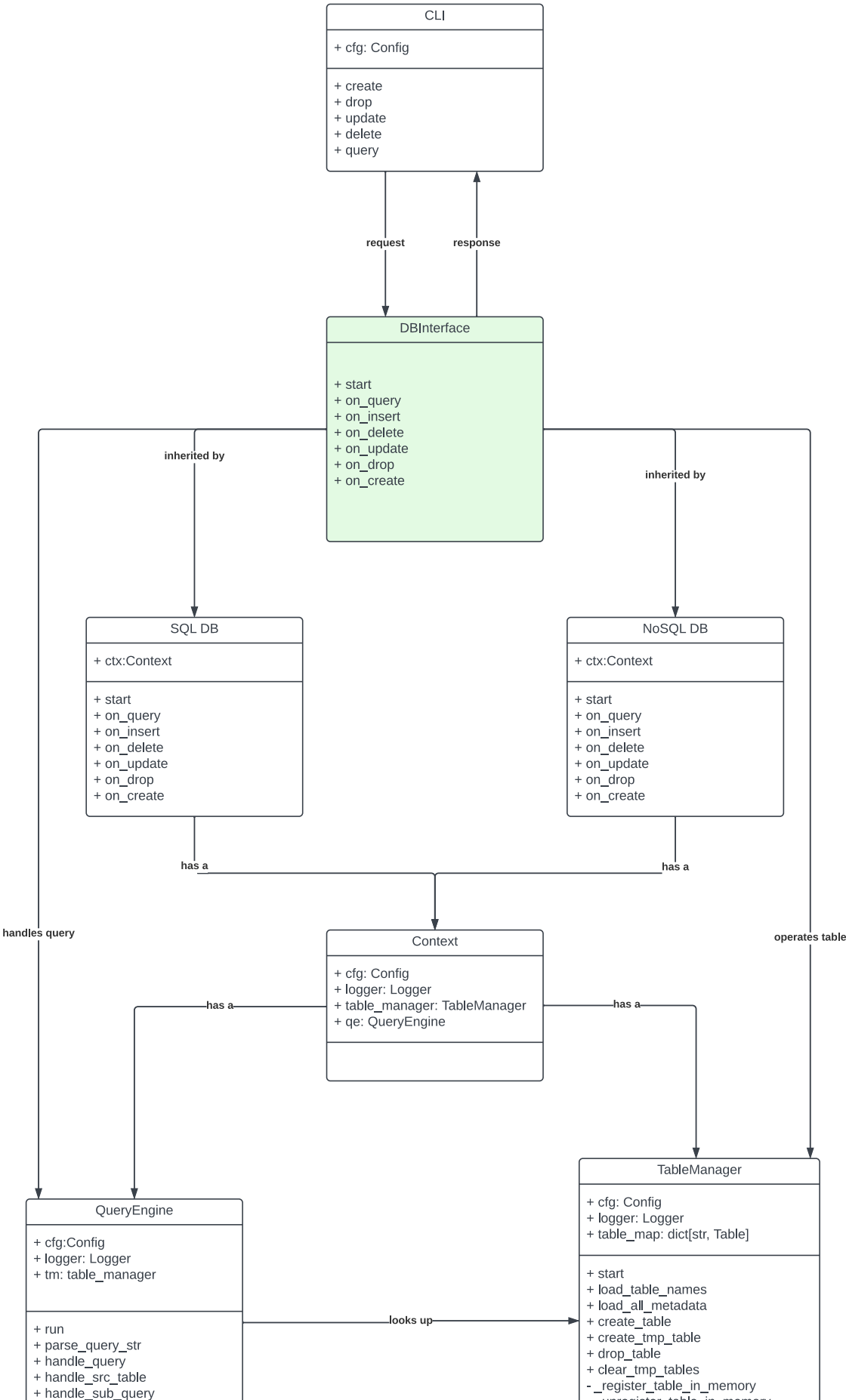
The project architecture is centered around the integration of SQL and NoSQL databases, providing a unified interface for data manipulation and query processing. It leverages Python for backend development, Flask for API handling, and a custom query language designed for optimized data handling in limited memory environments.

