





MAVIDSQL: A Model-Agnostic Visualization for Interpretation and Diagnosis of Text-to-SQL Tasks

Jingwei Tang , Guodao Sun , Jiahui Chen , Gefei Zhang , Baofeng Chang ,
Haixia Wang , *Member, IEEE*, and Ronghua Liang , *Senior Member, IEEE*

Abstract—Significant advancements in semantic parsing for text-to-SQL (T2S) tasks have been achieved through the employment of neural network models, such as LSTM, BERT, and T5. The exceptional performance of large language models, such as ChatGPT, has been demonstrated in recent research, even in zero-shot scenarios. However, the inherent transparency of T2S models presents them as black boxes, concealing their inner workings from both developers and users, which complicates the diagnosis of potential error patterns. Despite the fact that numerous visual analysis studies have been conducted in natural language processing communities, scant attention has been paid to addressing the challenges of semantic parsing, specifically in T2S tasks. This limitation hinders the development of effective tools for model optimization and evaluation. This article presents an interactive visual analysis tool, MAVIDSQL, to assist model developers and users in understanding and diagnosing T2S tasks. The system comprises three modules: the model manager, the feature extractor, and the visualization interface, which adopt a model-agnostic approach to diagnose potential errors and infer model decisions by analyzing input-output data, facilitating interactive visual analysis to identify error patterns and assess model performance. Two case studies and interviews with domain experts demonstrate the effectiveness of MAVIDSQL in facilitating the understanding of T2S tasks and identifying potential errors.

Index Terms—Error diagnosis, information visualization, text-to-SQL (T2S), visual analytics.

I. INTRODUCTION

TEXT-TO-SQL (T2S) tasks is a crucial subtask in the semantic parsing of natural language processing (NLP), as it involves mapping natural language (NL) utterances to structured query language (SQL) that can be executed on a relational database. These techniques bridge the semantic gap between a NL and database [1]. It can benefit various applications, such

Manuscript received 28 April 2023; revised 21 December 2023 and 6 March 2024; accepted 11 April 2024. Date of publication 18 April 2024; date of current version 15 October 2024. This work was supported in part by the Zhejiang Provincial Natural Science Foundation of China under Grant LR23F020003 and Grant LTGG23F020005, in part by the National Natural Science Foundation of China under Grant 62372411 and Grant 62036009, and in part by the Fundamental Research Funds for the Provincial Universities of Zhejiang under Grant RF-B2023006. (Corresponding author: Guodao Sun.)

The authors are with the College of Computer Science and Technology, Zhejiang University of Technology, Hangzhou 310023, China (e-mail: jwutang@zjut.edu.cn; guodao@zjut.edu.cn; chenjiahui@zjut.edu.cn; gefei@zjut.edu.cn; baofeng.chang@zjut.edu.cn; hxiwang@zjut.edu.cn; rhliang@zjut.edu.cn).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TCDS.2024.3391278>, provided by the authors.

Digital Object Identifier 10.1109/TCDS.2024.3391278

as enabling nonexperts to access data query and facilitating the development of human-computer intelligent interaction. In recent years, T2S tasks have embraced a new breakthrough with the Transformer architecture [2] and large language models (LLMs) [3]. As a result of pretraining and fine-tuning techniques, an increasing number of T2S models have demonstrated remarkable performance.

During the development of T2S models, NLP scientists are confronted with a series of challenges. First, deep learning models have complex internal mechanisms, making it difficult for developers to explain why and how the model derives a particular decision or prediction. Second, the evaluation metrics of the model do not possess the capacity to offer a comprehensive diagnostic assessment for the model. The evaluation metrics for the T2S task provide an overall accuracy score for a model, making it challenging to pinpoint specific areas where errors are likely to occur. Model developers are required to examine SQL statements individually to identify error patterns of the model. Extensive textual data present challenges in the analysis. Data scientists must possess an in-depth understanding of the context surrounding each question and SQL statement to make informed judgments regarding specific instances. Both of these challenges require more effective approaches that enable developers to interactively explore insights from model results and iterate novel models.

