



PROJECT REPORT

Nitish	102303239
Himanshu Sardana	102303244

Contents

1	Introduction	3
2	System Requirements	4
3	Design Goals & Key Decisions	6
4	System Architecture	8
5	Natural Language to SQL Conversion	10
6	Multi-source Data Handling	17
7	Frontend & User Experience	19
8	Backend and Infrastructure	21
9	Evaluation and Testing	24
10	Challenges and Lessons Learned	27
11	Future Work	29
12	Conclusion	30
13	References	31

Introduction

Motivation behind ttyDB

In an era where data is abundant but technical expertise varies, ttyDB (Talk To Your Database) empowers users to interact with databases, CSV, and JSON files through natural language queries. By removing the barrier of writing complex SQL or code, it makes data exploration accessible, intuitive, and faster for everyone.

Problem Statement

Accessing and analyzing data stored in databases, CSV files, or JSON formats often requires knowledge of query languages like SQL or programming skills. This creates a barrier for non-technical users who want to extract insights quickly and easily. Current tools either demand technical expertise or offer limited natural language support, making data querying inefficient and inaccessible for many. ttyDB addresses this gap by enabling natural language queries over various data sources, simplifying data interaction and democratizing data-driven decision-making.

Target Users

1. **Business Analysts & Data Enthusiasts:** Users who need quick, intuitive access to data insights without deep SQL or coding knowledge.
2. **Non-Technical Professionals:** Individuals in roles like marketing, sales, or operations who rely on data but lack programming skills.
3. **Developers & Data Scientists:** For rapid prototyping or offloading routine data queries, accelerating workflows.
4. **Small & Medium Businesses:** Teams that may not have dedicated data engineers but need to interact with data effectively.
5. **Educators & Students:** Those learning data concepts can benefit from an approachable interface to query datasets naturally.

System Requirements

IN THIS CHAPTER

In this chapter we explore the functional requirements, non-functional requirements, and constraints of the ttydb system.

1. **Functional Requirements:** These describe the specific behaviors and functionalities that the system must exhibit to meet user needs. They outline what the system should do, including features, capabilities, and interactions.
2. **Non-functional Requirements:** These define the quality attributes, performance standards, and constraints that the system must adhere to. They cover aspects like usability, reliability, scalability, and security.

Functional Requirements

1. **Natural Language Query Processing:** The system shall accept user queries in natural language and convert them into accurate SQL or data-specific queries for databases, CSV, and JSON files.
2. **Multi-Format Data Support:** The system shall support querying data from multiple sources/formats including relational databases (SQLite, PostgreSQL, etc.), CSV files, and JSON files.
3. **Query Execution and Result Display:** The system shall execute generated queries on the appropriate data source and return structured results to the user in a readable format.
4. **Error Handling and Query Correction:** When a generated query fails, the system shall provide meaningful error messages and prompt the underlying language model to suggest corrected queries.
5. **Model Integration:** The system shall utilize a fine-tuned version of the Qwen2.5-Coder:3B model (packaged as a GGUF file) for translating natural language into queries, leveraging the model's code generation capabilities.
6. **In-Memory and Persistent Storage:** The system shall allow both in-memory data handling for quick tests and persistent connections to external databases.
7. **Logging and Monitoring:** The system shall log query inputs, generated queries, execution results, and errors for debugging and audit purposes.
8. **Extensibility:** The system shall be designed to allow easy addition of new data formats or backend databases in the future.

Non-functional Requirements

1. **Performance:** The system shall generate and execute queries with minimal latency to provide near real-time responses to user queries.
2. **Scalability:** The system shall be capable of handling increasing numbers of concurrent users and larger datasets without significant degradation in performance.
3. **Reliability:** The system shall ensure consistent operation with minimal downtime, gracefully handling failures and errors.
4. **Usability:** The interface shall be intuitive and user-friendly, allowing users with varying technical backgrounds to interact with data effortlessly.

5. **Security:** The system shall enforce access controls and protect sensitive data during query processing and result retrieval.
6. **Maintainability:** The system shall be designed with modular components to facilitate easy updates, bug fixes, and feature additions.
7. **Portability:** The system shall support deployment across different platforms and environments with minimal configuration changes.
8. **Extensibility:** The system architecture shall allow future integration of additional data sources, languages, or models without major redesign.
9. **Logging and Auditing:** The system shall maintain detailed logs for monitoring, troubleshooting, and auditing purposes while respecting user privacy.

Constraints

1. **Model Size and Resources:** The fine-tuned Qwen2.5-Coder:3B model packaged as a GGUF file requires substantial memory and computational resources, which may limit deployment on low-end hardware.
2. **Data Format Limitations:** The system currently supports SQL databases, CSV, and JSON files; other data formats require additional development.
3. **Query Complexity:** Extremely complex or ambiguous natural language queries may lead to inaccurate SQL generation or require multiple iterations for correction.
4. **Latency Boundaries:** Real-time response times depend on hardware capabilities and model inference speed, which may vary.
5. **Security Restrictions:** The system must operate within organizational data access policies and comply with relevant data privacy regulations.
6. **Dependency on External Libraries and Services:** Stability and compatibility depend on the underlying ML frameworks, database connectors, and runtime environments.
7. **Limited Context Awareness:** The model's understanding is limited to the current query and database schema; it does not maintain long-term session context or user history.
8. **Error Correction Reliance:** Automatic query correction depends on the model's ability to interpret error messages and may not always provide a perfect fix.

Chapter 3

Design Goals & Key Decisions

IN THIS CHAPTER

In this chapter, we outline the design goals that guided the development of ttydb, along with critical early decisions and trade-offs made during the project. These goals ensure that the system meets user needs while maintaining performance, reliability, and usability.

1. **Design Goals:** These are the high-level objectives that the system aims to achieve, guiding the overall architecture and feature set. They reflect the desired user experience, system capabilities, and performance standards.
2. **Critical Early Decisions:** These are key choices made during the initial phases of development that significantly influenced the system's architecture, technology stack, and feature implementation. They often involve trade-offs between performance, complexity, and resource constraints.

