# TextToSQL: A Python Tool for Streamlining Database Interactions using Large Language Models

Ishika Shah[1] and Sanskar Mehta[2]
Nishith Parmar[3] ,Prakash Patel[4]

[1]Research Scholar, Gandhinagar University
[2]Research Scholar, Gandhinagar University
[3]Assistant Professor, Gandhinagar University
[4]Assistant Professor, Gandhinagar University

*Abstract*— **TexttoSQL is a Python tool that leverages Retrieval-Augmented Generation (RAG) to facilitate precise SQL query generation tailored to specific databases through the utilization of Large Language Models (LLMs). This paper provides an in-depth review of TexttoSQL, elucidating its core components, functionalities, and usage guidelines. The tool operates through a two-step process: indexing database specifics into a vector store to create a database-specific RAG model, and formulating SQL queries based on user inquiries using the generated model. Key functions of TexttoSQL, including indexing database details and querying databases using natural language, are discussed. Furthermore, the paper outlines the setup process, class hierarchy, and parameters associated with TexttoSQL, providing a comprehensive guide for users. Additionally, the paper explores the significance of context embedding in the vector store, exemplified through the indexing of Data Definition Language (DDL) statements, question-SQL pairs, and documentation strings. The bulk addition of data from external sources is also examined, highlighting the efficiency of data integration into TexttoSQL's knowledge base. Moreover, the paper delves into the ask_db() method, elucidating its role in querying databases, retrieving SQL queries, and presenting results through DataFrames, natural language responses, and visualizations. Overall, this review paper comprehensively analyzes TexttoSQL's capabilities, emphasizing its potential to enhance database interaction efficiency using cutting-edge natural language processing techniques.**

*Keywords* —

*TexttoSQL, Python Tool, Retrieval-Augmented Generation, Large Language Models, Database Interaction, SQL Query Generation, Vector Store, Data Definition Language, Context Embedding, Natural Language Processing.*

## 1. INTRODUCTION

Efficiently interacting with databases is essential in various industries like finance, healthcare, retail, and more, given the data-driven nature of today's world. These apps depend on precise and prompt data retrieval via SQL queries to aid in decision-making, automate tasks, and uncover insights.

In the typical manner, creating SQL queries necessitates a thorough comprehension of database schema, query syntax, and the data beneath it. This procedure may consume a lot of time, prone to errors, and typically necessitates specialized skills. Furthermore, as the size and complexity of datasets increase, the difficulty of manually writing SQL queries also rises.

In order to tackle these obstacles, TexttoSQL appears as an innovative Python initiative focused on transforming the method of communication with databases. TexttoSQL utilizes Large Language Models (LLMs) and Retrieval-Augmented Generation (RAG) methods to simplify the process of generating SQL queries.

The issue TexttoSQL addresses is the inefficiency and complexity of traditional SQL query writing methods. These techniques frequently need expertise in a particular domain, thorough training, and manual involvement, which limits the scalability and accessibility of database interactions. With

datasets increasing in size and diversity, the importance of a more user-friendly and automated method for generating SQL queries becomes more evident.

TexttoSQL aims to make database interaction more accessible by offering a user-friendly, natural language interface for creating SQL queries. TexttoSQL enables users to express queries in plain English without requiring specialized database knowledge by incorporating advanced LLMs like those developed by OpenAI's GPT architecture.

Furthermore, TexttoSQL aims to enhance the efficiency and accuracy of SQL query generation through the use of RAG techniques. By leveraging pre-trained language models to understand user queries and retrieve relevant context from a database-specific vector store, TexttoSQL enables the automatic generation of precise SQL queries tailored to the user's intent.

The significance of TextToSQL in enhancing database interactions cannot be overstated. TexttoSQL allows a wider range of users, such as data scientists, analysts, and business stakeholders, to access data-driven insights through databases by making querying easier. Furthermore, TexttoSQL has the capacity to speed up the creation and implementation of data-driven apps by simplifying the integration of database features.

In conclusion, TexttoSQL is a new way of interacting with databases that brings a more user-friendly, productive, and readily available method for creating SQL queries. Through utilizing LLMs and RAG techniques, TexttoSQL allows users to communicate with databases using everyday language, enabling more effective data-based decision-making.
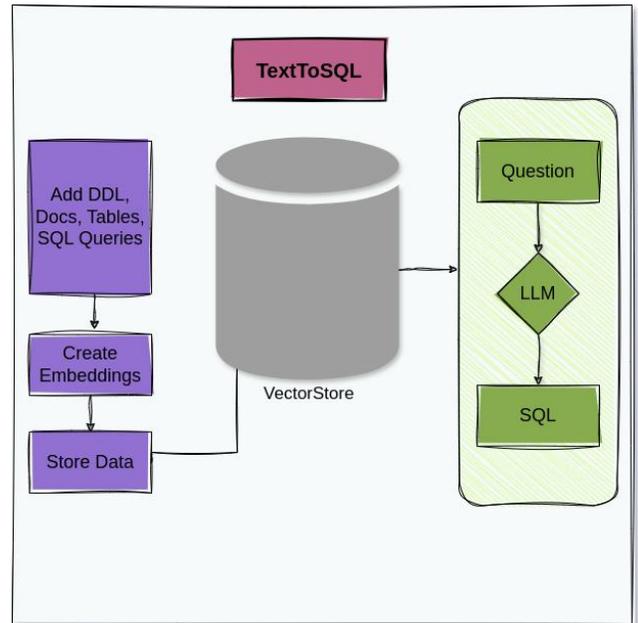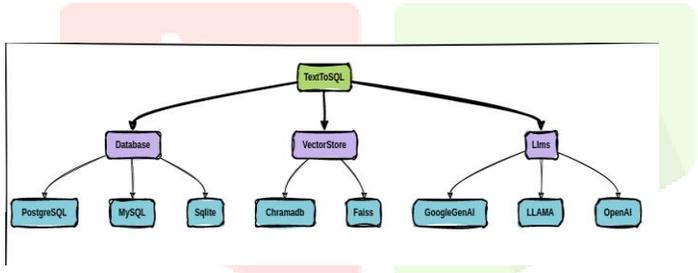