Visualization techniques have been employed to assist model developers in comprehending and enhancing deep learning models [4]. In the realm of NLP tasks, model developers confront challenges such as the unstructured nature of NL and its semantic diversity. Various visualization techniques have been proposed to facilitate the development of various deep learning based language models, such as LSTMVis [5], Seq2Seq-Vis [6], and RNNVis [7]. The incorporation of visual analysis can aid in the better understanding of the decision-making process and outcomes of reasoning models, assisting users in analyzing and comprehending the application scenarios and limitations of such models. However, applying existing techniques to T2S tasks poses challenges because of their exclusive design for specific neural network models, which might not be appropriate for diverse T2S parsers that adopt varying neural network frameworks.

In this article, we propose MAVIDSQL, a novel model-agnostic visual analysis tool that aids developers and users in comprehending and diagnosing for T2S tasks. The system does not require access to the internal logic of the model and only

relies on input instances and output results. This design enables it to support a wide range of model types, as long as they target the same machine learning task and have a consistent input–output format. We conducted an extensive review of the literature in the field to identify the design requirements of our system.

Inspired by the NLP community, we employ a sentence similarity comparison method based on semantic role labeling and syntactic dependency parsing to analyze the impact of input sentence semantics on prediction results. The projected view is generated using similarity measures, enabling users to filter input questions of interest interactively and explore the relationship between input questions and model error patterns at a global-class level in conjunction with model performance statics view. Through the control panel, users have the capability to select datasets and inspect databases for analysis. They may also filter specific attributes to assess the complexity of the model. To enhance the efficiency of analyzing large-scale SQL predictions by users, we employ a modeling technique to identify the differences between predicted and ground-truth SQL queries. These disparities are presented in the SQL comparison view, which users can interact with alongside the raw data view to explore instance-level details. We conducted two case studies and an expert interview to demonstrate the effectiveness and usability of MAVIDSQL in helping model developers understand and diagnose *T2S* tasks.

In summary, the major contributions of our work are as follows.

- 1) An effective method for extracting semantic attributes from NL texts and differentiating structural features among SQL statements, enhancing text comparison and interactive exploration in *T2S* tasks.
- 2) MAVIDSQL, a visual analysis tool that generates model-agnostic visualizations for *T2S* tasks designed to assist both model developers and users in understanding and diagnosing potential errors.
- 3) Case studies and expert interviews that demonstrate the effectiveness of our approach in assisting users with the identification and comprehension of prediction errors in *T2S* models through an interactive exploration of model input and prediction results.

II. RELATED WORK

This section discusses the relevant research of our approach. The related work of this article can be categorized into two groups: *T2S* tasks and visual analytics for deep learning models.

A. *T2S* Tasks

T2S tasks have been developed to bridge the gap between users and data. It enables nonexpert users to access and perform intelligent exploratory analysis on tabular data [8], [9], [10]. Recently, novel *T2S* parsers using deep learning methods have gained promising results.

Numerous techniques have been developed for *T2S* tasks, which can be categorized into four modules: *input encoding*, *schema linking*, *output decoding*, and *output refinement*. Seq2SQL [11] and SQLNet [12] were early attempts to

apply deep neural networks to *T2S* tasks, using bi-LSTM to encode both NL sentences and database column names. Pre-trained models based on Transformer architecture, such as BERT [13], RoBERTa [14], GraPPa [15], and TaBERT [16], are commonly used.

