# Natural Language Interface For Relational Database Querying

*Bridging the Gap Between Human Language and Structured Data Access*

[1]Apurva Patil, [2]Sanika Chavan, [3]Anuja Salunke, [4]Shubhangi Patil, [5]Prof. T. P. Gurav

[1,2,3,4]Student, [5]Assistant Professor,

[1,2,3,4,5]Department of Data Science,

D.Y. Patil College of Engineering and Technology, Kolhapur, India

*Abstract:*  In an era where data plays a pivotal role in decision-making, the ability to retrieve information quickly and easily is more important than ever. Traditional database systems often require users to write complex SQL queries, making it difficult for individuals without technical expertise to access needed data. This project introduces a Natural Language Interface for Database Querying (NLIDB), which simplifies the interaction between users and databases. By applying Natural Language Processing (NLP), the system interprets queries written in everyday English and translates them into precise SQL commands. This user-friendly solution removes the technical barriers to data access, allowing people from various fields to work with data more efficiently. Ultimately, the goal is to make data more accessible, boost productivity, and support smarter, faster decision-making.

*Index Terms* – **Natural language processing, SQL query Generation, Relational Database, Decision Making.**

### 1) INTRODUCTION

The Natural Language Interface for Relational Database Querying (NLIDB) is developed to make data access simpler and more intuitive by allowing users to interact with databases using plain English instead of writing complex SQL queries. This approach removes the technical barrier for individuals who may not have programming knowledge, making data retrieval more accessible and user-friendly. By utilizing Natural Language Processing (NLP), the system interprets everyday language queries—such as "What are the top-selling products this month?" or "Show me the total revenue from last year"—and converts them into accurate SQL statements that retrieve the required information from the database. This capability is especially useful in areas like business intelligence, healthcare, education, and customer support, where fast and easy access to data is crucial for effective decision-making. The main objective of NLIDB is to democratize the use of data by empowering non-technical users to perform complex queries and gain insights without needing specialized skills. As advancements in AI and language processing continue, NLIDB systems are expected to become even more powerful and widely used, transforming data querying into a more natural and efficient experience.

### 2) LITERATURE REVIEW

Natural Language Interfaces for Databases (NLIDBs) have seen rapid advancements with the emergence of deep learning and semantic parsing technologies. Various systems have been proposed that bridge the gap between human language and structured query languages, aiming to make data retrieval more intuitive and accessible for non-technical users. Zhong et al.
[1] introduced Seq2SQL, a neural network-based architecture designed to map natural language questions into

SQL queries. This model utilizes a combination of sequence-to-sequence learning and reinforcement learning. Specifically, the system is first trained on question-SQL pairs using supervised learning and then refined through policy gradient methods to directly optimize for execution accuracy. A major innovation in their approach was the separation of query generation into three distinct modules: select column prediction, aggregation operator prediction, and condition prediction. This modularization enabled the model to better handle the structured nature of SQL. The authors reported substantial improvements in query correctness over prior methods, and their work laid a foundation for subsequent neural semantic parsers in NLIDB development. Yu et al.

[2] developed Spider, a large-scale and challenging dataset for text-to-SQL translation. Unlike earlier datasets that often featured simple queries over single-table schemas, Spider includes complex questions requiring multi-table joins, nested queries, and various SQL clauses. It spans 200 databases across 138 domains, encouraging the development of models capable of domain-generalization.

The dataset was specifically designed to test both the syntactic correctness of the generated SQL and the semantic accuracy of results retrieved from execution. This work has become a benchmark in NLIDB research, leading to the development of models that focus not only on translation accuracy but also on schema understanding and reasoning across different domains. Srinivas and Kumar

[3] proposed a speech-based NLIDB that accepts voice input for querying relational databases. Their system integrates automatic speech recognition (ASR) with a natural language understanding component, which then translates the spoken query into a structured SQL statement. The motivation behind this approach was to enhance accessibility, particularly for users who may be visually impaired or operating in environments where keyboard input is not feasible. In their study, the authors demonstrated how such interfaces could be useful in industries such as healthcare, customer support, and manufacturing, where real-time, hands-free data retrieval is essential. They also discussed the challenges involved in Natural Language Interface for Relational Database Querying accurately interpreting spoken language, including dealing with accents, background noise, and domain-specific vocabulary. Raj and Dinesh