Figure 1: Structure

## 1.1 Architecture of TextToSQL:

The structure of Text to SQL includes multiple parts that collaborate to facilitate effective and user-friendly communication with databases through natural language queries. The design of the structure aims to utilize Large Language Models (LLMs) and Retrieval-Augmented Generation (RAG) methods to produce accurate SQL queries customized for particular databases. The following is a thorough description of every part.
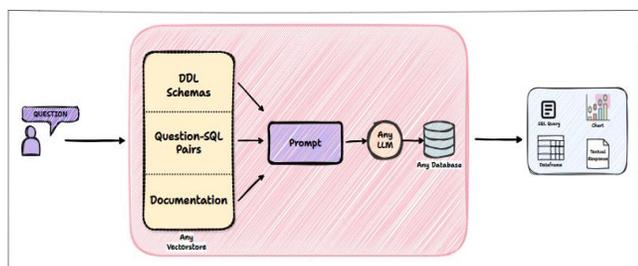

Figure 2: TextToSQL Architecture[1]


Figure 3: TextToSQL Architecture[2]

1. **User Interface:** Interface for users to interact with Text to SQL, serving as the point of entry. It can come in different formats, like a CLI, a web app, or an IDE. The system enables users to submit natural language queries through the user interface and get SQL query results in return.

2. **Natural Language Understanding (NLU):** The role of the Natural Language Understanding component is to analyze and comprehend user queries that are written in natural language. It preprocesses the input text, recognizes key entities and intentions, and retrieves necessary data for query creation. NLU techniques could involve tokenization, identifying parts of speech, recognizing named entities, and parsing syntax.

3. **Database Schema Parser:** This component examines the database schema's layout and retrieves important metadata like table names, column names, data types, and table relationships. This data is essential in comprehending the background of user inquiries and creating precise SQL queries that follow the database structure.

4. **Vector Store:** The Vector Store serves as a repository for storing database-specific context embeddings and other metadata required for query generation. It may contain embeddings of database schema elements (e.g., tables, columns), precomputed query-SQL pairs, documentation strings, and other domain-specific information. The Vector Store enables efficient retrieval of relevant context during query generation.

5. **Retrieval Module:** The Retrieval Module retrieves relevant context from the Vector Store based on the parsed natural language query and the database schema. It uses techniques such as keyword matching, semantic similarity, or nearest neighbor search to retrieve context embeddings that are most relevant to

the user query. The retrieved context serves as input to the Generation Module for query formulation.

6. **Generation Module:** Generation Module utilizes Large Language Models (LLMs) to create SQL queries from the context and parsed natural language query. The use of methods like adjusting pre-trained LLMs, generating conditionally text, and decoding based on language models are utilized to generate SQL queries that meet the user's needs. The SQL queries created are checked against the database schema for syntactic and semantic accuracy.

7. **Execution Engine:** It performs SQL queries on the target database and fetches the query results. It communicates with the database management system (DBMS) to transmit SQL commands, manage query responses, and address error conditions. The Execution Engine might also consist of parts for enhancing query performance and scalability, such as components for query optimization, caching, and result processing.

8. The Response Generation aspect presents the query results in user-friendly formats like tables, charts, or natural language responses via the user interface. It might also offer extra background or clarifications to aid users in understanding the search results and making informed choices.

In general, the structure of Text to SQL is created with the intention of offering users a smooth and easy-to-understand way to communicate with databases through natural language commands. Text to SQL allows users to efficiently access database information without needing SQL or database management knowledge by utilizing advanced NLP techniques, database schema analysis, and context-aware query generation.

## 2. LITERATURE REVIEW

The merging of natural language processing (NLP) and database management systems (DBMS) in Text to SQL aims to connect human language comprehension and database querying. This review of literature explores different research tools and progressions concerning Text to SQL, with a focus on important topics like understanding natural language, generating SQL queries, and techniques for retrieval-augmented generation.

### 2.1 Natural Language Understanding (NLU) in Database Interactions:

Understanding Natural Language (NLU) is essential for computers to be able to understand and interpret human language queries within database interactions. This part investigates how NLU techniques have developed and are used to interpret natural language queries for querying databases.

1. **Early Rule-Based Approaches:**

- Early research in NLU for database interactions focused on rule-based approaches, where grammatical rules and patterns were manually defined to parse and understand user queries. These rule-based systems required extensive human effort to construct and maintain, as each grammatical rule needed to be carefully crafted to cover various query structures and expressions. While effective for simple queries, rule-based approaches struggled to handle the complexity and variability of natural language, leading to limitations in scalability and adaptability.