To better clarify the concept of *schema-linking*, humans typically try to link key words in a NL question to corresponding elements in a given database when constructing SQL queries. Similarly, a *T2S* parser may benefit from employing a similar approach. For example, Schema-GNN [17] and Global-GNN [18] use graph-based approaches to represent database schemas and compute relevance probabilities for schema elements based on the question. RAT-SQL [19] uses a relation-aware transformer to explicitly encode relations between question words and schema elements, improving its ability to handle complex queries. The *output-decoding* in *T2S* models can be categorized into Sequence-based, Sketch-based, and Grammar-based approaches. Sequence-based methods consider SQL queries as a sequential set of sequences [11], while Sketch-based approaches [12] simplify the SQL generation task by decomposing it into multiple simple multicategory subtasks [20]. Grammar-based decoders [21] generate a sequence of grammar rules [22] instead of basic components and slots, allowing them to generate grammatically correct complex nested queries. IRNet [23] goes a step further by defining SemQL, an intermediate expression based on an abstract syntax tree to bridge NL and SQL generation. The *output-refinement*, which can be implemented during the decoding phase in order to reduce the possibility of errors [18] and attain better outcomes [24].

In addition, while LLMs, such as GPT-4, LLaMA, and Alpaca which are trained with reinforcement learning for human feedback (RLHF), have exhibited remarkable zero-shot capabilities [25]. There are still scenarios where these models may not perform as expected. In some cases, these models may encounter out-of-distribution data, which can lead to inaccurate or nonsensical predictions. Additionally, the performance of these models can be influenced by the quality and diversity of the training data. As a result, there is still a need for researchers and practitioners to carefully evaluate and fine-tune these models for specific applications and domains.

B. Visual Analytics for Deep Learning Models

There is a line of related research focusing on visualization to understand, interpret, and diagnose deep neural network models [26], [27]. Techniques such as direct inference and user interactivity in LSTMVis [5], and *what-if* explorations in Seq2SeqVis [6] have been commonly used for model interpretation. The mainstream research idea for model interpretation is to combine explainable models with visual analysis techniques. VBridge [28] and NLIZE [29] utilize popular explainable models such as SHAP and LIME to generate contribution-based explanations for a large number of features, which are then organized hierarchically to aid users in interactive refinement of models. For model diagnose, various works conduct performance analysis and provide support for common performance metrics such as accuracy, recall, true, and false prediction rates, including Squares [30] and Manifold [31], which are especially

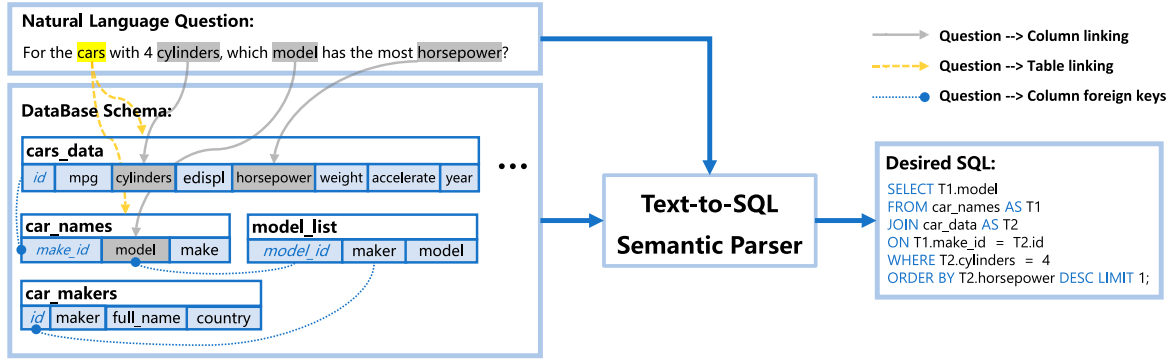


Fig. 1. Goal of T2S tasks is to convert the input NL question and database schema into its corresponding SQL queries. This process involves encoding the input question and database table/column names, and also linking the semantics of the question with the SQL schema, to improve the model's prediction robustness. The resulting output is the SQL statement that corresponds to the input question and schema.

effective in exploring multiclassification tasks. ConfusionFlow [32] proposes three levels of machine learning model exploration: global level, class level, and instance level. For the diagnosis of prediction errors in specific model task scenarios, GNNLens [33] facilitates the exploration and comprehension of prediction error patterns in graph neural network (GNN) models. This tool benefits both model developers and users in their endeavor to gain insights into the inner workings of GNNs.