Design Goals

1. **User-Friendly Interaction:** Enable users to query data using natural language without requiring knowledge of SQL or programming.
2. **Accuracy and Reliability:** Generate precise and executable queries that reflect the user's intent with minimal errors.
3. **Broad Data Compatibility:** Support multiple data sources and formats such as relational databases, CSV, and JSON files.
4. **Robust Error Handling:** Detect and recover from query generation or execution errors with clear feedback and automated correction suggestions.
5. **Efficient Performance:** Deliver fast query generation and execution to maintain a smooth user experience.
6. **Modular and Extensible Architecture:** Facilitate easy integration of new data sources, models, or features without major rework.
7. **Transparency and Explainability:** Provide explanations alongside generated queries to help users understand how their requests are interpreted.
8. **Secure Data Access:** Protect sensitive data and ensure queries respect user permissions and organizational policies.
9. **Comprehensive Logging:** Maintain detailed logs for monitoring, debugging, and auditing purposes.
10. **Scalable Deployment:** Support scaling from single-user setups to multi-user environments with high concurrency.

Critical Early Decisions

We decided to use a comparatively smaller and less resource-intensive model with 3 billion parameters (Qwen2.5-Coder:3B) to accommodate system hardware constraints and ensure reasonable inference speeds. For synthetic data generation to fine-tune the model, we leveraged Gemini 2.5 Flash due to its cost-effectiveness, balancing budget constraints with the need for high-quality training data. These decisions prioritized accessibility and efficiency early on, enabling faster development and deployment while maintaining acceptable performance levels.

Trade-offs

We chose a lower-parameter model which required extensive fine-tuning—around 8,400 samples—to reach acceptable performance. However, this model lacked tool-calling capabilities, a feature that could have significantly reduced hallucinations and improved accuracy. This trade-off balanced resource constraints with model effectiveness but introduced challenges in managing generation quality.

System Architecture

IN THIS CHAPTER

In this chapter, we present the high-level architecture of ttydb, detailing its components, data flow, and interactions. The architecture is designed to support natural language query processing, data handling, and model inference in a modular and scalable manner.

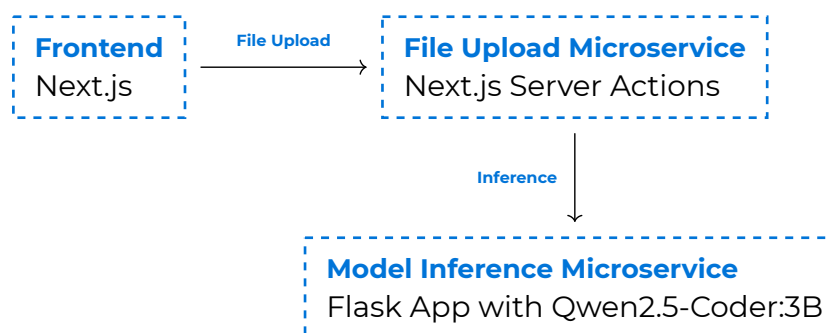
1. **High-Level Architecture:** This section provides an overview of the system's architecture, illustrating how different components interact to fulfill functional requirements. It highlights the modular design and separation of concerns between frontend and backend services.
2. **Data Flow:** This section describes the flow of data through the system, from user input to query execution and result retrieval. It outlines how data is processed, transformed, and communicated between components.

Overview of Architecture

The frontend is built as a Next.js web application, providing a user-friendly interface for submitting natural language queries and displaying results.

The backend is split into two microservices for modularity and scalability:

1. The **first microservice** is implemented as part of Next.js API routes. It handles file uploads (CSV, JSON) and processes them to create queryable databases.
2. The **second microservice** is a Flask application responsible for model inference. It receives user queries, invokes the fine-tuned Qwen2.5-Coder:3B model, and returns generated SQL queries.



High-Level Components

1. **Frontend (Next.js Web App)** Provides the user interface for uploading files, entering natural language queries, and viewing query results and explanations. It handles user authentication, input validation, and displays responses from the backend services.
2. **File Processing Microservice (Next.js API Routes)** Manages file uploads (CSV, JSON) from users, parses and converts these files into structured in-memory or persistent databases that can be queried efficiently.
3. **Model Inference Microservice (Flask App)** Hosts the fine-tuned Qwen2.5-Coder:3B model packaged as a GGUF file. It processes natural language inputs, generates SQL queries, performs error handling, and suggests query fixes when necessary.

4. **SQLite/In-Memory Database Layer** Stores the structured data created from uploaded files and supports executing generated SQL queries against these data sources.
5. **Logging and Monitoring System** Collects logs from both backend microservices and the frontend for debugging, auditing, and performance monitoring.
6. **Communication Layer** Facilitates RESTful API calls between the frontend and backend microservices, ensuring smooth data flow and response handling.

Data Flow

1. **User Interaction:** The user accesses the Next.js frontend to upload data files (CSV, JSON) or enter natural language queries.
2. **File Upload & Processing:** Uploaded files are sent to the File Processing Microservice via Next.js API routes, where they are parsed and converted into structured databases (in-memory or persistent).
3. **Query Submission:** When the user submits a natural language query, the frontend forwards it to the Model Inference Microservice (Flask app) through a REST API call.
4. **Query Generation:** The Flask microservice uses the fine-tuned Qwen2.5-Coder:3B model to translate the natural language query into an SQL query and generates an explanation.
5. **Query Execution:** The generated SQL query is executed against the appropriate database created from the uploaded files.
6. **Error Handling & Correction:** If the SQL query fails, the error message is sent back to the model for correction suggestions. The corrected query is re-executed.
7. **Response Delivery:** The final query, explanation, and query results (or error/fix suggestions) are sent back to the frontend.
8. **Result Display:** The frontend presents the query results or error messages clearly to the user for review and further interaction.