Key Components

1. **Database Configuration (DBConfig):** Manages database settings, including type (SQL or NoSQL), port, directories for tables and metadata, supported field types, and file extensions.
2. **Data Import (import_data.py):** Handles data import from CSV and JSON files, auto-detecting column types for CSV files and loading JSON files for NoSQL databases.
3. **Logging (logger.py):** Provides a logging mechanism for tracking and debugging. Custom loggers for different components ensure detailed logging.
4. **Server Management (app/server.py):** Uses Flask to handle API requests, enabling operations like file retrieval, data insertion, updating, deletion, table creation, and dropping.
5. **Context Management (app/common/context/context.py):** Manages the application's runtime context, including configurations, database instances, table managers, and query engines.
6. **Table Management (app/common/table/table_manager.py):** Manages tables, including creation, deletion, and metadata handling. Ensures data consistency and manages temporary tables for complex operations.
7. **Command Line Interface (cli/cli.py):**
Offers an interactive interface for users to directly communicate with the server. It translates command-line inputs into HTTP requests that are processed by `server.py`.
8. **Query Engine(app/common/query/query_engine.py):**
Processes and executes queries by interfacing with the SQL or NoSQL database. It plays a crucial role in the efficient handling of database queries.
9. **Chunk Management(app/common/table/chunk_manager.py) :**
Critical for handling large datasets in a limited memory space. This component ensures data is processed and stored in manageable chunks, allowing efficient data manipulation without overloading the system's memory. This approach is particularly important for operations like data import, query processing, and table management.

Flow Diagram

The flow diagram (presented below) illustrates the interaction between these components.