In summary, while various approaches have been proposed to enhance the interpretability of deep learning models, limited efforts have been devoted to incorporating interaction and visualizations in T2S tasks to improve their explainability. To address this research gap, this article proposes a visualization tool that leverages text-linking comparison to aid developers and users in diagnosing and comprehending T2S tasks.

III. DOMAIN CHARACTERIZATION

T2S tasks are typically considered to be end-to-end semantic parsing tasks that generate SQL statements based on NL questions, database schemas, and schema-linking features constructed by hidden layers, graph structures, or other techniques. This section introduces the related background about the formal problem definition of T2S parsing, evaluation metrics for verifying T2S parsing, and datasets used in this study.

A. T2S Task Formulation

T2S tasks are classic examples of end-to-end semantic parsing. These tasks take a NL question and a database schema as input and output their corresponding SQL without the need for additional manual intervention or processing. Fig. 1 showcases the parsing process of T2S tasks from the Spider dataset, which employs similar schema annotation as RAT-SQL [19].

Given a NL question Q and its corresponding database schema $S = \langle T, C \rangle$, the schema consists of table names $T = \{t_1, t_2, \dots, t_{|T|}\}$, and all column names of each table t_i are expressed as a set $C = \{c_1^{t_1}, c_2^{t_1}, \dots, c_{|C|}^{t_1}, c_1^{t_2}, c_2^{t_2}, \dots, c_{|C|}^{t_{|T|}}\}$. Each table name t_i is described by its name and is further composed of several words $[t_{i1}, t_{i2}, \dots, t_{i|t_i|}]$, and each column $c_j^{t_i}$ in table

t_i is represented by words or a phrase $[c_1^{t_i}, c_2^{t_i}, \dots, c_{|c_j^{t_i}|}^{t_i}]$. The whole input can be represented as $X = \langle Q, S \rangle$, the goal of the T2S task is to generate the SQL query Y .

Then T2S models typically learn to represent the input using LSTM or Transformer-based encoding models, along with schema-linking techniques. After extensive training and fine-tuning, the models generate SQL queries through a decoder, as introduced in Section II-A.

Our work combines statistical evaluation metrics with the semantic information of the input sentences to create an interactive visual analytics system. This system enables users to explore the impact of semantics on model performance in a dynamic and interactive manner. When users identify error patterns of interest, they can delve into the SQL source data to view and analyze the underlying data structures and relationships. We will provide a detailed discussion of the visualization design and interaction logic in Section V.

B. Motivation

Previous studies have summarized the overview of the challenges encountered by researchers and developers during the development of novel deep learning models. These challenges have encompassed tasks such as debugging coding errors, model comparison, and comprehending the inherent characteristics of these models [31]. Drawing inspiration from these studies, we expand this investigation to T2S tasks and delineate potential challenges as follows.

1) *Lack of Interpretability in Model Comprehension:* The internal mechanisms of many T2S models act like a "black box," making their predictive logic difficult to interpret. Model developers apply advanced network architectures and pretraining techniques such as Seq2Seq, Transformer, BERT, and GPT-3, to enhance the performance of T2S tasks. However, different mechanisms lead to a lack of clear rationale or evidence for developers, which is essential in guiding the development and debugging of these models. The high-dimensional representations within the internal layers of the model pose a significant challenge in comprehending and inferring the predictive behavior exhibited by the model. For

instance, when the inquiry takes the form as, “What is the model of the car with the smallest amount of horsepower?” The SQL query predicted by the model is as follows: “SELECT cars_data.Horsepower FROM cars_data ORDER BY cars_data.Horsepower LIMIT 1.” The model’s misinterpretation of the question’s semantic context resulted in inaccurate predictions [1].

