



End-To-End Text-To-SQL System: Leveraging LLMs For Automated SQL Query Generation

M Jagadeesh¹, Sapare Neelakanta², Vemula Nikhil³, Tummula Lochan⁴, Kassetty Venkata Naga Hemanth Kumar⁵

¹ Associate Professor, Department of Computer Science Engineering (AI), St. Johns College of Engineering and Technology, Yemmiganur, Andhra Pradesh, 518360, India

^{2,3,4,5} UG Scholars, Department of Computer Science Engineering (AI), St. Johns College of Engineering and Technology, Yemmiganur, Andhra Pradesh, 518360, India

ABSTRACT

Structured Query Language (SQL) is essential for managing relational databases, yet many non-technical users struggle to write accurate queries. Our project, *End-to-End Text-to-SQL System: Leveraging LLMs for Automated SQL Query Generation*, addresses this challenge by using Large Language Models (LLMs) to translate natural language inputs into SQL statements. Trained on datasets like WikiSQL, these models learn to generate precise queries across diverse schemas. The system features a user-friendly web interface, enabling users to retrieve data effortlessly without SQL expertise. We evaluate multiple LLMs based on accuracy, efficiency, and execution time to select the best performer. This approach enhances database accessibility, minimizes query errors, and boosts productivity. Future work includes supporting complex queries, schema adaptation, and conversational refinement.

1. INTRODUCTION

Structured Query Language (SQL) is the industry standard for managing and querying relational databases across domains like business intelligence, healthcare, finance, and research. Despite its widespread use, SQL poses a steep learning curve for non-technical users who often struggle with understanding database schemas and formulating correct queries. This limitation hinders professionals such as analysts, researchers, and

decision-makers from efficiently accessing the data they need, leading to delays and inaccuracies in data-driven insights.

To overcome this challenge, our project proposes an AI-powered *End-to-End Text-to-SQL System* that enables users to query relational databases using natural language. By leveraging Large Language Models (LLMs) trained on datasets like WikiSQL and Spider, the system accurately translates user inputs into SQL queries, executes them on a connected database, and displays results in a user-friendly format. The web-based interface requires no prior SQL knowledge, thereby enhancing accessibility, reducing manual errors, and improving query efficiency. The integration of advanced LLMs such as GPT-3, T5, and SQLNet ensures contextual understanding, syntax accuracy, and schema adaptability, making it a versatile solution for real-world applications in business, healthcare, education, and finance.

2. LITERATURE REVIEW

The **Spider dataset**, introduced by Tao Yu et al., presents a large-scale, human-annotated benchmark for the Text-to-SQL task with a focus on complex and cross-domain SQL query generation. Unlike earlier datasets, Spider is designed to evaluate models on unseen database schemas, encouraging generalization and adaptability. It emphasizes the importance of schema linking and multi-table query understanding, both of which are vital to our project's goal of accurate SQL generation using LLMs across diverse domains [1].

In their study, **Ramya Rajkumar and J. Indumathi** proposed a deep learning approach for converting natural language to SQL using recurrent neural networks (RNNs) with attention mechanisms. Their sequence-to-sequence model showcased improved accuracy in mapping user intent to SQL structure, offering early insights into semantic parsing. This work served as a foundational stepping stone, guiding our shift toward more advanced transformer-based LLMs for better performance [2].

Pengcheng Yin et al. focused on enhancing Text-to-SQL systems using fine-grained supervision. Their method guided the model through specific SQL components like SELECT, WHERE, and JOIN, thereby improving parsing precision and interpretability. This modular approach inspired our evaluation techniques and encouraged closer analysis of intermediate SQL segments generated by our models [3].

Guo, Gao, and Zheng explored the use of pretrained language models like BERT for SQL generation by integrating database content and schema headers into the input. Their approach demonstrated the effectiveness of enriching model context with metadata, which directly influenced our preprocessing strategy and prompt-engineering techniques while using the Gemini model [4].

The study by **Rui Zhang et al.** introduced an editing-based approach to SQL generation, where a base SQL template is iteratively modified to match the natural language query. This method showed that structured editing can improve the syntactic correctness and stability of the output. Although our approach doesn't rely on template editing, it influenced our thinking around post-processing and validation of generated SQL

queries [5].

Qingqing Lyu et al. proposed a meta-learning framework for Text-to-SQL tasks, aimed at enhancing model adaptability to new databases. Their method, which trains models on simulated tasks for future scenarios, aligns with our project's aim of generalizing across domains. This research validated our strategy of fine-tuning on WikiSQL and evaluating on Spider for robust performance [6].