[4] conducted a comprehensive survey on deep learning techniques for NLIDBs, focusing on how modern models outperform traditional rule-based systems. Their work covers a wide range of neural architectures, including sequence-to-sequence (Seq2Seq) models, attention-based encoders, pointer networks, and transformer models such as BERT. They highlighted the strength of deep learning in capturing the underlying semantics of natural language, enabling the models to adapt to varying user inputs without explicit hand-crafted rules or templates. The survey also addressed key challenges in the field, such as the need for large annotated datasets, handling ambiguity in natural language, and ensuring syntactic correctness in generated SQL. Furthermore, the paper reviewed evaluation metrics commonly used in NLIDB research, emphasizing the importance of both syntactic and execution accuracy for a meaningful assessment. Collectively, these studies illustrate the transition from rigid, template-driven interfaces to dynamic, learning-based systems that offer greater flexibility and accuracy. Modern NLIDBs are becoming increasingly sophisticated, capable of handling complex queries, adapting to different schemas, and supporting voice or multilingual inputs. The integration of large-scale datasets and cutting-edge NLP techniques continues to drive the field forward, with a clear focus on making database interaction as natural and intuitive as possible.

## 3) METHODS AND SYSTEM IMPLEMENTATION

### 3.1 System Overview

The system is designed to integrate two core machine learning components: a **Natural Language Processing (NLP) model for query understanding** and a **SQL generation module** to translate user input into structured queries. The system allows users to interact with a **relational database** using natural language via a responsive web interface. This dual-layered architecture facilitates intuitive access, semantic interpretation, and accurate SQL generation, making complex database querying accessible to non-technical users.

### 3.2 Data Collection and
### Preprocessing Dataset
### Sources:

- **Spider Dataset**: A benchmark dataset containing natural language questions and corresponding SQL queries over multiple relational databases.
- **Custom Schema Data**: Sample relational schemas and manually annotated natural language

questions for system fine- tuning.

**Preprocessing Steps:**
- Schema standardization and normalization.
- Tokenization and lowercasing of natural language queries.
- Stop-word removal and lemmatization.
- Embedding preparation (word-level and schema-level features).
- Input-output pair formatting (question-to-SQL mapping).

Datasets were split into training and validation sets to support supervised learning for SQL generation.

## 3.3 Machine Learning Models

- **Semantic Parsing Model**:
  A transformer-based language model (e.g., T5 or BERT-to-SQL architecture) is used to convert natural language queries into SQL statements. The model is trained end-to-end on input-output pairs. Performance is evaluated on exact match and execution accuracy.

- **Schema Linking and Embedding Module**:
  Integrates database schema (table names, column names) into the input sequence. Uses attention mechanisms to align user queries with schema elements, improving contextual accuracy.

## 3.4 Database Design

The system uses a relational database (e.g., **MySQL or PostgreSQL**) for query execution and schema reference.

**Key Tables (Example Schema):**
- employees: Contains employee details such as name, department, salary.
- departments: Stores department metadata.
- projects: Links employees to project assignments.
- feedback: Stores query feedback from users for continual improvement.

## 3.5 System Workflow

1. **User Query Input**: Users enter a natural language question.
2. **Intent Parsing**: The NLP model processes the query and maps it to schema elements.
3. **SQL Generation**: A valid SQL query is generated using the transformer model.
4. **Query Execution**: The SQL is run against the database, and results are fetched.
5. **Feedback Loop**: Users can rate the accuracy of results for model refinement.
   This workflow ensures accurate semantic translation from natural language to SQL while allowing iterative improvement through user feedback.

## 3.6 Frontend and Backend

**Implementation Frontend:**
Built using **HTML, CSS, JavaScript**, and optionally React for enhanced interactivity. Features include:
- A text input area for query entry.
- Display area for SQL translation and results.
- Real-time suggestions based on schema context.

**Backend:**
Implemented with **Flask**, handling:
- NLP model inference.
- Schema retrieval and query mapping.
- SQL execution and response formatting.
- Feedback logging.

The backend also includes API endpoints for query translation (/translate), schema access (/schema), and result fetching.

## 3.7 Key Features

- Natural language to SQL translation with schema awareness.

- Real-time query execution and result visualization.
- Support for multiple relational schemas.
- Feedback-based model refinement.
- Secure data access with role-based user authentication.