2) *Lack of Diagnosability in Evaluation Metrics*: Existing evaluation metrics for T2S tasks include accuracy, recall, precision, and F1-score, which serve as valuable and commonly used measures to assess the performance of various models. Nevertheless, these metrics may not provide a comprehensive understanding of the model’s performance. They also do not provide insights into the model’s errors and weaknesses, nor do they facilitate focused exploration of the model’s predictions [34]. In T2S tasks, execution accuracy (EX) serves as a metric to assess whether the execution result of the model-predicted SQL in the database matches the ground truth. Additionally, exact set match accuracy (ESM) evaluates the precision of the output SQL structure in relation to the ground truth SQL clauses. The result can be considered a correct model prediction only if all SQL subcomponents match. However, these metrics provide only numerical accuracy values and may not offer a systematic diagnosis of the errors underlying the model. For example, relying on evaluation metrics makes developers challenging to discern where the model is prone to prediction failures and in which query patterns the model is more likely to generate accurate SQL queries.

3) *Lack of Interactivity in Model Improvement*: Model developers encounter challenges during the analysis of model behavior due to handling a substantial volume of unstructured NL input and structured SQL output. On one hand, aligning and correlating a substantial volume of textual data in the model’s input and output. On the other hand, the dataset includes hundreds of databases and thousands of column names. Understanding the diverse contexts of these databases is essential and adds complexity to the data exploration and analysis process. Furthermore, a large volume of textual data also present challenges in analyzing and uncovering the primary data patterns within the constraints of limited space, thereby limiting the utilization of domain knowledge by developers. Designing an interactive system that enables developers to integrate domain knowledge and visualize the distribution of model input semantics and model output comparison could aid in gaining a comprehensive understanding of the strengths and weaknesses of models. This approach has the potential to provide insights to developers, improving their efficiency when iterating on new models, and it reduces the burden on developers while mitigating the risk of errors.

While the challenges faced by the model have been outlined, addressing these issues may not be straightforward. It requires a comprehensive understanding of the task’s workflow, input–output data, and domain knowledge. Visual analytics techniques can assist in both of these aspects. Model developers can apply appropriate visualizations to explore the large-scale and complex input and output data, thereby gaining additional insights into the data and uncovering relationships between the data

and the model output. The next section presents a design requirement analysis specifically tailored to visual analysis for T2S tasks.

IV. DESIGN REQUIREMENT ANALYSIS

MAVIDSQL is a tool designed to interpret and diagnose the error pattern patterns between model inputs and corresponding prediction results of T2S models. This tool aims to enhance the efficiency of model fine-tuning. The general goal of MAVIDSQL is to understand the factors that cause these models predict failure and to interactive discover common patterns of model prediction errors for improving the efficiency of fine-tuning. This could involve identifying patterns in the types of errors that the models tend to make, or identifying areas of the input text where the models have difficulty interpreting the meaning. By gaining insights from these factors, it may be possible to improve the performance of T2S models and make them more effective at generating precise queries. Based on the discussions with domain experts, a review of current literature, and our own attempts to reproduce the SOTA models one Spider leaderboard, the specific requirements are formulated as follows:

A. R1: Associate and Align T2S Data

Neural network models are characterized by their black-box nature, which refers to the embedding of input and output features within the hidden state of the model. Consequently, the interpretation and application of such models are hindered by a lack of intuitive understanding of the underlying data. In the context of model input data, the test dataset comprises a substantial number of text messages that exhibit similar structural characteristics but differ in terms of their semantic content. Identifying which question patterns are beneficial for model prediction is important for enhancing model performance. For model output data, SQL statements have a complex and flexible structure, as well as the absence of precise evaluation metrics to evaluate the accuracy of model output data, a collaborative approach involving aligned model predictions, ground truth data, human feedback can be instrumental in identifying and interpreting patterns of model error.

B. R2: Provide an Overview of T2S Prediction

Experts and literature research show that an overview of the model performance is crucial for T2S task results analysis. To gain an overview of the dataset and predict results, the system needs to summarize various types of information, such as performance statistics and ground truth label distribution. This information, covering various aspects of a T2S model, needs to be organized and presented in a clear manner. Meanwhile, users should be provided with the interrelation between this information to aid in forming preliminary hypotheses regarding potential error patterns in T2S results. This entails identifying a group of incorrect predictions that exhibit comparable question patterns or SQL structures.