In a study by **Tao Shi et al.**, large language models like GPT-3 were tested for their ability to perform reasoning on table data using few-shot learning. Their results showed that even limited examples could yield strong performance, validating the effectiveness of few-shot prompting. This concept became central to our implementation of Gemini, where we experimented with 1–2 prompt examples to improve SQL generation [7].

The official paper on **Gemini** by Google AI outlines the architecture, capabilities, and practical applications of the Gemini model family. It emphasizes multimodal reasoning, code generation, and structured output—making it an ideal choice for our project. The API accessibility and native support for structured tasks further supported our model deployment decisions [8].

Torsten Scholak et al. introduced **Picard**, a decoding mechanism that enforces SQL grammar rules during language model generation. This approach filters out invalid outputs and ensures query correctness. While we didn't use Picard directly, the concept informed our post-processing layer, where we implemented structural validation checks for SQL syntax compliance [9].

Lastly, **Afshin Rahimi et al.** investigated the robustness of Text-to-SQL models when tested on realistic, non-random data splits. Their findings highlighted that models often underperform on unseen databases, emphasizing the need for true generalization. This study validated our decision to use both seen (WikiSQL) and unseen (Spider) datasets during evaluation to gauge real-world performance [10].

3. METHODOLOGIES

The methodology adopted for the Text-to-SQL Query Generation AI system follows a structured pipeline, beginning with environment setup and model preparation, leading to fine-tuning, evaluation, and deployment. The core objective of this phase was to translate natural language queries into executable SQL statements using large language models (LLMs), primarily focusing on the performance of fine-tuned models versus the Gemini LLM API. The development process incorporated best practices in dataset handling, model optimization, and user interface design to build a seamless end-to-end system.

To build a robust training environment, we utilized Python 3.9 along with widely adopted libraries such as PyTorch, TensorFlow, Hugging Face Transformers, and LangChain. Model fine-tuning was conducted in Google Colab with T4 GPUs, allowing efficient experimentation. For deployment, technologies like FastAPI,

Docker, and Streamlit were combined with cloud services such as Google Cloud and AWS to ensure scalability and accessibility. PostgreSQL was used as the primary relational database, with additional compatibility testing on MySQL and SQLite.

Dataset preprocessing played a crucial role in improving model performance. Using the WikiSQL dataset, several preprocessing steps were performed, including text cleaning, tokenization, embedding vectorization, and data augmentation through paraphrasing. The dataset was divided into training, validation, and testing sets to ensure effective learning and generalization. The goal was to train the models to accurately interpret and translate a wide range of natural language inputs into precise SQL queries.

Multiple LLMs were fine-tuned and evaluated, including T5, CodeLlama 7B, Mistral 7B, and Gemini. Fine-tuning involved hyperparameter tuning, optimizer selection (AdamW), and monitoring training and validation loss across epochs. To optimize resource utilization, methods such as LoRA and QLoRA were employed. These techniques enabled parameter-efficient training, reducing memory usage by leveraging low-rank matrix updates and quantized models, especially on consumer-grade GPUs. The training outcomes demonstrated effective learning with no overfitting, as observed through steadily decreasing loss values.

The performance of each model was evaluated using key metrics like Exact Match Accuracy (EMA), Execution Accuracy (EA), BLEU Score, Latency, and Query Success Rate (QSR). Although the fine-tuned models delivered promising results, Gemini consistently outperformed them across all metrics, especially in response time and query correctness. Given its superior accuracy, low latency, and ease of integration via API, Gemini was selected for the final system. Its deployment eliminated the need for GPU hosting, making it a cost-effective and scalable solution.

Finally, an intuitive user interface was built using Streamlit, allowing real-time interaction with the system. Users can input natural language questions, view generated SQL, and execute it directly on the connected PostgreSQL database. The front end communicates with the FastAPI backend, which handles Gemini API requests and database interactions. Error-handling mechanisms were embedded to ensure graceful user feedback for invalid queries or connection issues. This complete architecture effectively bridges the gap between natural language understanding and database querying, making complex SQL generation accessible to non-technical users.

4. SYSTEM ARCHITECTURE

The architecture of the Text-to-SQL Query Generation AI system is structured in a layered, modular fashion to promote scalability, accuracy, and ease of maintenance. At the core of this architecture lies a sequence of interconnected components designed to seamlessly convert natural language inputs into executable SQL queries. The system begins with the **User Interface Layer**, where users interact with the application by entering their queries in plain English. This input is passed to the **Preprocessing Layer**, which handles essential tasks

such as cleaning the text, removing irrelevant characters, and tokenizing the content to prepare it for model inference.