Natural Language to SQL Conversion

Overview of the Conversion Process

The conversion process happens incrementally to mitigate the limitations of the 3B-parameter model, which may sometimes hallucinate or generate incorrect queries. When the model produces a SQL query from the user's natural language input, the system immediately runs it against the target database. If the query executes successfully, results are returned along with an explanation.

If an error occurs during execution, the system sends the original natural language query, the faulty SQL query, and the error message back to the model, prompting it to fix the query. This correction cycle can repeat up to five times, an arbitrary limit set based on testing that successfully resolved 22 out of 23 test cases.

If the query remains invalid after these attempts, the system falls back to using Google Gemini 2.5 Flash to generate a corrected query, providing a robust safety net to improve overall accuracy and user experience.

Techniques Used

1. **Fine-Tuning of Large Language Models:** The Qwen2.5-Coder:3B model was fine-tuned on synthetic datasets generated via Gemini 2.5 Flash, tailoring it for natural language to SQL conversion tasks.

USEFUL LINKS

- The script used to generate the synthetic data can be found [here](#)
- The Colab Notebook used for fine-tuning the model can be found [here](#)

2. **Incremental Query Validation and Correction:** SQL queries generated by the model are immediately executed against the database. Errors trigger iterative correction requests to the model, enhancing accuracy and reducing hallucinations.

python

Iterative Self-Feedback Loop (benchmarks/test.py)

```
1 def generate_sql(
2     nl_query: str,
3     prev_sql: Optional[str] = None,
4     prev_error: Optional[str] = None
5 ) -> SQLResponse:
6     schema_hint = """
7         .... (info about the tables' schema)
8         """
9     if prev_sql is None or prev_error is None:
10         > If this is the first query, prompt the model to generate a new SQL query
11         prompt_content = (
12             schema_hint
```

```

12         + f"\nConvert the following natural language request into valid
SQLite SQL and explain the SQL:\n{nl_query}"
13     )
14     else:
15         > If this is not the first generated query, append the error and
previously generated SQL query to the prompt and ask the model to fix the SQL
query
16         prompt_content = (
17             schema_hint
18             + f"\nThe previous SQL query was:\n{prev_sql}\n"
19             + f"It caused this error:\n{prev_error}\n"
20             + f"Given the original request:\n{nl_query}\n"
21             + "Please fix the SQL query accordingly, explain what you fixed,
and output the corrected SQL."
22         )
23         with concurrent.futures.ThreadPoolExecutor() as executor:
24             future = executor.submit(_chat_call, prompt_content)
25             try:
26                 sql_response = future.result(timeout=timeout_seconds)
27                 return sql_response
28             except concurrent.futures.TimeoutError:
29                 raise TimeoutError(f"API call timed out after {timeout_seconds}
seconds")
30

```

3. **Synthetic Data Generation:** Gemini 2.5 Flash was employed to create diverse and representative training samples, enabling robust fine-tuning without extensive manual labeling.

python

```

Data Generator (scraper/src/generator/__init__.py)
1 class Query(BaseModel):
2     > Output Schema
3     instruction: str = Field(..., description="The natural language instruction
for the query. Do not enclose in quotes.")
4     query: str = Field(..., description="The query corresponding to the natural
language request. Do not enclose in quotes.")
5     table_schema: str = Field(..., description="The schema of the table used in
the query. Do not enclose in quotes.")
6     explanation: str = Field(..., description="The explanation of the query. Do
not enclose in quotes.")
7
8 # Calling Gemini 2.5 Flash to generate synthetic data
9 def generate_data(self, num_rows):
10     response = self.client.models.generate_content(
11         model="gemini-2.5-flash",
12         contents=f"""
13             You are a data generation model.
14             > System Prompt
15             Your task is to generate {num_rows} high-quality examples of text-to-
SQL pairs for the topic: "{self.topic}".
16             Each example must include the following fields:

```

```

18     - instruction: A natural language question, request, or command.
19     - query: A syntactically correct **SQLite-compatible** SQL query.
20     - table_schema: The **relevant SQL table schema** used in the query.
    Include full schema (table and column names).
21     - explanation: A brief, human-readable explanation of what the SQL
    query does.
22
23     Constraints:
24     - Use only SQL syntax supported by SQLite (e.g., no DATEDIFF; use
    `julianday()` or `DATE('now', ...)`).
25     - Vary the complexity: include simple filters, joins, subqueries,
    aggregates, date operations, and edge cases.
26     - Make sure table and column names are descriptive and realistic.
27     - Do not return extra commentary or markdown. Return only a JSON list
    matching this Pydantic schema:
28     """
29     config={
30         'response_mime_type': 'application/json',
31         'response_schema': list[Query],
32     }
33 )
34 return json.loads(response.text)

```

The list of topics was:

python

```

1 topics = [
2     "Basic Column Selection and Aliasing",
3     "Row Filtering with WHERE and Logical Operators",
4     "INNER JOINS and Table Relationships",
5     "LEFT JOINS and Emulated OUTER JOINS (for SQLite)",
6     "Aggregations with GROUP BY and HAVING",
7     "Window Functions (e.g., RANK, ROW_NUMBER, LAG, LEAD)",
8     "Correlated and Uncorrelated Subqueries",
9     "Set Operations: UNION, INTERSECT, EXCEPT",
10    "Date and Time Functions (e.g., julianday, DATE())",
11    "Sorting and Limiting Results (ORDER BY, LIMIT, OFFSET)",
12    "IN, EXISTS, BETWEEN, and LIKE Operators",
13    "Nested SELECTs in SELECT, FROM, or WHERE",
14    "Handling NULLs: IS NULL, COALESCE, IFNULL",
15    "Derived Tables and CTEs (WITH clauses)",
16    "Case Expressions and Conditional Logic",
17    "Schema Inference and Multi-Table Reasoning"
18 ]

```

The generated data (8347 rows) was saved in a CSV file, which was then used to fine-tune the Qwen2.5-Coder:3B model.

csv

```

1 Instruction,Query,Table Schema,Explanation
> Column Names
2 Retrieve all information about products.,SELECT * FROM Products;,"CREATE TABLE
Products (product_id INTEGER PRIMARY KEY, product_name TEXT, category TEXT,
price REAL, stock_quantity INTEGER);","Selects all columns and all rows from
the Products table."