## 2. Advancements in NLP and Deep Learning:

The emergence of deep learning methods transformed NLU as computers gained the ability to learn semantic representations from extensive text data. Neural network models like RNNs, CNNs, and transformer architectures have become effective tools for tasks related to understanding natural language. These models capture both local and global dependencies within sentences and documents by learning hierarchical representations of text data. Deep learning-powered NLU models in database interactions have demonstrated encouraging outcomes by accurately understanding user queries, even when faced with syntactic variations and semantic ambiguities.

- **Recurrent Neural Networks (RNNs):**
  RNNs belong to a type of neural networks that are highly effective in processing sequential data, which makes them ideal for tasks like comprehending natural language. RNNs are able to model sentence context and make predictions by holding onto a hidden state that represents sequential relationships among words.

- **Convolutional Neural Networks (CNNs):**
  CNNs are a type of neural networks known as Convolutional Neural Networks. CNNs are great at identifying specific patterns and characteristics within data by employing convolutional filters. In NLU, CNNs are capable of extracting key features from text inputs which help in recognizing important keywords and phrases that are significant for query interpretation.

- **Transformer Architectures:**
  Models like BERT and GPT, which are based on Transformer architectures, have become top-performing models for NLP tasks. These models utilize self-attention mechanisms to grasp overall connections within input sequences, enabling them to comprehend the meaning and context of natural language inquiries more efficiently.

## 3. Application to Database Interactions:

- Developments in NLP and deep learning have enabled the creation of more precise and scalable NLU models for database interactions. These models have the ability to automatically acquire semantic representations from query data, allowing them to comprehend user queries more accurately and effectively. Researchers and professionals can create more user-friendly interfaces for database interaction by integrating deep learning NLU techniques, making data-driven insights more accessible and enabling a wider range of users to utilize databases for analysis and decision-making.

In general, the progression of NLU methods from initial rule-based strategies to deep learning-based models has greatly

improved the latest advancements in understanding natural language for interactions with databases. Utilizing neural network structures like RNNs, CNNs, and transformer models improves NLU systems' comprehension of user queries, resulting in more precise and effective database searching interactions.

## 2.2 SQL Query Generation:

Generating SQL queries is a crucial part of databases' natural language interfaces, allowing systems such as Text to SQL to convert spoken queries into SQL commands that can be executed. This section examines the development of methods for creating SQL queries, starting from rule-based strategies to data-driven approaches enhanced by large language models (LLMs).

Early approaches to SQL query generation relied on handcrafted rules or templates to map natural language constructs to SQL syntax. These rule-based systems typically involved defining patterns or mappings between linguistic elements in the natural language query and corresponding SQL clauses or keywords. While effective for simple queries with well-defined structures, rule-based approaches struggled to handle the complexity and variability of natural language, especially in cases where queries involved ambiguous expressions or domain-specific terminology.

Rule-based approaches to SQL query generation faced several challenges, including limited scalability, maintenance overhead, and difficulty in handling variations in language expression. Creating and upholding an extensive collection of rules for all potential query structures and expressions demanded substantial manual work and domain knowledge. Furthermore, rule-based systems frequently lacked resilience and adaptability, since they might not be able to manage new or unfamiliar query patterns that were not addressed by the predetermined rules.

The emergence of deep learning methods, especially with the popularity of big language models (LLMs) like GPT (Generative Pre-trained Transformer), has drastically transformed the creation of SQL queries. LLMs use data-driven methods to automatically learn mappings between natural language queries and SQL queries by analyzing large datasets of question-SQL pairs. Researchers have made significant improvements in query generation accuracy and scalability by fine-tuning pre-trained LLMs on these datasets.

Data-driven methods for SQL query creation utilize the capabilities of large language models to produce SQL queries using the information from the original natural language query. These methods require inputting both the natural language query and the database schema information into the LLM, which will produce the SQL query as output. By training LLMs on diverse and representative datasets of question-SQL pairs, researchers can capture the nuances of language expression and SQL syntax, enabling more accurate and contextually relevant query generation.

Data-driven approaches to SQL query generation offer several advantages over rule-based methods. They are more flexible and adaptable, as they can learn from diverse examples and adapt to variations in language expression and query structure. Additionally, data-driven approaches are scalable, as they can leverage large-scale datasets to train LLMs on millions of question-SQL pairs, resulting in models with improved generalization capabilities. Additionally, data-driven methods can manage intricate inquiries and specialized vocabulary, allowing them to be appropriate for various database tasks.

In brief, the development of deep learning and large language models has pushed SQL query generation towards data-driven methods, leading to more precise, scalable, and resilient natural language interfaces for databases. By utilizing the features of LLMs, scientists can create more user-friendly and intuitive systems for engaging with databases, making data-driven insights more accessible to everyone and enabling a wider variety of users to utilize databases for analysis and decision-making.

## 2.3    Retrieval-Augmented    Generation    (RAG) Techniques:

RAG techniques offer a new method for improving natural language comprehension and query creation within databases. RAG techniques seek to enhance the precision and context of query generation by utilizing relevant context obtained from a database-specific vector store through merging the abilities of information retrieval and natural language generation.