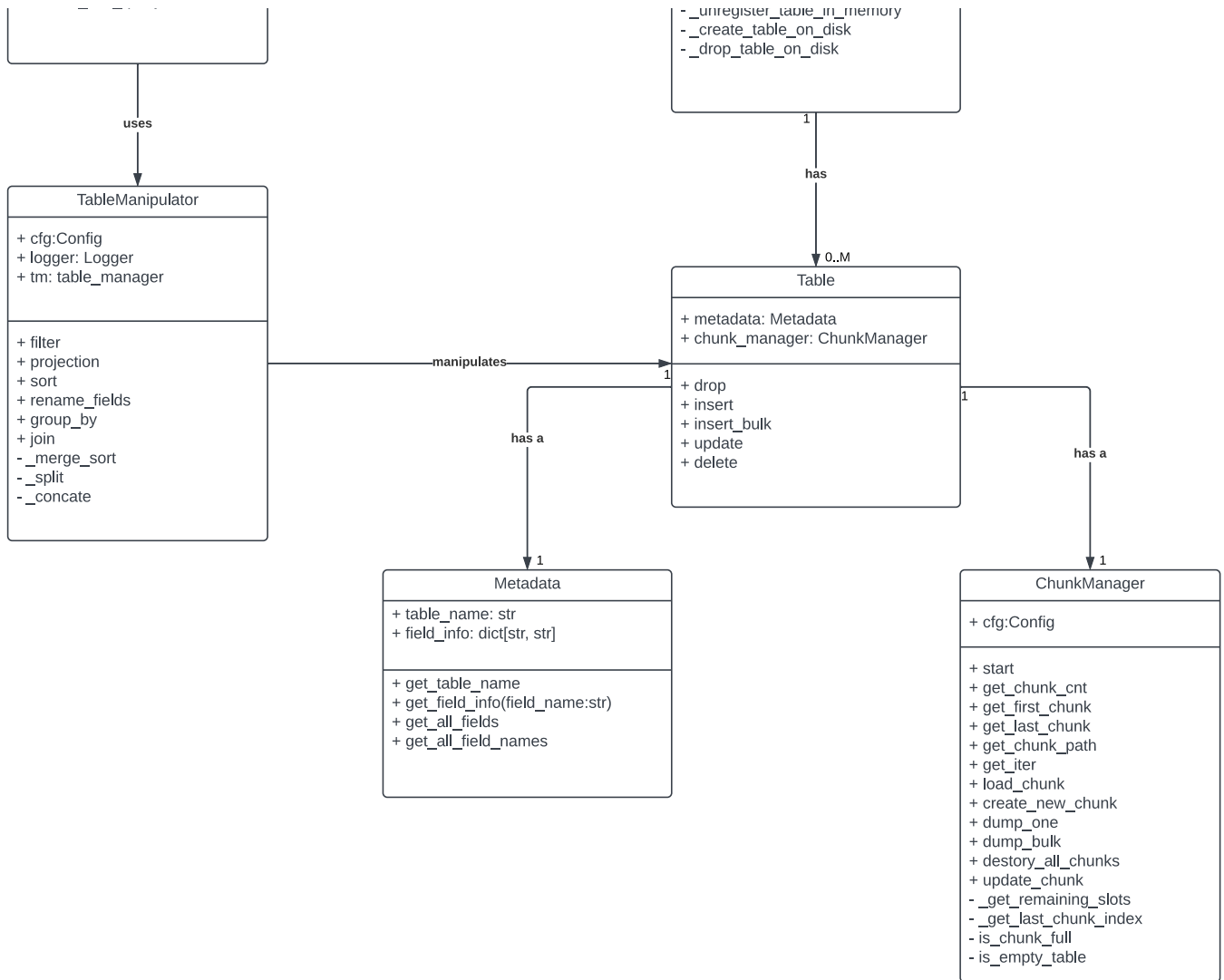


Diagram Description

1. CLI

- interacts with users and send table operations or queries to flask server
- formats query result for human to read

2. DBInterface

- it's a flask server
- parses requests (CRUD operations, table creation/dropping)
- decides which database the CLI is talking to and send the request to corresponding database to handle
- formats the query result before sending back to CLI

3. Context

- holds instance of query engine, table manager, logger, and config
- gets initialized when starting the flask server
- be passed to all other modules to use

4. TableManager

- loads all info of tables from the disk and registers them in memory

- handles creation, deletion of tables
- offers apis to create a temporary table, all temporary tables will be deleted when starting the flask server

5. QueryEngine

- parses queries and decides the execution plan
- uses `TableManipulator` to execute each step in the execution plan

6. TableManipulator

- implements all operations including: projection/renaming/filtering/ for both SQL and NoSQL
- uses `ChunkManager` to create/load/update/delete chunks of tables on the disk

7. ChunkManager

- manages all chunks of a table on the disk
- offers CRUD apis for chunks

8. Logging:

- logs all operations across modules
- useful for debugging

Implementation

Functionalities

Projection: This functionality involves selecting specific columns from a table. In the context of the AvA Database, the projection operation allows for tailoring the data output to include only the desired columns. This is particularly useful when dealing with large datasets, where you might only need a subset of the available data fields.

Filtering: Filtering refers to the capability to select data records based on specific criteria. The AvA Database can perform efficient filtering operations, allowing users to retrieve only those records that meet specified conditions. This is crucial for querying and analyzing data based on specific requirements.

Grouping: Grouping is the process of collating data records based on one or more columns. This feature in the AvA Database project is essential for organizing data into subsets for further analysis, such as aggregation or detailed examination of grouped data.

Aggregation: Aggregation involves computing summary statistics or other analytical operations on groups of data. The AvA Database supports various aggregation functions like count, sum, average, minimum, and maximum. These operations are particularly valuable in data analysis, allowing for the extraction of meaningful insights from grouped data.

Ordering: This function sorts data based on one or more columns, either in ascending or descending order. In the AvA Database, ordering is implemented to organize data in a specific sequence, which is crucial for reporting, data visualization, and further analytical processing.

Insertion: The database system allows for the insertion of new records into a database table. This functionality is fundamental to any database system, enabling the addition of new data into the existing dataset.

Update: The AvA Database provides the capability to update existing records in the database. This feature is essential for maintaining the accuracy and relevance of the data stored in the database.

Deletion: This functionality pertains to the removal of existing records from the database. The ability to delete data is crucial for data management and ensuring that the database only contains current and relevant information.

Sorting (Ordering): The AvA Database's sorting functionality allows for the organization of data based on specified criteria. This can be done on one or more columns and can be configured to sort data in either ascending or descending order. Sorting is essential for data analysis, as it arranges data in a meaningful order, making it easier to understand, compare, and visualize. The database might employ algorithms like merge sort, which are particularly efficient for large datasets, ensuring that the sorting process is both fast and reliable.

Join: Join operations are crucial in relational database management systems, and the AvA Database supports them to combine data from two or more tables. This feature is based on a common field shared between the tables. The types of joins can include inner joins, left/right outer joins, and full outer joins. Joining is particularly important for relational databases as it allows for the relational model to be effectively utilized, enabling users to construct complex queries that gather data across multiple tables. This capability is vital for scenarios where data is normalized and spread across different tables but needs to be viewed or analyzed together.