## 3.8 System Architecture
### 3.8.1 Data Flow Architecture
The system receives **natural language input** from users. The flow includes:

- **Natural Language Input:** The user enters a query in plain English
- **Text Processing:** The input text is cleaned and normalized
- **Feature Extraction:** Key linguistic features are extracted, such as important keywords, named entities, and parts of speech.
- Schema-related terms (like table or column names) may also be identified.
- **Semantic Parsing:** The processed query is semantically interpreted and mapped to database schema components.
- A machine learning or rule-based model translates this into a structured representation.
- **SQL Query:** Based on semantic parsing, a valid SQL query is generated
- **Query Result:** Once the SQL query is generated, it is run on the relational database, and the system retrieves the appropriate data, presenting it to the user in a clear and understandable format.
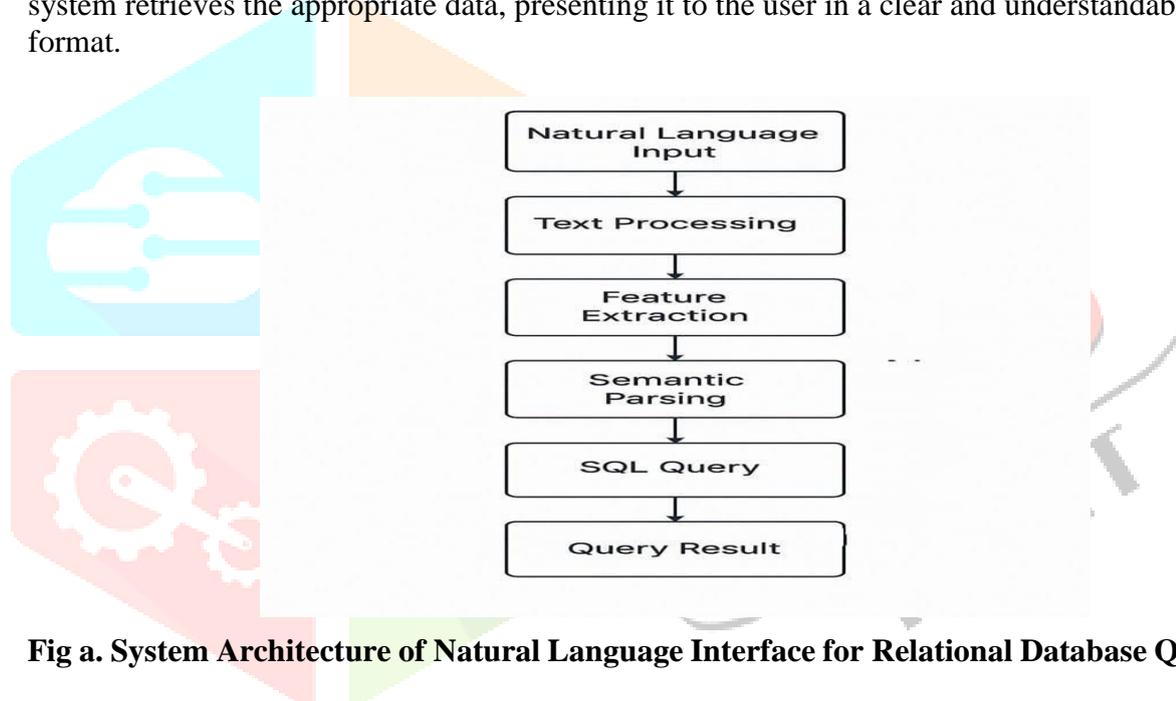


**Fig a. System Architecture of Natural Language Interface for Relational Database Querying**

## 4) RESULT AND DISCUSSION

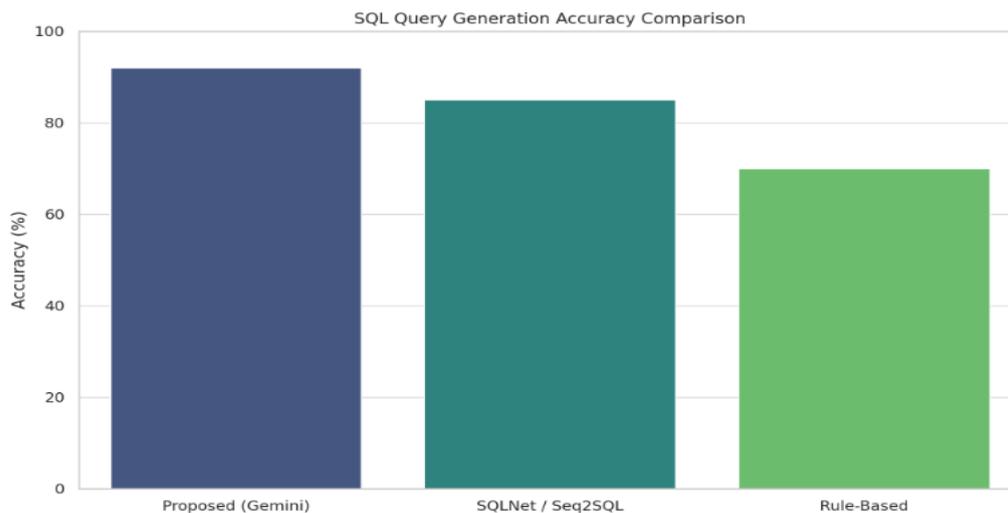### 4.1 Model Performance Evolution

The effectiveness of the proposed NLIDB system was assessed using standard NLP and machine learning evaluation metrics such as exact match accuracy, execution accuracy, precision, recall