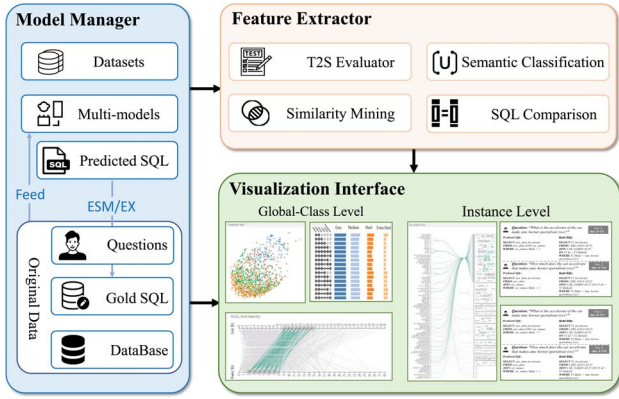


Fig. 2. Overview of the system architecture. Model manager: it executes T2S models, generates predictions, and manages datasets. Feature extractor: it evaluates the model provided by the model manager, extracts relevant semantic features, and passes them into the visualization interface. Visualization Interface: it allows for interactive exploration from global-class level to instance level.

C. R3: Identify T2S Model Failure Patterns

After developing initial hypotheses about the error patterns, users need more detailed information to verify them. Specifically, users need to examine the question pattern or SQL structure shared by a set of wrong predictions and verify whether error patterns formed by these patterns make sense in analyzing T2S based on their domain knowledge. During the literature review and expert interview, domain experts agreed that there are some relatively complex SQL structures that often make the model wrong [1]. For example, the model could encounter an issue of robustness, resulting in its inability to identify the SQL table name that corresponds to a given statement. Therefore, the system should support users in examining the influences of these association characteristics and identifying error patterns.

D. R4: Support Multilevel Exploration of T2S Model Prediction

For a comprehensive understanding and analysis of the T2S model's predictions, visualization should empower users to explore the model input and output. The exploration should be conducted on multiple levels, comprising an overview of the model's performance at the global level, and a detailed analysis of the raw data at the instance level. Specifically, users infer the causes of error patterns by interacting with a filtering mechanism for global performance, which allows them to pinpoint the raw question and SQL structures. Aggregating sentences of the same category together makes it convenient for users to quickly search and infer them.

V. MAVIDSQL

As shown in Fig. 2, MAVIDSQL consists of three major modules: the model manager, feature extractor, and visualization interface. The model manager module is responsible for storing and managing T2S datasets and model predictions. The feature extractor model carries out the necessary feature construct procedures for analyzing the T2S model predictions.

The processed data feature is then passed to the visualization module, which supports the interactive visual analysis of the T2S. The model manager and feature extractor modules are developed using the Python, while MongoDB is employed as the data storage engine. The system is integrated into a backend web server built using flask. We implement the visualization module as a front-end application using Vue, JavaScript, and D3. The visualization interface is illustrated in Fig. 3.

A. Methods

In this section, we introduce the primary data mining and layout optimization techniques employed in this article for extracting features from input and output data. For model input NL question data, we perform semantic classification and sentence similarity mining to assess the similarity characteristics. For model output SQL queries, we calculate the text-level differences between the model-predicted SQL and the ground truth gold SQL to facilitate large-scale comparisons of p-g SQL pairs in visualization.

1) *Input Natural Language Question Similarity Mining*: In the context of T2S models, the performance of prediction results is influenced by the variability in NL descriptions and syntax structures. To address the challenge of measuring semantic similarity in interrogative and imperative sentences, our approach was informed by the NLP field and has undergone multiple experimental iterations to refine and optimize the methodology. The approach integrates multigranularity semantic information, including lexical semantics, dependency syntax, and sentence structure patterns, to accurately capture the nuances of sentence similarity. T2S models usually trains an input-encoding model to map NL query tokens into a high-dimensional embedding vector space. We initially considered whether it could directly use the results of the model input encoding to calculate the similarity of the questions. However, it is experimentally proven that mapping input sentences into high-dimensional vector space after encoding by various models can introduce uncertainty, make it difficult to reflect the desired syntactic structure and semantic information.