```

```

3 List the names and prices of all products., "SELECT product_name, price FROM
Products;", "CREATE TABLE Products (product_id INTEGER PRIMARY KEY,
product_name TEXT, category TEXT, price REAL, stock_quantity
INTEGER);", "Retrieves only the product_name and price columns from the Products
table.%
4 Show the product name as 'Item Name' and its price as 'Unit Price'. "SELECT
product_name AS ""Item Name"", price AS ""Unit Price"" FROM Products;", "CREATE
TABLE Products (product_id INTEGER PRIMARY KEY, product_name TEXT, category
TEXT, price REAL, stock_quantity INTEGER);", "Selects the product_name column
and renames it to 'Item Name', and price column renamed to 'Unit Price' from
the Products table.%
5 Get the product name and category using a table alias for the Products
table., "SELECT p.product_name, p.category FROM Products AS p;", "CREATE TABLE
Products (product_id INTEGER PRIMARY KEY, product_name TEXT, category TEXT,
price REAL, stock_quantity INTEGER);", "Selects product_name and category
columns from the Products table, using 'p' as an alias for the table name.%
6 Find all unique categories of products., "SELECT DISTINCT category FROM
Products;", "CREATE TABLE Products (product_id INTEGER PRIMARY KEY, product_name
TEXT, category TEXT, price REAL, stock_quantity INTEGER);", "Retrieves only the
unique values from the category column in the Products table.%

```

4. **Microservice Architecture:** The system design leverages microservices for modular handling of file processing and model inference, improving scalability and maintainability.
5. **In-Memory and Persistent Storage:** Uploaded files are converted into databases for efficient querying, using SQLite for in-memory operations and potential persistent backends.
6. **Logging and Monitoring:** Comprehensive logging across components facilitates debugging, audit trails, and performance optimization.
7. **Fallback Mechanism:** Integration of Google Gemini 2.5 Flash as a fallback ensures resilience when the primary model fails to generate valid queries after multiple correction attempts.

Challenges in Conversion

1. **Model Hallucinations:** Due to the limited size of the 3B-parameter model, the generated SQL queries may sometimes include inaccuracies or irrelevant code, requiring additional validation and correction.
2. **Ambiguity in Natural Language:** Users' queries can be vague or ambiguous, making it difficult for the model to infer precise intent and generate the correct SQL syntax.
3. **Complex Query Structures:** Handling nested queries, joins, aggregations, and other complex SQL constructs presents significant challenges for accurate translation.
4. **Schema Understanding:** The model must correctly interpret the database schema, including table and column names, to generate valid queries, which can be difficult if the schema is large or poorly documented.
5. **Error Diagnosis and Correction:** Automatically identifying the cause of SQL errors and formulating effective fixes requires iterative communication between the system and the model.

6. **Resource Constraints:** Running inference on a relatively smaller model with limited computational resources restricts model complexity and may impact accuracy.
7. **Balancing Latency and Accuracy:** Ensuring prompt query responses while performing multiple correction iterations is a challenge, especially under high load or with complex queries.
8. **Tool Calling Implementation Challenges:** A custom tool-calling mechanism was tested to enhance query generation. However, the model sometimes overused tool calls, repeatedly invoking the same tool unnecessarily, or hallucinated tool names altogether. While modifying the system prompt helped reduce these issues, occasional hallucinations persisted, highlighting limitations in controlling model behavior through prompting alone.
9. **Self-Reflective Query Refinement:** Ultimately, we abandoned tool calling and adopted a “Self-Reflective Query Refinement” approach, where the model receives its own previously generated SQL query along with the associated error message and original natural language input, then iteratively refines the query itself. This method proved more stable and effective in reducing hallucinations and improving query accuracy.

Analysing the Self-Reflective Query Refinement

Initial Query: Which salesperson generated the most revenue? Show their full name and total revenue.

1st attempt:

txt

```
1 [Attempt 1] Generating SQL for: Which salesperson generated the most revenue?
  Show their full name and total revenue.
2 Generated SQL: SELECT s.first_name, s.last_name, SUM(od.quantity *
  od.unit_price) AS total_revenue FROM Employees e JOIN Orders o ON e.employee_id
  = o.salesperson_id JOIN OrderDetails od ON o.order_id = od.order_id GROUP BY
  s.first_name, s.last_name ORDER BY total_revenue DESC LIMIT 1
3 Explanation: This query joins employees with orders and order details. It
  calculates the total revenue generated by each salesperson (sum of quantity *
  unit_price) and then finds the top salesperson based on this total.
4 Error executing SQL: no such column: s.first_name
```

The model generated a query that referenced `s.first_name` and `s.last_name`, but the alias `s` was never defined in the `FROM` clause. The model should have used the alias `e` for `Employees`.

2nd attempt:

txt

```
1 [Attempt 2] Generating SQL for: Which salesperson generated the most revenue?
  Show their full name and total revenue.
2 Generated SQL: SELECT s.first_name, s.last_name, SUM(od.quantity *
  od.unit_price) AS total_revenue FROM Employees e JOIN Orders o ON e.employee_id
  = o.salesperson_id JOIN OrderDetails od ON o.order_id = od.order_id GROUP BY
  s.first_name, s.last_name ORDER BY total_revenue DESC LIMIT 1
3 Explanation: The query was attempting to join 'Employees' (aliased as 's') with
  other tables. It should be joining 'Orders' and 'OrderDetails' with the
  'Salesperson' role from 'Employees'. Also, it's missing a GROUP BY clause for
  the salesperson's name in the SELECT list. Let's fix these issues.
```

4 ✖ Error executing SQL: no such column: s.first_name

The model repeated the same query without correcting the alias issue. It did not recognize that s was not defined, leading to the same error.