Tech Stack

1. The project is mainly developed with `Python`. Some `Bash` scripts are also used for starting servers and cli.
2. `Flask` framework is used to serve both SQL and NoSQL database engines.
3. We use `Postman` to test apis offered by servers.
4. `Pytest` is used for unittests.
5. Version control system is `Git`. We use `Github` to collaborate.
6. We use standard libs of `Python` to implement the entire project. No other third-party libs are used. (e.g. pandas)

Implementation Screenshots

CLI


```

cli > cli.py
1 import json
2 import requests
3 import config
4 import constant
5 import logger
6 import tempfile
7
8 logger = logger.get_logger(constant.CLI_NAME)
9
10
11 > def download_file(url, query) -> str:--
31
32
33 > def print_file(file_name: str):--
40
41
42 supported_apis = [constant.REQUEST_KEY_QUERY, constant.REQUEST_KEY_CREATE, constant.REQUEST_KEY_DROP, constant.REQUEST_KEY_DELETE, constant.REQUEST_KEY_UPDATE]
43
44
45 > def is_request_valid(request_json) -> bool:--
59
60
61 > def parse_json_input(user_input: str):--
72
73
74 > def get_cfg(request_json) -> config.DBConfig | None:--
78
79
80 def main():
81     logger.info("User connects")
82     print(constant.CLI_WELCOME)
83     while True:
84         user_input = input(constant.CLI_PROMPT)
85         if user_input.lower() == 'exit':
86             logger.info("User exits")
87             break
88         try:
89             request_json = parse_json_input(user_input)
90             if request_json is None:
91                 logger.error("Invalid user input: {}".format(user_input))
92                 continue
93             cfg = get_cfg(request_json)
94             if constant.REQUEST_KEY_QUERY in request_json:
95                 tmp_file = download_file('http://localhost:{}'.format(cfg.get_port()), json.dumps(request_json.get(constant.REQUEST_KEY_QUERY)))
96                 if tmp_file == "":
97                     print("Get result from database failed")
98                     continue
99                 print_file(tmp_file)
100             else:
101                 method = [key for key in request_json.keys() if key in supported_apis][0]
102                 response = requests.get("http://localhost:{}/{}".format(cfg.get_port(), method), json=json.dumps(request_json[method]))
103                 print(response.text)
104             except Exception as e:
105                 logger.error("An unexpected error occurred: {}".format(e))
106                 continue
107
108
109 if __name__ == "__main__":
110     main()
111

```

Flask Server

```

app > server.py
13 from app.common.context.context import Context
14 from typing import cast
15
16 app = Flask(__name__)
17
18
19 @app.route('/')
20 def get_file():
21     query = request.json
22     result, status = g.ctx.get_db().on_query(query)
23
24     if not status.ok():
25         return "failed to process query: {}, due to {}".format(query, status.code()), 400
26
27     if result is None:
28         return "query result is empty for query {}".format(query), 200
29
30     result = cast(QueryResult, result)
31     g.ctx.logger.info("result dat can be found under {}".format(result.get_result_file_path()))
32
33     return send_file(result.get_result_file_path(), as_attachment=True, mimetype='text/plain')
34
35
36 @app.route('/insert')
37 > def insert(): ...
38
39
40
41
42
43
44 @app.route('/update')
45 > def update(): ...
46
47
48
49
50
51
52 @app.route('/delete')
53 > def delete(): ...
54
55
56
57
58
59
60 @app.route('/drop')
61 > def drop(): ...
62
63
64
65
66
67
68 @app.route('/create')
69 > def create(): ...
70
71
72
73
74
75
76 > def check_args(): ...
77
78
79
80
81
82 def get_db(logger: logger.Logger, cfg: config.DBConfig):
83     if cfg.is_sql():
84         return SQLDBFactory.instance()
85     elif cfg.is_nosql():
86         return NosqlDBFactory.instance()
87     else:
88         logger.error("unrecognized database type {}, unable to start".format(ctx.get_cfg().get_db_type()))
89         return None
90
91
92 def start_db(ctx: Context) -> Status:
93     status = ctx.get_db().start(ctx)
94     if not status.ok():
95         ctx.logger().error("failed to start {}, due to {}".format(ctx.get_cfg().get_db_type(), status))
96         return START_FAILED
97
98     return OK
99
100
101 @app.before_request
102 def before_request():
103     # global context which could be used in entire life cycle
104     g.ctx = ctx
105
106
107 if __name__ == '__main__':