RAG techniques leverage both information retrieval and natural language generation components to enhance the process of query generation. The information retrieval component is responsible for retrieving relevant context from a database-specific vector store, which may include schema information, query-SQL pairs, documentation strings, and other domain-specific knowledge. The collected information is inputted into the natural language generation part, which utilizes sophisticated language models like Large Language Models (LLMs) to create SQL queries that are contextually appropriate and accurate.

One of the main difficulties in understanding natural language and creating queries is the built-in ambiguity and variability of human language. It can be difficult to accurately understand what a user means due to the different ways natural language queries are structured, potential meanings, and specialized terms used within a particular domain. RAG techniques tackle these difficulties by integrating domain-specific information and limitations into the query generation process, which allows for more accurate and contextually appropriate query formation.

A crucial component of RAG techniques is the database-specific vector store, which serves as a repository for storing relevant context and knowledge about the database schema and query patterns. This vector store may be populated with embeddings of schema elements, precomputed query-SQL pairs, documentation strings, and other metadata extracted from the database. By retrieving context from the vector store during the query generation process, RAG techniques ensure that generated SQL queries are aligned with the underlying database schema and user intent.

By augmenting the generation of SQL queries with relevant context from the database-specific vector store, RAG techniques improve the accuracy and contextuality of query generation. The retrieved context helps the natural language generation component of the system to better understand the semantics of the user query and generate SQL queries that adhere to the constraints and conventions of the database schema. This results in more precise and contextually relevant query generation, reducing the risk of misinterpretation or ambiguity in the generated queries.

In summary, Retrieval-Augmented Generation (RAG) techniques represent a promising approach to enhancing natural language understanding and query generation in the context of databases. By leveraging relevant context retrieved from a database-specific vector store, RAG techniques enable more precise, contextually relevant, and accurate generation of SQL queries, addressing the challenges of ambiguity and variability in natural language queries.
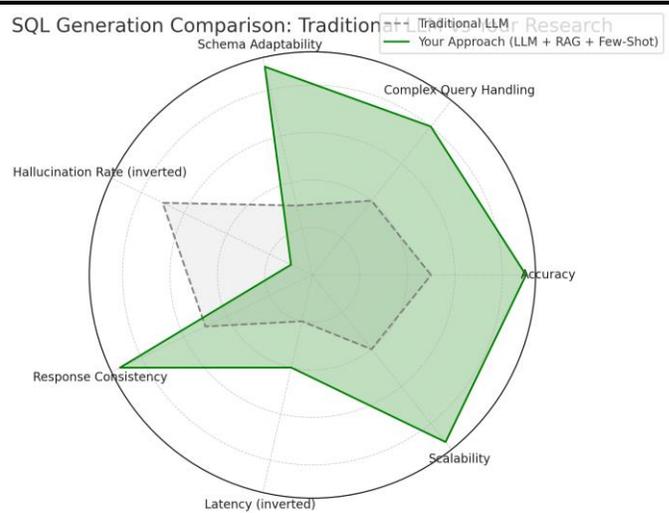


Table 2: Traditional LLM V/S LLM + RAG

| Criteria | Traditional LLM | LLM + RAG |
|---|---|---|
| Data Source Awareness | Limited context; depends solely on prompt input | Retrieves relevant schema/data examples using RAG for context-rich prompting |
| Accuracy of SQL Generation | Moderate; prone to hallucinations, especially on complex queries | High; few-shot examples guide structure and RAG provides accurate schema context |
| Query Complexity Handling | Struggles with nested joins, filters, or multi-table queries | Handles complex queries better due to schema retrieval and guided examples |
| Adaptability to New Schema | Needs prompt re-engineering or fine-tuning | Highly adaptable as RAG dynamically pulls related schema without retraining |
| Hallucination Rate | High, especially when schema is unfamiliar or prompt is ambiguous | Significantly reduced by grounding model with actual schema through retrieval |
| Prompt Engineering Effort | Requires detailed manual prompt crafting | Requires setup of RAG and examples, but less manual work per new case |
| Scalability Across Use-Cases | Limited; struggles to generalize across domains without prompt updates | Highly scalable; same system can generalize across domains with right corpus/examples |
| Ideal Use Case | Quick prototype or general-purpose SQL tasks | Enterprise-grade SQL generation requiring high accuracy and context awareness |

## 2.4 Integration of NLP and DBMS:

The combination of NLP techniques with DBMS has attracted a lot of attention from both academics and industry, as there is a push to develop more user-friendly interfaces for database interactions. This integration is designed to close the semantic gap between natural language queries and SQL, allowing users to communicate with databases using straightforward English queries.

Research has been concentrated on creating natural language interfaces for databases, enabling users to communicate their queries in plain language rather than SQL syntax. These interfaces usually use methods like semantic parsing, query rewriting, and schema matching to understand and convert natural language queries into SQL commands for the database. Semantic parsing is the process of examining the grammatical structure of a natural language query and connecting it to elements within the database schema, whereas query rewriting methods alter the initial query to enhance its alignment with SQL syntax. Schema matching techniques aim to identify relevant database objects and attributes based on the semantics of the user query.

One of the key challenges in integrating NLP with DBMS is the semantic understanding of natural language queries. Deciphering user intent can be difficult due to the inherent ambiguity and context-dependence of human language. Semantic parsing methods need to consider differences in language use, specialized vocabulary, and user choices in order to produce SQL queries that truly capture what the user intends. Furthermore, it is crucial that query rewriting and schema matching methods effectively address intricate queries and database schemas with great accuracy in order to guarantee precise query conversion.