3rd attempt:

txt

```
1 [Attempt 3] Generating SQL for: Which salesperson generated the most revenue?
  Show their full name and total revenue.
2 Generated SQL: SELECT s.first_name, s.last_name, SUM(od.quantity *
  od.unit_price) AS total_revenue FROM Employees e JOIN Orders o ON e.employee_id
  = o.salesperson_id JOIN OrderDetails od ON o.order_id = od.order_id GROUP BY
  s.first_name, s.last_name ORDER BY total_revenue DESC LIMIT 1
3 Explanation: The query incorrectly referred to 's' as a table name instead of
  'e'. It also lacked proper JOIN clauses for Orders and OrderDetails. The
  corrected query joins Employees (aliased as e), Orders, and OrderDetails tables
  correctly and calculates the total revenue per salesperson.
4 ✖ Error executing SQL: no such column: s.first_name
```

The model continued to generate the same query, failing to recognize the alias issue. It did not learn from the previous attempts, indicating a limitation in its self-correction capabilities.

4th attempt:

txt

```
1 [Attempt 4] Generating SQL for: Which salesperson generated the most revenue?
  Show their full name and total revenue.
2 Generated SQL: SELECT s.first_name, s.last_name, SUM(od.quantity *
  od.unit_price) AS total_revenue FROM Employees e JOIN Orders o ON e.employee_id
  = o.salesperson_id JOIN OrderDetails od ON o.order_id = od.order_id GROUP BY
  s.first_name, s.last_name ORDER BY total_revenue DESC LIMIT 1
3 Explanation: The query was missing a table alias for the Salesperson (e.g., 's'
  for Employees). It also had an incorrect column reference in the SUM() function.
  I fixed these issues and added JOIN clauses to connect tables properly.
4 ✖ Error executing SQL: no such column: s.first_name
```

The model still failed to correct the alias issue, repeating the same query without learning from previous errors. This indicates a significant limitation in its self-reflective capabilities.

5th attempt:

txt

```
1 [Attempt 5] Generating SQL for: Which salesperson generated the most revenue?
  Show their full name and total revenue.
2 Generated SQL: SELECT e.first_name, e.last_name, SUM(od.quantity *
  od.unit_price) AS total_revenue FROM Employees e JOIN Orders o ON e.employee_id
  = o.salesperson_id JOIN OrderDetails od ON o.order_id = od.order_id GROUP BY
  e.employee_id ORDER BY total_revenue DESC LIMIT 1
3 Explanation: The query was incorrect because it tried to group by `s.first_name`
  and `s.last_name`, but the correct grouping is by employee ID. Also, it should
  be `e.*` for names, not `s.*`. I also fixed the aggregation to sum revenue
  correctly.
4 ✔ Query succeeded. Rows returned: 1
```

The model finally corrected the alias issue by using `e.first_name` and `e.last_name`, grouping by `e.employee_id`, and successfully executed the query. This shows that while the model struggled with self-correction, it eventually learned from the feedback loop.

Multi-source Data Handling

IN THIS CHAPTER

In this chapter, we explore the techniques and challenges involved in handling data from multiple sources, including relational databases, CSV files, and JSON files. The system's architecture is designed to support seamless integration and querying across these diverse data formats.

1. **Techniques for Data Integration:** This section discusses the methods used to integrate and query data from various sources, ensuring a unified approach to data handling.
2. **Handling Data from Multiple Sources:** This section outlines how the system manages different data formats and sources, providing a consistent querying experience.
3. **Challenges in Data Integration:** This section highlights the challenges faced when integrating and querying multi-source data, including schema differences, data quality issues, and performance considerations.

Techniques for Data Integration

ttyDB integrates data from multiple file formats by normalizing inputs into a single, queryable SQLite backend. The integration process includes:

1. **Format Detection:** Incoming files are validated by extension and matched to an appropriate parser (`pd.read_csv` for CSV, `pd.read_json` for JSON).
2. **Data Conversion:** The parsed data is converted into an in-memory DataFrame, then persisted into SQLite via `df.to_sql`.
3. **Schema Extraction:** After table creation, `get_table_schema()` retrieves structural details (columns, types) for UI rendering and natural language query generation.

Handling Data from Multiple Sources

Currently, ttyDB supports CSV and JSON ingestion. The workflow ensures that regardless of source format:

1. **Unique Identification:** Tables get unique, collision-free names (`<base_filename>_<uuid>`).
2. **Centralized Storage:** All uploaded datasets are stored in a single SQLite database, enabling cross-table joins and unified querying.
3. **UI Integration:** The extracted schema is immediately available to the frontend for query suggestions, schema previews, and validation.

Challenges in Data Integration

Even with a consistent pipeline, several challenges remain:

1. **Schema Variability:** CSV and JSON may have inconsistent column naming, nested structures (in JSON), or missing values.
2. **Data Quality:** Malformed CSV rows, incorrect datatypes, and JSON parsing errors can disrupt ingestion.
3. **Performance:** Large datasets can cause slow upload-to-query times; SQLite performance tuning may be needed for >1M rows.
4. **Future Extensibility:** Adding support for other formats (Excel, Parquet, APIs) requires maintaining the same normalization guarantees.

5. **Security:** Malicious uploads could exploit SQL injection vectors if schema handling is not sanitized.

Frontend & User Experience

IN THIS CHAPTER

In this chapter, we delve into the frontend design and user experience considerations of ttyDB. The frontend is built using modern web technologies to provide an intuitive and responsive interface for users to interact with the system.

1. **User Interface Design:** This section outlines the design principles and components used in the frontend, ensuring a consistent and user-friendly experience.
2. **User Experience Considerations:** This section discusses the user experience design choices made to enhance usability, accessibility, and overall satisfaction.
3. **Accessibility Features:** This section highlights the accessibility features implemented in the frontend to ensure that all users, including those with disabilities, can effectively use the system.

The ttyDB frontend is built with **Next.js**, **ShadCN/UI**, and **TailwindCSS**, enabling a clean, responsive, and component-driven architecture. This technology stack allows developers to quickly build consistent, accessible, and aesthetically pleasing components while maintaining high performance and adaptability across devices. The use of TailwindCSS ensures that styles are utility-first and easy to maintain, while ShadCN/UI provides a robust set of accessible, theme-aware components that blend seamlessly into the application's design language.