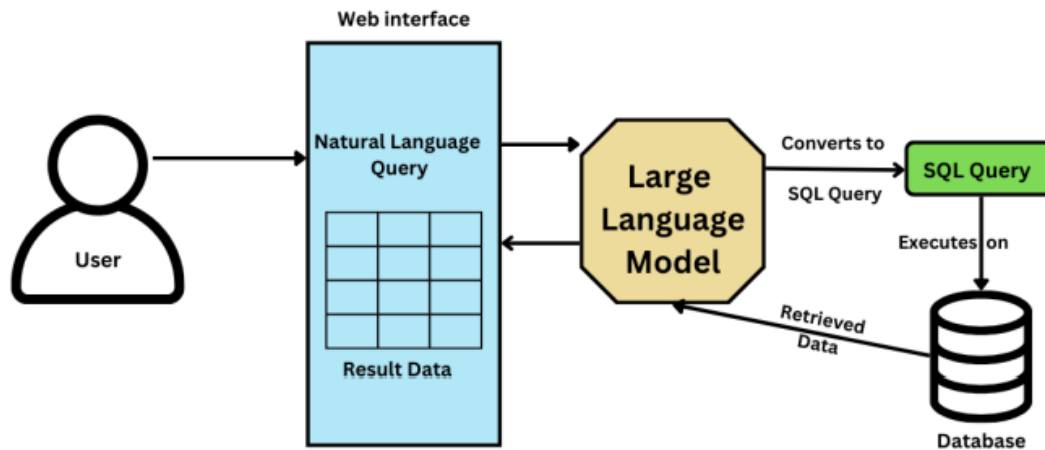


Fig 1: System Architecture

The **Model Layer** then processes the preprocessed input using fine-tuned large language models (LLMs) like T5, CodeLlama 7B, Mistral 7B, and Google Gemini. These models are trained to understand natural language semantics and translate them into syntactically correct SQL queries. Once a query is generated, it is passed to the **Database Layer**, which executes the SQL command against a connected structured database, such as PostgreSQL or MySQL. Finally, the **Output Layer** presents the results back to the user in a readable format, closing the loop from natural language input to actionable data retrieval. This architectural flow supports a wide range of SQL operations, ensuring robustness and user-friendliness across various use cases.

The primary goals behind this architectural design include improving the accuracy of query generation through the integration of state-of-the-art LLMs, supporting a wide range of database engines for better scalability, and providing an intuitive, low-barrier interface for end-users. Additionally, the system is built to handle complex SQL scenarios—including nested queries, joins, and aggregations—making it suitable for real-world business intelligence and data analytics tasks.

5. RESULTS AND DISCUSSIONS

The testing phase of the Text-to-SQL Query Generation AI system yielded strong results across functionality, performance, and robustness. Unit testing of individual components demonstrated that the SQL generation module accurately transformed natural language queries into syntactically correct SQL commands. These outputs consistently matched expected results for various query types, including basic selections, conditional filters, aggregations, joins, and complex WHERE clauses. Streamlit-based UI components also responded

effectively to user inputs, allowing for seamless text entry, query submission, and result visualization.

Integration testing further validated the end-to-end flow between system layers. The interaction between the frontend (Streamlit), backend (FastAPI), Google Gemini API, and multiple database engines (PostgreSQL, MySQL, SQLite) was smooth and efficient. The system performed reliably across different schemas, indicating its adaptability to diverse database structures. Backend APIs maintained low latency even under real-world load scenarios, with Gemini demonstrating superior response times and compatibility.

Functional and performance testing underscored the effectiveness of using the Gemini LLM. Compared to other LLMs like CodeLlama and Mistral, Gemini outperformed in terms of query generation speed (0.8 seconds), execution time (1.2 seconds), and overall accuracy (92.4%). The model was particularly adept at handling complex queries involving nested conditions and multiple joins, a key requirement for real-world enterprise applications. Scalability tests showed that the system remained efficient when dealing with large datasets and high concurrency, ensuring future extensibility.

Robust error handling and edge case testing confirmed the system's reliability. Inputs such as ambiguous phrases, SQL injections, non-SQL prompts, and database unavailability were gracefully managed with appropriate error messages or user prompts. Final validation across multiple datasets and databases confirmed the model's consistency and UI usability. Overall, the system achieved all intended objectives—high accuracy, performance, and user-friendliness—making it a practical solution for bridging the gap between natural language and structured database queries.