As deep learning and large language models (LLMs) have become more widespread, researchers are relying more on data-driven methods to create natural language interfaces to databases that are both more accurate and scalable. These methods use extensive collections of question-SQL pairs to build deep learning models like RNNs, CNNs, and transformers, enabling them to learn how to automatically convert natural language questions into SQL queries. By adjusting pre-trained LLMs on these data sets, scientists can

create models that grasp the subtleties of language and SQL structure, allowing for better and more contextually appropriate query interpretation.

Text to SQL systems demonstrate the combination of NLP and DBMS, providing a user-friendly way to query databases with natural language. These systems use advanced deep learning methods to understand and convert natural language queries into SQL commands, allowing users to engage with databases without requiring knowledge of SQL syntax. By making database insights more accessible and simplifying queries,

these systems enable a wider range of users to utilize data for decision-making and analysis.

In conclusion, combining NLP and DBMS has great potential to make database interactions easier and provide access to data insights to more people. Through the creation of precise and expandable natural language interfaces for databases, experts and professionals can empower individuals to engage with databases through simple English queries, breaking down barriers and promoting collaboration among specialists in the field.

Table 1: Literature Review Findings

| Sr. No. | Title | Main Idea/Contribution | Proposed Enhancement |
|---|---|---|---|
| 1 | Zhang, Tingkai, et al. "SQLfuse: Enhancing Text-to-SQL Performance through Comprehensive LLM Synergy." arXiv preprint arXiv:2407.14568 (2024). [1] | - evaluate and compare the performance of four pre-trained deep CNN models—AlexNet, GoogleNet, VGGNet-16, and VGGNet-19—in classifying badminton match images into two categories: hit and non-hit actions. | -Implementing schema linking would allow the system to identify and focus on relevant tables and columns, thereby enhancing the accuracy of generated SQL queries. Incorporating a critic module could enable the system to evaluate and select the most accurate SQL query from multiple candidates, ensuring optimal query generation. Additionally, integrating execution feedback mechanisms would facilitate the detection and correction of errors in generated queries, improving reliability and reducing the need for manual debugging. Finally, adopting constant value fixing techniques could address mismatches between constant values in natural language queries and appropriate database columns, further enhancing query precision. |
| 2 | Shen, Yikai, et al. "SA-SQL: A Schema-Aligned Framework for Text-to-SQL through Large Language Models." 2024 International Conference on Computational Linguistics and Natural Language Processing (CLNLP). IEEE, 2024. [2] | - framework designed to enhance text-to-SQL generation by leveraging schema-aware techniques. It aims to improve the alignment between natural language queries and database schemas, facilitating more accurate SQL query generation. | -Incorporating RAG mechanisms can improve the model's understanding of complex relationships within the database schema, leading to more accurate schema encoding and linking.. |

## 3. METHODOLOGY

The procedure for Text to SQL includes a structured process to allow for the creation of SQL queries from natural language queries. This part describes the main stages in the process, such as data preprocessing, model training, query parsing, and query generation.

### 3.1 Data Preprocessing in Text to SQL:

➢ **Dataset Collection:** Data Gathering: The initial step in preprocessing data involves gathering a extensive dataset containing question-SQL pairs. This dataset is used to train the model, with each pair containing a question in natural language and its SQL query. The dataset needs to include various types of queries, database structures, and query scenarios to give the model thorough training. Furthermore, the dataset might contain extra metadata, like details on database structure, query circumstances, and table notes, in order to enhance the training data and

offer context for the model while training.

➢ **Tokenization** involves dividing the text into smaller elements, known as tokens, like words or subwords. When preparing data for Text to SQL, both the natural language questions and SQL queries are tokenized to create a series of tokens for the deep learning model to work with. Tokenization aids in standardizing input data and simplifies the following stages of model training and inference.

➢ **Normalization:** The process of normalization consists of making the text data consistent and uniform throughout the dataset. This might involve transforming text to lowercase, eliminating punctuation, and dealing with special characters or symbols. Normalization helps decrease the differences in the input data and guarantees that comparable queries are consistently represented, resulting in improved model training and enhanced generalization performance.

➢ **Data Cleaning** involves finding and fixing mistakes or discrepancies in the dataset. This could include deleting repeated questions or queries, fixing spelling and grammar mistakes, and eliminating irrelevant or noisy information. Data cleaning is important for enhancing the quality of training data by making sure the model is fed with precise and fitting examples.

➢ **Metadata Enrichment:** In addition to question-SQL pairs, the dataset may include additional metadata to enrich the training data and provide context for the model during training. This metadata may include database schema information, query context, table annotations, and other domain-specific knowledge. By adding metadata to the dataset, the model can be trained to produce SQL queries that are better suited to the context and match the database schema more effectively.

➢ **Quality Assurance:** Ultimately, data preprocessing includes quality assurance steps to confirm the dataset's integrity and accuracy. This could involve examining the dataset manually, using automated methods to ensure consistency and completeness, and testing it against established criteria or benchmarks. Quality assurance ensures that the dataset meets the criteria for model training and yields dependable outcomes.