```

QueryEngine: handle query

```
app > common > query > query_engine.py
86
87 def handle_query(self, q: dict) -> (Table | None, Status):
88     # 1. handle src table, may contain subquery and joining
89     src_table, status = self.handle_src_table(q)
90     if not status.ok():
91         self.logger.error("failed to handle query {} due to failed to parse src_tables".format(q))
92         return None, INVALID_ARGUMENT
93
94     res_table = src_table
95     # 2. handle group by
96     if constant.QUERY_GROUP_BY in q:
97         columns = q[constant.QUERY_GROUP_BY]
98         if len(columns) > 1:
99             self.logger.error("failed to handle more than one group by in query {}".format(q))
100             return None, NOT_IMPLEMENTED
101         reduce_options = []
102         for column in q[constant.QUERY_DESIRED_COLUMNS_KEY]:
103             if FieldNameProcessor.get_suffix(column) != "":
104                 reduce_options.append(ReduceOption(column, ReduceOperation[FieldNameProcessor.get_suffix(column)]))
105         self.logger.info("grouping by table {}".format(res_table.name))
106         res_table = TableManipulator.group_by(res_table, GroupByOption(columns[0], reduce_options))
107         self.logger.info("group by on table {} is finished".format(res_table.name))
108
109     # 3. filter data
110     if constant.QUERY_FILTER_KEY in q:
111         expression = q[constant.QUERY_FILTER_KEY]
112         self.logger.info("filtering table {}".format(res_table.name))
113         res_table = TableManipulator.filter(res_table, Selector(expression))
114         self.logger.info("table is filtered, new_table is {}".format(res_table.name))
115
116     # 4. sorting
117     if constant.QUERY_ORDER_BY_KEY in q:
118         sort_options: list[dict] = q[constant.QUERY_ORDER_BY_KEY]
119         if len(sort_options) > 1:
120             self.logger.error("doesn't support order by 2 columns".format(q[constant.QUERY_ORDER_BY_KEY]))
121             return None, NOT_IMPLEMENTED
122         self.logger.info("sorting table {}".format(res_table.name))
123         res_table = TableManipulator.sort(res_table, [
124             SortOption(
125                 option.get(constant.QUERY_ORDER_BY_COLUMN_KEY),
126                 option.get(constant.QUERY_ORDER_BY_ASC_KEY)) if constant.QUERY_ORDER_BY_ASC_KEY in option else True # default is asc
127             for option in sort_options
128         ])
129         self.logger.info("table is sorted, new table: {}".format(res_table.name))
130
131     # 5. projection
132     # keep all columns by default
133     if constant.QUERY_DESIRED_COLUMNS_KEY in q:
134         columns = q[constant.QUERY_DESIRED_COLUMNS_KEY]
135         modified_columns = []
136         for column in columns:
137             if FieldNameProcessor.get_inner_prefix(column) in self.prefix_map:
138                 modified_columns.append(FieldNameProcessor.replace_inner_prefix(column, self.prefix_map[FieldNameProcessor.get_inner_prefix(column)]))
139             else:
140                 modified_columns.append(column)
141         if len(modified_columns) == 0:
142             self.logger.warn("empty table due to empty desired_columns field in query: {}".format(q))
143             return None, INVALID_ARGUMENT
144         self.logger.info("projecting table {} to columns: {}".format(res_table.name, modified_columns))
145         res_table = TableManipulator.projection(res_table, modified_columns)
146         self.logger.info("table is projected to columns: {}, new table is {}".format(modified_columns, res_table))
147
148     return res_table, OK
149
150 def run(self, query: str) -> (Table | None, Status):
151     q, status = self.parse_query_str(query)
152     if not status.ok():
153         return None, INVALID_ARGUMENT
154
155     return self.handle_query(q)
156
```