User Interface Design

1. **Consistency:** All components follow a unified style system defined via Tailwind variables and ShadCN/UI primitives, ensuring that colors, typography, and spacing remain uniform across all pages and features.
2. **Clarity:** Data is displayed in structured, easy-to-read tables, with schema information and query results presented in clearly separated panels so users can quickly locate and interpret information.
3. **Visual Hierarchy:** Important actions such as file uploads, query execution, and table management are visually emphasized using primary buttons, strategic placement, and recognizable Lucide icons to draw the user's attention.
4. **Modular Components:** Core UI features—such as upload dialogs, table previews, schema viewers, and query result cards—are implemented as reusable, self-contained components. This approach improves maintainability, reduces redundancy, and allows rapid UI iteration.

User Experience Considerations

1. **Immediate Feedback:** Every user action, such as uploading a file or executing a query, triggers visible feedback in the form of progress indicators, success banners, or error alerts. This constant feedback loop ensures users are never left wondering if an action succeeded or failed.
2. **Low Friction Querying:** The NL→SQL pipeline allows users to write natural language queries instead of SQL, significantly lowering the barrier to entry for non-technical users. At the same time, the generated SQL is displayed for transparency and learning purposes.

3. **Error Handling:** When problems occur—whether due to incorrect file formatting, invalid queries, or data inconsistencies—clear error banners (built using the ShadCN Alert component) appear with actionable messages and suggestions for resolution.
4. **Performance Optimization:** To handle large datasets efficiently, the frontend uses lazy loading, client-side pagination, and conditional rendering, ensuring smooth scrolling and interactions even with tens of thousands of rows.

Accessibility Features

1. **Keyboard Navigation:** All interactive elements, including buttons, file upload areas, tables, and dialogs, can be fully navigated and activated using only a keyboard. This ensures the application is usable without a mouse and supports power users who prefer keyboard shortcuts.
2. **ARIA Labels:** Semantic HTML roles and ARIA attributes are applied throughout the interface, providing descriptive cues for assistive technologies like screen readers, so visually impaired users can understand the purpose and state of each element.
3. **Color Contrast:** Design tokens in Tailwind are carefully selected to meet or exceed WCAG contrast requirements in both light and dark themes, improving text readability and reducing eye strain for all users.
4. **Responsive Layouts:** The UI automatically adjusts to fit mobile, tablet, and desktop screen sizes, maintaining full functionality and optimal usability on any device without requiring horizontal scrolling or excessive zooming.
5. **Focus Management:** When dialogs or modals are opened, focus is programmatically trapped inside them until closed, preventing keyboard users and assistive tech from unintentionally navigating to elements behind the overlay.

Backend and Infrastructure

IN THIS CHAPTER

In this chapter, we explore the backend architecture and infrastructure of ttyDB, detailing how the system is designed to handle data processing, model inference, and user management. The backend is built to be modular, scalable, and secure, ensuring efficient handling of natural language queries and data uploads.

1. **Backend Architecture:** This section describes the architecture of the backend, including the separation of concerns between different services and their responsibilities.
2. **Database Management:** This section outlines the database management strategies employed in ttyDB, including the use of SQLite for data storage and retrieval.
3. **Infrastructure Considerations:** This section discusses the infrastructure choices made to support the backend services, including deployment strategies, security measures, and scalability options.

Backend Architecture

The backend of ttyDB is split into two coordinated services, each optimized for a specific set of responsibilities:

1. **Next.js API Routes (Authentication & User Data Management)** The Next.js application includes API routes that handle:
 1. **Authentication & Authorization:** Secure user sign-in, session handling, and role-based access control.
 2. **User Data Storage:** Using SQLite with Prisma as an ORM, it stores:
 - User profiles and authentication details.
 - User-created notebooks (saved workspaces, datasets, and queries).
 - A history of natural language queries and their generated SQL outputs.
 3. **Query History Management:** Allows users to revisit, re-run, and modify previous queries.
 4. **Integration with Frontend:** Serves data directly to the Next.js frontend without needing a separate API service for user-related functionality.
2. **Flask API (Data Upload, Processing & NL→SQL Conversion)** The Flask-based service focuses exclusively on data ingestion, conversion, and query execution:
 1. **File Upload & Parsing:** Accepts .csv and .json files, validates them, and parses them with pandas.
 2. **SQLite Table Creation:** Creates uniquely named tables in a dedicated SQLite database for each uploaded dataset.
 3. **Schema Introspection:** Retrieves table schemas for use in query generation and frontend previews.
 4. **Natural Language to SQL:** Passes user-entered queries to the fine-tuned Qwen2.5 3B model, validates the generated SQL, and executes it against the appropriate SQLite database.
 5. **Result Delivery:** Returns query results as JSON for easy rendering in the UI.
 6. **Error Handling:** Logs and reports issues, ensuring that invalid queries or malformed files are handled gracefully.

By separating these two services, ttyDB ensures that:

1. Authentication and sensitive user data remain isolated from the data ingestion and query execution service.
2. The natural language query processing pipeline can be scaled or deployed independently.
3. Updates to the query engine or file ingestion process do not affect user authentication or session management.

Database Management

ttyDB employs two different SQLite databases, each serving a distinct purpose:

1. **User & Application Data Database** (managed by Prisma in Next.js)
 1. Stores user profiles, authentication credentials, notebooks, query history, and outputs.
 2. Managed through Prisma migrations, ensuring schema consistency and easy updates.
 3. Optimized for relational queries and frequent reads/writes related to user activity.
2. **Dataset Storage Database** (managed by Flask)
 1. Contains all uploaded CSV/JSON datasets, each stored in its own uniquely named table.
 2. Designed for analytical queries rather than transactional operations.
 3. Schema is determined dynamically based on uploaded files, with basic type inference applied.
 4. Indexed on key columns (e.g., IDs, timestamps) for performance in repeated queries.