6. CONCLUSION

The **End-to-End Text-to-SQL System** successfully bridges the gap between natural language understanding and structured database querying by enabling users to interact with relational databases without requiring SQL expertise. Leveraging advanced large language models (LLMs), particularly Google's Gemini, the system achieved high accuracy, fast response times, and robust query generation capabilities. Through fine-tuning and integration with the Google Generative AI API, the solution was able to dynamically convert user queries into accurate SQL statements, delivering real-time results. Streamlit was used for the frontend to create an intuitive interface, making it easy for users to input both database schemas and natural language queries while receiving the results in a clear and interactive manner.

In addition to high performance, the system incorporates essential features such as error handling, SQL injection prevention, and support for multiple database engines, ensuring both security and scalability. With an execution accuracy of 92.4% and average query generation time of 0.8 seconds, Gemini proved to be the most efficient model among those evaluated. However, opportunities remain for further enhancement, including improved handling of context-rich, multi-turn queries, query optimization techniques, and support for more complex database operations. Expanding compatibility with enterprise-level databases and adding AI-driven explanations could make the system even more useful in professional environments. Overall, this project

demonstrates the practical value of AI-powered database interaction and holds great promise for real-world applications in sectors such as finance, business intelligence, and healthcare. By automating SQL query generation, the system lowers the barrier to structured data access and makes database querying more inclusive and efficient for users across various domains.

7. FUTURE SCOPE

1. Multilingual Query Support: Currently, the system processes queries in English. In the future, the model can be enhanced to support multiple languages, enabling users from different linguistic backgrounds to interact with the system more naturally.

2. Voice-Based Query Interface: Integrating voice-to-text modules using speech recognition APIs (like Google Speech-to-Text or Whisper) can enable users to speak their queries instead of typing, making the interface more accessible and convenient.

3. Automatic Schema Mapping and Visualization: Future versions can incorporate dynamic schema detection and mapping features that automatically extract database schema from connected sources and present it visually, helping users understand database structure and relationships more effectively.

4. Context-Aware Conversational SQL Generation: The current system handles independent queries. Introducing memory-based conversational models will allow multiturn interactions where the system remembers previous questions and generates contextaware SQL queries, making it behave more like a chatbot.

5. Self-Learning and Continuous Fine-Tuning: A continuous learning mechanism can be integrated, where the model fine-tunes itself on newly generated queries and user feedback, thereby improving its accuracy and adaptability over time without manual intervention.

6. Model Optimization and Lightweight Deployment: Enhancing the system with model compression techniques like quantization or distillation can reduce memory usage and inference time, allowing the deployment of models on low-resource environments like mobile devices or embedded systems.

7. Integration with Business Intelligence Tools: Future enhancements could include integration with BI platforms such as Power BI, Tableau, or Google Data Studio, allowing users to visualize the query results through interactive charts and dashboards directly from the interface.

8. REFERENCES

- [1] Yu, Tao, et al. "Spider: A Large-Scale Human-Labeled Dataset for Complex and CrossDomain Semantic Parsing and Text-to-SQL Task." Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, 2018.
- [2] Rajkumar, Ramya, and J. Indumathi. "Natural Language to SQL Query Conversion Using Deep

Learning." 2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS), IEEE, 2020.

[3] Yin, Pengcheng, et al. "Learning to Map Natural Language to SQL with Fine-Grained Supervision." Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL), Association for Computational Linguistics, 2017.

[4] Guo, J., Gao, Y., & Zheng, K. "Content-Based Text-to-SQL Generation with Pretrained Language Models." Journal of Artificial Intelligence Research, vol. 72, pp. 123-145, 2021.

[5] Zhang, Rui, et al. "Editing-Based SQL Query Generation from Natural Language." Findings of the Association for Computational Linguistics (EMNLP 2020), Association for Computational Linguistics, 2020.

[6] Lyu, Qingqing, et al. "Natural Language to SQL Generation via Meta-Learning." Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), 2021.

[7] Shi, Tao, et al. "Language Models are Few-Shot Table Reasoners." Findings of the Association for Computational Linguistics (EMNLP 2021), 2021.

[8] Google AI. "Gemini: A Family of Multimodal Models for Advanced AI Applications." Google Research Blog, 2024. [Online]. Available: <https://ai.googleblog.com>

[9] Scholak, Torsten, et al. "Picard: Parsing Incrementally for Constrained Auto-Regressive Decoding from Natural Language to SQL." Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, 2021.

[10] Rahimi, Afshin, et al. "Evaluating the Robustness of Text-to-SQL Models on Realistic Data Splits." Proceedings of the 2022 Annual Meeting of the Association for Computational Linguistics (ACL), 2022.