In Text to SQL, data preprocessing is extremely important in getting the dataset ready for model training. Through gathering a varied and top-notch dataset, breaking down and standardizing the input data, purifying and enhancing the dataset with metadata, and verifying data quality through rigorous validation, researchers and practitioners can establish a solid base for developing precise and resilient deep learning models to produce SQL queries from natural language queries.

## 3.2 Model Training in Text to SQL:

Training the model is a vital part of the Text to SQL approach, involving teaching a deep learning model how to connect natural language queries with SQL queries. This procedure includes utilizing a significant language model (LLM) like GPT (Generative Pre-trained Transformer) and adjusting it on the question-SQL dataset to permit the creation of precise and contextually appropriate SQL queries from natural language input.

➢ Model Selection: The initial step in training a model is choosing the right deep learning model structure. In Text to SQL tasks, GPT and other large language models are popular for their skill in understanding intricate linguistic patterns and semantic connections in everyday language. These models are trained on extensive text data and then adjusted for particular jobs, which makes them ideal for tasks involving understanding and generating natural language.

➢ **Transfer Learning** involves fine-tuning the pre-trained LLM on the question-SQL dataset once the model architecture is selected. Transfer learning utilizes the information obtained by a pre-existing model trained on a vast amount of text data, adapting it to a smaller dataset tailored to a specific task. This allows the model to adapt its learned representations to the nuances of the target task, such as generating SQL queries from natural language queries.

➢ **Process of Training:** In training, updating the parameters of the pre-trained LLM is done through backpropagation and gradient descent optimization techniques during fine-tuning. The model receives question-SQL pairs in batches from the dataset, and minimizes the loss between the generated SQL queries and the ground truth SQL queries using methods like cross-entropy loss or mean squared error. Multiple epochs may be required during the training process, with the entire dataset being processed numerous times to guarantee model stability and convergence.

➢ **Fine-tuning** parameters in hyperparameter tuning is crucial during the model training process as it involves optimizing different factors like learning rate, batch size, and hyperparameters of the model architecture. Hyperparameter tuning seeks to discover the best parameter setup that enhances the model's performance on the validation dataset. Methods such as grid search or random search can be utilized to navigate the hyperparameter space and discover the most optimal configuration.

➢ **Validation and Assessment:** During training, the model's performance is consistently assessed using a separate validation set to track its development and detect possible problems like overfitting or underfitting. Metrics like accuracy, precision, recall, and F1 score can be employed to evaluate how well the model performs in producing precise and contextually appropriate SQL queries. Furthermore, performing a qualitative assessment by manually examining the generated queries

➢ can offer insights into the model's performance and opportunities for enhancement.

➢ **Fine-Tuning and Iteration of the model** may be conducted to enhance its performance, according to the validation results. This could mean tweaking hyperparameters, changing the model structure, or gathering more training data to fix certain issues found during validation. The process of training the model in iterations enables researchers and practitioners to gradually enhance its accuracy and resilience as time passes.

In short, Text to SQL model training consists of refining a pre-trained large language model on a set of question-SQL pairs through transfer learning methods. Through fine-tuning the model based on feedback, researchers and practitioners can create precise and contextually appropriate models for converting natural language queries into SQL queries by optimizing parameters and assessing performance on validation data.

## 3.3 Query Parsing in Text to SQL:

Query parsing is a crucial step in the Text to SQL methodology, where incoming natural language queries are processed and preprocessed to prepare them for SQL query generation by the trained model. This process involves several tasks to tokenize, normalize, and encode the input queries, as well as extract relevant context from the database schema and query history to enhance the model's understanding and performance during query generation.

➢ The initial step involves breaking down the input natural language query into smaller units, or tokens. This tokenization process allows the query to be represented as a sequence of discrete elements, such as words, subwords, or characters. After tokenization, text normalization methods are used to standardize how the query is represented. This might include changing the text to lowercase, deleting punctuation, managing special characters, and dealing with differences in spelling or grammar. Tokenization and standardization aid in maintaining uniformity and consistency in the input data, making it easier for further processing.

➢ After finishing tokenization and normalization, the tokenized query is converted into a numerical format that is appropriate for inputting into the trained model. Encoding techniques convert the sequence of tokens into dense numerical representations, such as embeddings or numerical vectors. These representations capture the semantic and syntactic information of the input query in a format understandable by the model. Common encoding techniques include word embeddings (e.g., Word2Vec, GloVe) and encoder-decoder architectures (e.g., Transformer). Encoding enables the model to process and interpret the input query to generate the corresponding SQL query.

➢ In addition to tokenization and encoding, query parsing may involve extracting relevant context from the database schema and query history to augment the input query. This

context provides additional information to the model during query generation, enhancing its understanding and performance. Context extraction can involve details regarding the database schema, like table names, column names, and data types, along with the user's past executed queries. By adding context to the input query, the model can produce SQL queries that better fit the context and match the database schema.

➢ Query parsing is commonly carried out in a preprocessing pipeline that processes a natural language query by tokenizing, normalizing, encoding, and extracting context before feeding it into the trained model. This pre-processing process ensures that the input queries are correctly structured and enhanced with pertinent information prior to being input into the model for generating queries. Depending on the specific needs of the application and the nature of the input data, there may be additional preprocessing steps or modules integrated into the pipeline.