Researchers commonly consider both lexical and semantic factors when evaluating the semantic similarity between sentences in the NLP community [35]. Various corpus-based methods have been proposed for this purpose, including N-gram, WordNet, Jaccard similarity, and Word Embedding. These approaches primarily concentrate on the semantic representation of words, exhibiting limited engagement with the contextual associations within the text. For instance, two syntactically similar NL questions, "How much water is available?" and "How many countries are listed?" may exhibit substantial semantic differences in their vector representations, but they correspond to very similar SQL query structures.

Additionally, dependency parsing is another widely used technique in NLP that involves identifying the syntactic structure of a sentence by analyzing the relationships between words in a sentence. The output of a dependency parser is a treelike structure that represents the syntactic relationships between words in a sentence. On the other hand, role labeling is a

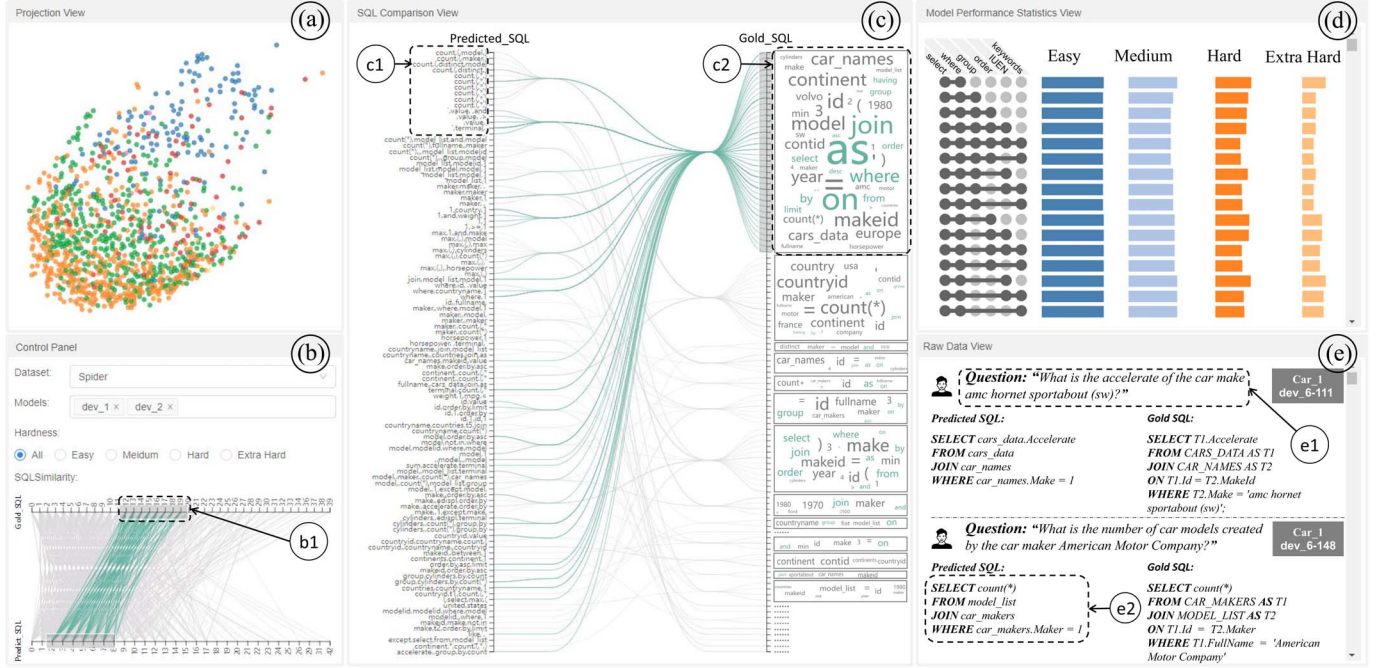


Fig. 3. Explanatory interface of MAVIDSQL consists of five views. (a) Projection view provides the similarity among input NL questions, and supports lasso-selection interactions for exploring global-level model performance and instance-level SQL prediction results. (b) Control panel enables users to interactively configure basic parameters (e.g., the dataset). Users can also swipe a particular set of SQL predictions of interest to obtain further details. (c) SQL comparison view further visualizes the specific category of user-selected SQL comparison results. (d) Model performance statistics view provides users with an overall evaluation of SQL prediction performance. (e) Raw data view presents detailed information on the original data used in the model.