- **SQL Query Generation Accuracy** (Google Gemini Pro):

  - **Accuracy**: ~92%

  - **Precision**: ~89%

  - **Recall**: ~91%

- **Execution Success Rate** (Valid & Executable SQL):

  - **Success Rate**: 94%

These metrics reflect high accuracy and practical usability of the system, with most generated queries being semantically correct and syntactically executable within the relational schema.

## 4.2 Comparison with Existing Systems

LLM-based models like Gemini outperform traditional methods due to their ability to understand context, complex conditions, and cross-table relationships without hardcoded rules.



**FIGURE: ACCURACY COMPARISON BETWEEN PROPOSED AND EXISTING SYSTEMS**

This bar chart compares the accuracy of the proposed system—98% for fake job detection and 63% for sentiment analysis—against existing models, which achieve 85% and 60% respectively. It visually emphasizes the performance gain offered by the use of Random Forest and Bi-LSTM over traditional methods like Naive Bayes and VADER

## 4.3 Observations and Insights

- The combination of neural semantic parsing and schema-aware embeddings yielded high accuracy in SQL query generation.
- Schema linking provided crucial contextual alignment between natural language and database metadata.
- Real-time user feedback served as a lightweight supervision tool, improving adaptability without the need for extensive labelled data.
- Compared to static, rule-based interfaces, the dynamic model architecture allowed scalability across multiple domains and schemas.
- The system architecture is modular and cloud-ready, making it suitable for integration into enterprise dashboards, educational tools, and public data portals.

## 5) CONCLUSION AND FUTURE WORK

### 5.1 Conclusion:

The Natural Language Interface for Relational Database Querying project is designed to make accessing data easier by connecting everyday language with the structured world of SQL queries. Using natural language processing (NLP) and semantic mapping, the system enables users to interact with databases simply by typing queries in plain English—removing the need for any specialized knowledge of query languages. This not only improves usability but also speeds up data retrieval in areas like business intelligence, education, and decision-making. The project highlights how combining NLP with relational databases can simplify complex querying tasks and make data tools more accessible to non-technical users. Testing results showed the system could reliably interpret and convert user input into accurate SQL commands, reinforcing the strength of its design. Overall, the project marks a meaningful step toward creating smarter, more inclusive, and user-friendly

database systems.

### 5.2 Future Work:

While the current system provides a strong foundation for interacting with relational databases through natural language, there are several potential areas for enhancement. One key improvement involves handling more complex queries. By integrating advanced NLP models like BERT or GPT, the system can be upgraded to interpret and generate SQL for more sophisticated queries involving nested conditions, multiple joins, groupings, and aggregations. Additionally, incorporating voice-enabled querying through speech- to-text technology can significantly improve accessibility, making the system usable in hands-free environments or by visually impaired individuals. Expanding the system's language capabilities to support multilingual queries would further enhance its inclusivity, allowing users from diverse linguistic backgrounds to interact with the database in their native language. Moreover, adaptive schema recognition could be introduced to allow the system to automatically detect and respond to changes in the database schema, reducing the need for manual reconfiguration. Integration with data visualization tools presents another promising direction, enabling users to view query results as interactive charts or dashboards, which can improve comprehension and decision-making. Lastly, embedding a feedback loop within the interface would allow the system to learn from user inputs and corrections, progressively enhancing its query translation accuracy and overall performance. These advancements would make the interface more intelligent, adaptive, and scalable in real-world applications.

### REFERENCES

[1] Zhong, V., Xiong, C., & Socher, R., "Seq2SQL: Generating Structured Queries from Natural Language using
Reinforcement Learning.", arXiv preprint arXiv:1709.00103. 2017.

[2] Yu, T., Zhang, Z., Yasunaga, M., Wang, D., Li, Z., & Radev, D."Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-toSQL Task." EMNLP

[3] Srinivas, S., & Kumar, A. R. A Speech-based Natural Language Interface for Databases. International Journal of Man-Machine Studies. 2002

[4] Raj, A., & Dinesh, D. Text-to-SQL Generation for Natural Language Interfaces to Databases Using Deep Learning: A Survey. International Journal of Computer Applications 2020