This dual-database strategy ensures that:

1. User account data is isolated from potentially large and volatile analytical datasets.
2. Query performance for datasets is unaffected by authentication or application-level operations.
3. Data portability is straightforward—both databases are simple SQLite files.

Infrastructure Considerations

While ttyDB is lightweight by design, several infrastructure choices improve its performance, scalability, and maintainability:

1. **Service Separation:** The authentication/user-data API (Next.js) and data ingestion/query API (Flask) can be deployed on separate servers or containers, allowing independent scaling.
2. **Containerization:** Both services can be packaged into Docker containers, simplifying deployment and dependency management.
3. **Security & Data Isolation:** User authentication runs in a separate service from dataset uploads and query execution, reducing attack surface.
4. **Stateless Flask API:** The Flask service is stateless aside from the SQLite dataset storage, making horizontal scaling easier.

5. **Database Backups:** Automated backup routines for both SQLite databases ensure recoverability in case of corruption or accidental deletion.
6. **Potential Future Upgrades:** The architecture allows an easy switch to more scalable databases (e.g., PostgreSQL for user data, DuckDB for analytics) if needed without major code rewrites.

This backend design balances simplicity with extensibility, keeping deployment friction low while ensuring that the system can evolve as data volumes and user demands grow.

Evaluation and Testing

IN THIS CHAPTER

In this chapter, we detail the evaluation and testing methodology used to assess the performance and accuracy of the ttyDB system in converting natural language queries into SQL statements. The testing process is designed to ensure that the system can handle a variety of query types and data structures effectively.

1. **Testing Methodology:** This section outlines the approach taken to evaluate the system, including the setup, test cases, and evaluation criteria.
2. **Performance Metrics:** This section presents the metrics used to measure the system's performance, including success rates, error handling, and query execution times.
3. **Results of Testing:** This section summarizes the results of the testing process, highlighting the system's strengths and areas for improvement.

Testing Methodology

The system is evaluated using an in-memory SQLite database populated with representative sample data across multiple related tables including Employees, Customers, Products, Orders, Reviews, and OrderDetails. This schema reflects a realistic business domain to challenge the model with diverse query types.

A set of natural language test queries covering aggregation, filtering, joins, subqueries, and conditional logic is defined to benchmark the model's ability to generate accurate SQL statements. Each query is processed through an incremental retry mechanism that allows up to five attempts to correct generated SQL in case of execution errors.

The test harness works as follows:

1. For each natural language query, the system sends the query, optionally with the previous failed SQL and error message, to the fine-tuned text2sql model.
2. The model returns a SQL query along with an explanation.
3. The SQL is executed against the in-memory database.
4. If the query succeeds, the test counts it as passed.
5. If execution fails, the error and faulty query are fed back to the model for correction.
6. This retry loop continues up to five times, providing the model chances to self-correct its SQL output.
7. If after five attempts the query still fails, the test is marked as failed.
8. This approach mimics real-world usage where natural language inputs may require iterative refinement to generate valid queries, and it measures the robustness of the model's conversion and self-correction capabilities.

Sample queries used test common analytical requests such as:

- Finding highest priced products
- Listing products with filters
- Aggregations by category or customer
- Joining orders with customer and employee data
- Handling review data with conditional logic

The final benchmark reports the number of queries correctly executed out of the total, giving a quantitative measure of system accuracy and resilience.

The logs for the test runs can be found [here](#)

Performance Metrics

The system was evaluated on a benchmark suite consisting of 23 natural language queries covering a variety of database operations including selection, aggregation, joins, filtering, and nested queries.

1. **Success Rate:** The system successfully generated and executed valid SQL for 22 out of 23 queries, yielding a success rate of approximately 95.7%.
2. **Retries:** On average, queries requiring correction were resolved within 2 retries, with a maximum retry limit set to 5 attempts per query.
3. **Error Types:** Initial failures predominantly involved SQL syntax errors such as incorrect aliases, malformed clauses, or typographical mistakes.
4. **Correction Mechanism:** The iterative feedback loop—feeding back error messages and prior SQL to the model—proved effective in resolving errors without human intervention.
5. **Timeouts:** No API call timeouts occurred during testing, with a per-call timeout threshold of 60 seconds.
6. **Failure Analysis:** The single failure involved a persistent syntax issue related to a subquery with NOT IN clause, which the model was unable to correct within the retry limit.
7. **Execution:** Successful queries returned meaningful result sets consistent with the expected database schema and semantics.

Overall, these metrics demonstrate the system's high reliability in converting natural language queries into syntactically correct and executable SQL statements, validating the efficacy of the combined fine-tuned model and iterative error correction approach.

Results of Testing

The testing revealed several key insights about the system's performance in converting natural language queries into SQL:

1. The system generated SQL queries that matched the intended semantics of the input questions with a high degree of accuracy.
2. Most queries were correctly handled on the first attempt, demonstrating the fine-tuned model's capability to understand diverse query types.
3. Queries involving straightforward selection, aggregation, and grouping operations consistently produced correct results.
4. Complex queries requiring multiple joins, nested subqueries, or advanced filtering often required one or two iterations to fix minor syntactic or semantic errors.
5. The iterative self-correction mechanism effectively improved query accuracy by leveraging runtime error feedback.
6. The only query that failed after the maximum retry attempts involved a NOT IN subquery, indicating a potential limitation in handling certain SQL constructs.

7. The model demonstrated strong understanding of the database schema, correctly referencing table and column names in most cases.
8. Output explanations accompanying the generated SQL queries helped in diagnosing errors and understanding query logic during testing.

These results affirm the practical viability of the proposed approach for real-world natural language querying of structured data, with room for improvement in handling rare edge cases.

Challenges and Lessons Learned

IN THIS CHAPTER

In this chapter, we reflect on the technical challenges encountered during the development and testing of ttyDB, as well as the lessons learned from these experiences. The project faced several hurdles related to model limitations, error handling, and system architecture, which provided valuable insights for future improvements.