TableManipulator: merge_sort

```
app > common > table > manipulator.py
```

```
107
108     @staticmethod
109 > def filter(src_table: Table, selector: Selector) -> Table: ...
142
143     # TODO support nosql, which has nested fields i.e. a.b.c
144     @staticmethod
145 > def projection(src_table: Table, desired_column: list[str]) -> Table: ...
189
190     @staticmethod
191 > def _sort_each_chunk(src_table: Table, column: str, is_asc: bool) -> Table: ...
207
208     @staticmethod
209 def _merge_sort(src_table: Table, ways: int, column: str, is_asc: bool) -> Table:
210     if ways <= 1:
211         raise RuntimeError("invalid ways {} for merge sort".format(ways))
212
213     total_chunks = src_table.chunk_manager.get_chunk_cnt()
214     # if there are 0 - 100 chunks, then 1 run is needed
215     # if there are 101 - 10000 chunks, then 2 run is needed
216     # ...
217     passes = math.ceil(math.log(total_chunks, ways))
218     input_table = src_table
219     step = 1
220     generated_runs_cnt = total_chunks
221     for _ in range(passes):
222         cm = input_table.chunk_manager
223         # total_chunks must be updated for each pass, as the size might shrink
224         total_chunks = cm.get_chunk_cnt()
225         generated_runs_cnt = math.ceil(generated_runs_cnt / ways)
226         out_put_table, status = get_table_manager().create_tmp_table(src_table.metadata)
227         if not status.ok():
228             raise RuntimeError("failed to create new tmp table")
229         idx = 0
230         for i in range(generated_runs_cnt):
231             runs = []
232             # load at most ways chunks into memory
233             for j in range(ways):
234 >                 if idx >= total_chunks: ...
236                 chunk, status = cm.load_chunk(idx)
237                 if not status.ok():
238                     raise RuntimeError("failed to load records")
239                 runs.append([chunk, max(min(step, total_chunks - idx) - 1, 0), min(idx + step, total_chunks)])
240                 idx += step
241             idx_per_way, finish_cnt, new_records = [0 for _ in range(len(runs))], 0, []
242             while finish_cnt != len(runs):
243                 # if any way has run out of element, it should be ignored
244                 elements = [(runs[j][0][idx_per_way[j]], j) for j in range(len(runs)) if idx_per_way[j] != len(runs[j][0])]
245
246                 def sort_key(item):
247                     return item[0][column]
248
249                 next_idx = min(elements, key=sort_key)[1] if is_asc else max(elements, key=sort_key)[1]
250                 # notice, here we don't have to flush chunk into memory whenever it's full
251                 new_records.append(runs[next_idx][0][idx_per_way[next_idx]])
252                 idx_per_way[next_idx] += 1
253                 if idx_per_way[next_idx] == len(runs[next_idx][0]):
254                     if runs[next_idx][1] != 0:
255                         chunk, status = cm.load_chunk(runs[next_idx][2] - runs[next_idx][1])
256                         if not status.ok():
257                             raise RuntimeError("failed to load records")
258                         runs[next_idx][0] = chunk
259                         runs[next_idx][1] -= 1
260                         idx_per_way[next_idx] = 0
261                     if idx_per_way[next_idx] == len(runs[next_idx][0]) and runs[next_idx][1] == 0:
262                         finish_cnt += 1
263                 status = out_put_table.insert_bulk(new_records)
264                 if not status.ok():
265                     raise RuntimeError("failed to flush onto disk")
266             step *= ways
267             input_table = out_put_table
268
269     return input_table
270
```

TableManager

```
app > common > table > table_manager.py
```

```
22
```

```

23 # Manager of all tables, avoid data racing
24 class TableManager:
25
26 > def __init__(self, ctx: Context): ...
35
36 > def load_table_names(self) -> (list[str], Status): ...
49
50 > def load_all_metadata(self) -> (dict[str, Metadata], Status): ...
67
68 > def _register_table_in_memory(self, table_name: str, metadata: Metadata) -> (Table | None, Status): ...
91
92 > def _unregister_table_in_memory(self, table_name: str) -> Status: ...
98
99 # will override
100 > def _create_table_on_disk(self, table_name: str, metadata: Metadata) -> Status: ...
111
112 # fault tolerate
113 > def _drop_table_on_disk(self, table_name: str): ...
120
121 > def drop_table(self, table_name: str) -> Status: ...
132
133 > def create_table(self, table_name: str, metadata: Metadata) -> (Table | None, Status): ...
155
156 > def check_consistency_between_metadata_and_table(self) -> bool: ...
172
173 > def create_tmp_table(self, metadata) -> (Table | None, Status): ...
181
182 > def is_tmp_table(self, table_name: str): ...
184
185 > def clear_tmp_tables(self) -> Status: ...
191
192 def start(self) -> Status:
193 > | if self.is_started(): ...
196
197 | status = self.clear_tmp_tables()
198 > | if not status.ok(): ...
201
202 | # check all tables under tables subdir
203 | # any table under tables subdir must have a metadata file under metadata (no dangling table is allowed)
204 > | if not self.check_consistency_between_metadata_and_table(): ...
207
208 | metadata_dict, status = self.load_all_metadata()
209 > | if not status.ok(): ...
212
213 | for name, metadata in metadata_dict.items():
214 | | if name in self.table_map:
215 | | | self.logger.error("duplicated table {} detected", name)
216 | | | return DUPLICATED_TABLE_CREATION_REQUEST
217 | | table_path = os.path.join(self.tables_dir, name)
218 | | if not os.path.exists(table_path):
219 | | | os.makedirs(table_path)
220 | | self.table_map[name] = Table(table_name=name, metadata=metadata, ctx=self.ctx)
221
222 | self.state = TableManagerState.RUNNING
223 | self.logger.info("table manager started successfully")
224 | return OK
225
226 > def is_started(self) -> bool: ...
228
229 > def get_table(self, table_name: str) -> Table | None: ...
231
232
233 table_manager_singleton = None
234
235
236 > def get_table_manager(ctx: Context = None) -> TableManager: ...
241
242
243 > if __name__ == "__main__": ...
278