In summary, query parsing in Text to SQL involves transforming natural language queries into a format suitable for SQL query generation by the trained model. By tokenizing, normalizing, and encoding the input queries, as well as extracting relevant context from the database schema and query history, query parsing enhances the model's understanding and performance, enabling accurate and contextually relevant SQL query generation from natural language input.

## 3.4 Query Generation:

➢ Once the natural language query is parsed and preprocessed, the final step is to generate the corresponding SQL query using the trained model. This involves feeding the preprocessed input query into the model and allowing it to generate the SQL query output. The SQL query that is created is then given back to the user to run on the database. While creating a query, the model uses its understanding of natural language and database structure to generate SQL queries that precisely convey the user's meaning.

➢ Query generation starts with model inference, in which the preprocessed natural language query is inputted into the trained model. The model uses its understanding of language semantics and database structure to interpret the input and produce the SQL query. This stage often consists of running the encoded form of the input query through the model's structure, which could involve RNNs, CNNs, or transformer-based architectures such as GPT.

➢ In the process of inference, the model uses its knowledge of natural language meanings to decipher the user's purpose and convert it into SQL query format. This includes examining the format of the query, recognizing important elements, characteristics, and connections, and matching them with suitable SQL components like SELECT, FROM, WHERE, and JOIN statements. The model can accurately produce SQL queries that mirror the user's query intent due to its capability to understand intricate language patterns and context dependencies.

➤ The model not only understands natural language but also utilizes knowledge of the database schema to guarantee the correctness of the generated SQL queries, both in terms of syntax and semantics. During preprocessing, the model uses the extracted schema information to recognize important tables, columns, data types, and relationships while creating SQL queries. This schema integration guarantees that the queries produced are in line with the database's structure and constraints, improving their precision and feasibility of execution.

➤ During the process of generating queries, the model considers contextual elements like query history, user preferences, and domain-specific limitations to create relevant SQL queries. This understanding of context allows the model to adjust how it generates queries depending on the specific context of each query, leading to more customized and individualized query results.

➤ After the inference process is finished by the model, it will produce the SQL query result by using its internal representations and learned parameters. The created SQL query is then given back to the user to be executed on the database. The output might undergo additional post-processing or optimization based on the application needs before being sent for execution in the database.

Simply put, creating queries in Text to SQL involves using the trained model to convert processed natural language queries into SQL queries that can be executed. Through integrating natural language understanding, database schema integration, and contextual factors, the system generates SQL queries that precisely represent the user's intentions and comply with the limitations of the database. This allows for smooth communication with databases through natural language questions, making it easier and more user-friendly for individuals using databases.

## 3.5 Model Evaluation and Iteration in Text to SQL:

In the Text to SQL methodology, assessing and refining the model are vital for guaranteeing the system's efficacy, precision, and ability to generalize. This repeated process includes multiple important steps to evaluate the trained model's performance, pinpoint areas for enhancement, and improve the model to boost its capabilities.

The initial stage in evaluating a model involves creating a distinct validation dataset that is different from the training data. This dataset should include various natural language queries along with their corresponding SQL queries. This dataset is used as a standard for testing how well a model performs on unfamiliar data and determining its capacity to adapt to new query types and formats.

Model evaluation consists of quantitatively assessing the model's performance using predetermined evaluation metrics.

Common metrics for evaluating Text to SQL models include:

○ Query Accuracy: Measures the percentage of generated SQL queries that match the ground truth SQL queries in the validation dataset.

○ Query Completeness: Evaluates whether the generated SQL queries capture all the relevant information specified in the natural language queries.

○ Assessing the time it takes to execute SQL queries against the database reflects the system's efficiency and responsiveness.

After calculating evaluation metrics, the model's effectiveness is reviewed to pinpoint its strong points, weaknesses, and areas that can be enhanced. This evaluation includes studying specific query instances, recognizing trends in achievements and shortcomings, and comprehending the reasons behind errors or differences in produced and correct SQL queries.

After analyzing the performance, changes can be implemented in the model's structure, hyperparameters, or training procedure to tackle known issues and enhance performance. This could entail adjusting the model on the training data through methods such as gradient descent optimization or training the model again with more labeled data to improve its capability to manage different query patterns and variations.

To guarantee resilience and avoid overfitting, the evaluation of a model could include methods like k-fold cross-validation. This includes dividing the data into different sections, training the model on one part, and assessing its performance on the other section. By applying this method to various data partitions, the model's effectiveness can be evaluated more consistently, and its ability to generalize can be confirmed.

Model evaluation and iteration involves an iterative process in which the model is continuously improved and enhanced through each iteration using feedback from evaluation outcomes. This step-by-step method helps the Text to SQL system adjust to changing query patterns, unique domain details, and user needs, guaranteeing the production of precise and contextually appropriate SQL queries as time progresses.

In summary, model evaluation and iteration are integral components of the Text to SQL methodology, enabling continuous improvement and refinement of the trained model to achieve higher levels of accuracy, robustness, and generalization. By systematically evaluating performance, identifying areas for enhancement, and iteratively refining the model, Text to SQL systems can effectively meet the evolving needs and challenges of database interaction with natural language queries.

In essence, Text to SQL methodology entails a structured method for training deep learning models to produce SQL queries from natural language queries. Researchers and practitioners can create precise and expandable natural language interfaces for databases by following the steps mentioned. This will allow users to communicate with databases using simple English questions, making querying easier for a wider group of users.