1. **Technical Challenges:** This section discusses the technical challenges faced during development, including model limitations, error handling, and system architecture issues.
2. **Lessons Learned:** This section summarizes the key lessons learned from the project, including insights into model performance, error correction strategies, and system design considerations.

Technical Challenges

During development and testing, several challenges emerged:

1. **Model Size Limitations:** Using a comparatively smaller 3B parameter model constrained the complexity of queries the system could confidently generate, necessitating extensive fine-tuning and iterative correction.
2. **Hallucinations and Syntax Errors:** The model occasionally hallucinated incorrect SQL syntax, table names, or column names, leading to frequent execution errors. Mitigating this required building a robust feedback loop to detect and correct errors automatically.
3. **Tool Calling Difficulties:** A custom implementation of tool calling was tested but abandoned due to excessive and irrelevant calls, including hallucinated tool names. Although prompt engineering reduced hallucinations somewhat, tool calling remained unreliable.
4. **Error Correction Strategy:** The final strategy—where the model receives its own generated SQL and the execution error as input for correction (termed Iterative Self-Feedback Loop)—proved effective in improving query quality without external intervention.
5. **Complex Query Handling:** Queries involving complex SQL constructs like nested subqueries or specific JOIN conditions sometimes required multiple correction cycles, revealing the need for further model enhancements or complementary heuristics.
6. **Schema Awareness:** Explicitly providing the database schema as context was critical in reducing hallucinations and improving the accuracy of generated SQL.

Lessons Learned

1. **Importance of Iterative Correction:** Implementing an iterative self-feedback loop where the model reviews its own SQL errors significantly improves query accuracy and robustness, especially when working with smaller models.
2. **Schema Context is Crucial:** Providing a detailed schema description upfront helps reduce hallucinations and guides the model to generate syntactically and semantically correct SQL queries.

3. Limitations of Tool Calling with Smaller Models: Despite its potential, tool calling can introduce instability and hallucinations in smaller models, making a simpler error-feedback approach more reliable in practice.
4. Fine-Tuning Depth Matters: Extensive fine-tuning with a sizable dataset (8.4k samples) was necessary to compensate for the smaller model size and enhance domain-specific SQL generation.
5. Fallback Mechanisms Enhance Reliability: Having a fallback to a stronger model (e.g., Google Gemini 2.5 flash) ensures robustness for edge cases where the primary model fails repeatedly.
6. Balancing Model Complexity and System Constraints: Selecting a smaller model allowed deployment within resource limits but required compensatory architectural design choices to maintain accuracy.

Future Work

IN THIS CHAPTER

In this chapter, we outline the planned enhancements and potential features for ttyDB, focusing on improving the system's capabilities, user experience, and overall performance. The future work aims to address the current limitations and expand the system's functionality to better serve users' needs.

1. **Planned Enhancements:** This section describes the planned enhancements to the system, including improvements to the tool calling mechanism, adaptive query refinement, and multi-modal data support.
2. **Potential Features:** This section outlines potential features that could be added to enhance the user experience, such as natural language query history, visual query builders, and real-time query suggestions.

Planned Enhancements

1. **Improved Tool Calling Mechanism:** Develop a more reliable tool calling system with better prompt engineering to reduce hallucinations and irrelevant calls.
2. **Adaptive Query Refinement:** Introduce smarter retry logic that prioritizes fixes based on error types to optimize query correction efficiency.
3. **Multi-Modal Data Support:** Expand support for querying non-SQL data sources like CSV and JSON through natural language.
4. **Contextual and Multi-Turn Querying:** Enable context-aware conversations that remember previous queries for more complex interactions.
5. **Enhanced Explainability:** Provide clearer and more user-friendly explanations of generated SQL queries to aid understanding and trust.

Potential Features

1. **Natural Language Query History:** Maintain a user query history for quick access and reuse of past queries.
2. **Visual Query Builder:** Interactive UI to construct and refine SQL queries visually alongside natural language input.
3. **Query Result Visualization:** Charts, graphs, and tables to help users interpret query results easily.
4. **Real-Time Query Suggestions:** Provide dynamic autocomplete and query suggestions as users type their natural language input.

Conclusion

Summary of Achievements

1. Developed a natural language to SQL conversion system fine-tuned on a 3B-parameter Qwen2.5-coder model.
2. Implemented an iterative SQL validation and self-correction loop, reducing query errors and hallucinations.
3. Successfully created and tested a representative multi-table in-memory SQLite database with realistic sample data.
4. Achieved a high accuracy benchmark, correctly answering 22 out of 23 diverse natural language queries.
5. Designed a scalable backend architecture with separated microservices for file handling and model inference.

Final Thoughts

This project demonstrates the feasibility of converting natural language queries into executable SQL statements using a fine-tuned 3B-parameter model. Despite inherent challenges such as hallucinations and syntax errors, an iterative correction mechanism coupled with fallback strategies ensured high accuracy and robustness. The system successfully handled a diverse set of complex queries against a representative database schema, achieving a benchmark accuracy of 22 out of 23 queries. Future improvements in tool calling, query refinement, and multi-modal support will further enhance usability and performance, making natural language database querying accessible to a broader audience.

References

IN THIS CHAPTER

In this section, we provide a list of references and resources that were instrumental in the development of ttyDB. These resources include documentation, tutorials, and articles related to natural language processing, SQL generation, and the technologies used in this project.

Fine-tuning LLMs

1. Fine-Tuning LLMs with Unsloth and Ollama: A Step-by-Step Guide
2. Tutorial: How to Finetune Llama-3 and Use In Ollama
3. Qwen2.5-Coder:3B on Ollama

Synthetic Data Generation

1. Gemini API | Google AI for Developers
2. Pydantic Models Documentation
3. ChatGPT

SQLite and Pandas

1. SQLite3 Documentation
2. Pandas Documentation

Frontend Development

1. Next.js Documentation
2. ShadCN/UI Component Library
3. TailwindCSS Documentation

Backend Development

1. Flask Documentation
2. Prisma Documentation
3. Ollama Documentation
4. Basic Usage | BetterAuth