```

Table

```
0 # iterate through a by chunks
1 > class TableIterator: ...
2
3
4
5 # The abstraction of actual tables storing on the disk
6 # It should support two different data formats: json or csv
7 # Supported operations: projection, filtering, sorting, deduplication, grouping by, joining two tables
8 # It also contains all information of the table including metadata, location of metadata and actual table on the disk
9 class Table:
10
11     def __init__(self, table_name: str, metadata: Metadata, ctx: Context):
12         self.name = table_name
13         self.metadata = metadata
14         self.chunk_cnt = 0
15         self.logger = ctx.get_logger()
16         self.cfg = ctx.get_cfg()
17         self.chunk_manager = ChunkManager(os.path.join(self.cfg.get_tables_dir(), table_name), metadata, ctx)
18         self.chunk_manager.start()
19
20     # drop this table, delete from disk
21 > def drop(self) -> Status: ...
22
23     # table should know how many chunks are on the disks, where to find each chunk, and how to update/delete/insert data into these trunks
24     # record should be a json object
25 > def insert(self, record: dict) -> Status: ...
26
27     def insert_bulk(self, records: list[dict]) -> Status:
28         # self.logger.info("inserting a record {} to table {}".format(record, self.name))
29         status = self.chunk_manager.dump_bulk(records)
30         if not status.ok():
31             self.logger.warn("failed to insert record #{} to table {}".format(len(records), self.name))
32             return status
33         # self.logger.info("record {} is inserted to table {}".format(record, self.name))
34         return OK
35
36     # TODO add lock for all these operations on table
37 > def update(self, selector, new_record: dict) -> Status: ...
38
39     # TODO add lock for all these operations on table
40 > def delete(self, selector) -> Status: ...
41
```

ChunkManager

```
app > common > table > chunk_manager.py
34 class ChunkManager:
35
36     # metadata is only used for SQL
37 > def __init__(self, table_path: str, metadata: Metadata, ctx: Context): ...
38
39 > def start(self) -> Status: ...
40
41 > def get_chunk_cnt(self): ...
42
43 > def get_fist_chunk(self) -> (list[object], Status): ...
44
45 > def get_last_chunk(self) -> (list[object], Status): ...
46
47 > def get_chunk_path(self, chunk_idx) -> str: ...
48
49 > def get_iter(self) -> ChunkIterator: ...
50
51 > def get_remaining_slots(self, occupied_cnt: int) -> int: ...
52
53 > def get_last_chunk_index(self) -> int: ...
54
55     # return a list of objects which can be iterated through
56     # for sql, we have to transfer each row (which is simply a list to an object) to use it with ease
57 > def load_chunk(self, chunk_idx: int) -> (list[dict[str, object]], Status): ...
58
59 > def create_new_chunk(self) -> Status: ...
60
61 > def is_empty_table(self) -> bool: ...
62
63 > def is_chunk_full(self, chunk: list[object]) -> bool: ...
64
```

```

140 # always append new record to the last chunk
141 # if the last chunk is full i.e. len(chunk) == config.chunk_size, then create a new chunk as the last chunk
142 > def dump_one(self, record: dict) -> Status: ...
143
144
145
146 def dump_bulk(self, records: list[dict]) -> Status:
147     if self.is_empty_table():
148         status = self.create_new_chunk()
149         if not status.ok():
150             return INTERNAL
151
152     chunk, status = self.get_last_chunk()
153     if not status.ok():
154         self.logger.error("failed to append new record as unable to load last chunk")
155         return status
156
157     # fill the last chunk
158     remaining = self.get_remaining_slots(len(chunk))
159     insert_cnt = min(remaining, len(records))
160     chunk += records[:insert_cnt]
161     status = self.update_chunk(self.get_last_chunk_index(), chunk)
162     if not status.ok():
163         self.logger.error("failed to append records to the last chunk")
164         return status
165     records = records[insert_cnt:]
166
167     # create new chunks and fill
168     while len(records):
169         status = self.create_new_chunk()
170         if not status.ok():
171             return INTERNAL
172         size = min(self.max_chunk_size, len(records))
173         status = self.update_chunk(self.get_last_chunk_index(), records[:size])
174         if not status.ok():
175             return INTERNAL
176         records = records[size:]
177     return OK
178
179 > def destroy_all_chunks(self) -> Status: ...
180
181
182
183 > def update_chunk(self, chunk_idx: int, chunk: list[dict[str, object]]) -> Status: ...
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205