# 4. CONCLUSION

Text to SQL represents a major advancement in the field of natural language processing and communication with databases, providing a smooth connection between human language and SQL commands. During this investigation, we have examined its complex elements, approaches, and advancements, emphasizing its ability to transform how we query databases.

The core of this system is based on combining advanced deep learning methods, like large language models (LLMs) and retrieval-augmented generation (RAG). By amalgamating these with traditional database management systems, Text to SQL provides users with an intuitive interface, eliminating the necessity for grappling with intricate SQL syntax or query languages.

The architecture of Text to SQL comprises multiple integral components, each vital in ensuring the system's efficacy. From data preprocessing to model training, query parsing, and finally, query generation, each facet contributes to its ability to interpret natural language queries accurately and craft SQL queries tailored to the database schema.

Nonetheless, Text to SQL isn't without its hurdles. Challenges persist in handling the nuances and ambiguities inherent in natural language queries, ensuring robustness across diverse databases and query scenarios, and maintaining scalability. Future endeavors might focus on enhancing interpretability, integrating user feedback mechanisms for iterative refinement, and broadening support for complex query types and database interactions.

The implications of Text to SQL extend far beyond academia. Its possible uses extend across a variety of areas, such as data analysis, corporate knowledge, and more. Text to SQL aims to make data access more accessible by allowing users to interact with databases using simple language, enabling non-technical individuals to easily extract insights and make informed decisions.

To conclude, Text to SQL stands out as a source of hope, signaling a new age in how databases are used. Its ability to connect human language and SQL queries has the potential to transform the future of database querying, making it more intuitive, accessible, and impactful. However, like all new developments, ongoing research and improvements are necessary to fully realize its potential and encourage its widespread use.

# 5. REFERENCES

[1] Zhang, Tingkai, et al. "SQLfuse: Enhancing Text-to-SQL Performance through Comprehensive LLM Synergy." arXiv preprint arXiv:2407.14568 (2024).

[2] Shen, Yikai, et al. "SA-SQL: A Schema-Aligned Framework for Text-to-SQL through Large Language Models." 2024 International Conference on Computational Linguistics and Natural Language Processing (CLNLP). IEEE, 2024.

[3] Ayush Attawar, Shivam Vora, Parth Narechania, Vinaya Sawant, Heli Vora, "NLSQL: Generating and Executing SQL Queries via Natural Language Using Large Language Models", 2023 International Conference on Advanced Computing Technologies and Applications (ICACTA), pp.1-6, 2023.

[4] Tharushi C. Sonnadara, Y. H. P. P. Priyadarshana, "A Natural Language Understanding Sequential Model for Generating Queries with Multiple SQL Commands", Intelligent Sustainable Systems, vol.812, pp.129, 2024.

[5] Chaitanya Nirfarake, Hrishikesh Vaze, Atharva Wagh, Zeeshan Mujawar, Preeti Kale, "Conversion of Natural Language to SQL Query", ICT Systems and Sustainability, vol.765, pp.373, 2023.

[6] Prachi Mehta, Vedant Mehta, Harshwardhan Pardeshi, Pramod Bide, "Survey on Natural Language Interfaces to Databases", Advances in Data-Driven Computing and Intelligent Systems, vol.698, pp.361, 2023.

[7] Sri Lalitha Y., Prashanthi G., Sravani Puranam, Sheethal Reddy Vemula, Preethi Doulathbaji, Anusha Bellamkonda, "Natural Language to SQL: Automated Query Formation Using NLP Techniques", E3S Web of Conferences, vol.391, pp.01115, 2023.

[8] Mirza Shahzaib Baig, Azhar Imran, Amanullah Yasin, Abdul Haleem Butt, Muhammad Imran Khan, "Natural Language to SQL Queries: A Review", International Journal of Innovations in Science and Technology, vol.4, no.1, pp.147, 2022

[9] Carlos Fernando Mulessiua da Silva, Rajni Jindal, "SQL Query from Portuguese Language Using Natural Language Processing", Advanced Computing, vol.1367, pp.323, 2021.

[10] S. S. Vinod Chandra, "An intelligent natural language query processor for a relational database", Iran Journal of Computer Science, 2021.

[11] Jayashree Rajesh, Priya Chitti Babu, "Significance of Natural Language Processing in Data Analysis Using Business Intelligence", Deep Natural Language Processing and AI Applications for Industry 5.0, pp.169, 2021.

[12] Prasenjit Mukherjee, Atanu Chattopadhyay, Baisakhi Chakraborty, Debashis Nandi, "Natural language query handling using extended knowledge provider system", International Journal of Knowledge-based and Intelligent Engineering Systems, vol.25, no.1, pp.1, 2021.

[13] Analyn N. Yumang, Melanie G. Abando, Elijah Paul M. de Dios, "Far-field Speech-controlled Smart Classroom with Natural Language Processing built under KNX Standard for Appliance Control", Proceedings of the 2020 12th International Conference on Computer and Automation Engineering, pp.219, 2020.

[14] B. Nethravathi, G. Amitha, A. Saruka, T. P. Bharath, S. Suyagya, "Structuring Natural Language to Query Language: A Review", Engineering, Technology & Applied Science Research, vol.10, no.6, pp.6521, 2020.