```

Learning Outcomes

1. Designing a Unique Query Language

- **Challenge:** Creating a query language distinct from existing ones, which needed to be expressive, clear, and easy to parse, presented a unique set of difficulties.
- **Solution:** We chose JSON as the format for our query language, striking a balance between expressiveness and simplicity, ensuring ease of parsing and clarity.

2. Handling Multiple Source Tables and Subqueries

- **Challenge:** Managing references to specific fields across multiple tables and nested subqueries posed a significant challenge, particularly in maintaining clarity and accuracy for users and in the backend processing.
- **Solution:** We devised a set of naming conventions that could be seamlessly integrated into our query language and the query engine. This innovation greatly simplified referencing fields in complex query scenarios.

3. Integration of SQL and NoSQL

- **Challenge:** A major challenge was integrating SQL and NoSQL databases to handle both CSV (structured) and JSON (possibly nested) data in a unified manner.
- **Solution:** We developed a system where components like the TableManipulator and ChunkManager are versatile enough to support both data formats. This approach minimized code duplication and streamlined our development process.

4. Memory Management and Chunk-Based Data Storage

- **Challenge:** Given the constraint of limited memory, efficiently managing large datasets by storing and processing them in chunks was a significant technical challenge.
- **Solution:** We carefully designed the ChunkManager, Table, and TableManipulator classes, simplifying the implementation of complex operations like merge sorting on chunks. To ensure the reliability of each API, extensive unit testing was conducted. This comprehensive testing approach reduced the likelihood of errors when integrating different system components.

Individual Contribution

Tasks	Assignee
Design Query Language	Weihao Zhao
Design Data Model & Import Data	Zhenyu Xiong
Design Testcases	Weihao Zhao
Implement SQL Database	Zhenyu Xiong
Implement NoSQL Database	Weihao Zhao
Implement CLI	Weihao Zhao
Draft Project Proposal	Weihao Zhao
Draft Midterm Report	Zhenyu Xiong
Draft Final Report	Weihao Zhao
Make Demo Video	Zhenyu Xiong

Conclusion

This project has been an exemplary journey in database management system design, especially in the context of an educational setting. By successfully integrating SQL and NoSQL databases and developing a custom JSON-based query language, our team not only met but exceeded the project's requirements. This accomplishment underscores our ability to address complex challenges, such as handling various data formats and implementing efficient memory management in a constrained environment.

The collaboration throughout this project was solid and effective, with each team member contributing significantly to the project's success. We navigated the complexities of database systems, developing solutions for managing multiple source tables, subqueries, and implementing chunk-based data storage. These experiences have greatly enhanced our understanding of databases, showcasing our adaptability and problem-solving skills in the face of technical challenges.

In essence, this project has not only fulfilled its academic objectives but has also provided us with invaluable insights into the practical aspects of database system design and development, laying a strong foundation for our future endeavors in this field.

Future Scope

1. **Concurrency Control**

Introducing locks for operations on tables is an essential aspect of concurrency control in database systems, ensuring data integrity and consistency during concurrent access.

2. **Implementation of Mock Server in CLI:**

A mock server setup in `cli/mock_server.py` is planned but not yet implemented. Completing this would enhance testing and simulation capabilities for the CLI.

3. **Type Conversion Handling:**

Handling potential type conversions (e.g., int to float) during data manipulations is an area for development. This addition would increase the robustness and accuracy of data processing.

4. **Completion of Test Cases:**

Developing comprehensive test cases is crucial for ensuring the system's reliability and performance under various scenarios.