technique that involves identifying the semantic relationships between words in a sentence, such as who and what is the subject, object, or indirect object of a sentence. However, relying solely on dependency syntax trees for evaluation may prove inadequate in capturing the nuances of sentence structure. For example, consider the following inquiries: “In which years cars were produced weighing no less than 3000 and no more than 4000?” and “What are the different years in which there were cars produced that weighed less than 4000 and also cars that weighed more than 3000?” Although these two sentences exhibit semantic proximity, they manifest significant differences in their syntactic structures.

We comprehensively integrate both the semantic information and syntactic structure of the model input questions. Furthermore, to correlate questions from similar databases, we also incorporate additional meta features such as sentence length and the count of shared words. The process involves generating dependency parsing trees for two sentences, identifying corresponding roles through specific formulas. We utilize Hanlp [36] for the extraction of relevant features, which comprise categorically labeled identifiers based on the NLP task. The dependency parsing trees of two sentences are obtained using Hanlp, followed by the computation of identical dependency parsing roles’ positions in the sentences.

First, we calculate sentence semantic similarity between two NL question sentences Q_a and Q_b , denoted as $\text{SenSemSim}(Q_a, Q_b)$. We tokenize sentence Q to obtain tokens T_1, \dots, T_n , where each token represents an individual word in the sentence.

Then, we utilize the Word2Vec module from the gensim library to train a Skip-gram model and determine the correlation tokens denoted as $\text{SemSim}(T_{ai}, T_{bj})$. The calculation for $\text{SenSemSim}(Q_a, Q_b)$ is as follows:

$$\text{SenSemSim}(Q_a, Q_b) = \left(\frac{\sum_{i=1}^m S_{ai}}{m} + \frac{\sum_{j=1}^n S_{bj}}{n} \right) / 2 \quad (1)$$

$$S_{ai} = \max(\text{Sem}(T_{ai}, T_{b1}), \dots, \text{Sem}(T_{ai}, T_{bn})) \quad (2)$$

$$S_{bj} = \max(\text{Sem}(T_{a1}, T_{bj}), \dots, \text{Sem}(T_{am}, T_{bj})) \quad (3)$$

$$\text{SemSim}(T_{ai}, T_{bj}) = \frac{V(T_{ai}) \cdot V(T_{bj})}{\|V(T_{ai})\| \cdot \|V(T_{bj})\|} \quad (4)$$

where S_{ai} and S_{bi} are intermediate variables that measure the semantic similarity between words in sentences Q_a and Q_b , respectively. $V(T_{ai})$ and $V(T_{bj})$ are the vector representation of tokens T_{ai} and T_{bj} obtained from the word2Vec model. $\|V(T_{ai})\|$ and $\|V(T_{bj})\|$ are the Euclidean norms (or magnitudes) of the vector representations for T_{ai} and T_{bj} , respectively.

To enhance the assessment of sentence syntactic similarity between questions Q_a and Q_b , denoted as $\text{SenSynSim}(Q_a, Q_b)$, we focus on the hierarchy distance, $\text{HierDis}(r_{ai}, r_{bi})$, between tokens. In this context, r_{ai} and r_{bi} are defined as the i th dependency parsing roles corresponding to the tokens in questions Q_a and Q_b , respectively. The measurement of $\text{HierDis}(r_{ai}, r_{bi})$ is pivotal as it quantifies the hierarchical distance between each node and the root in the dependency syntactic tree of the

sentences. The formula for calculating the semantic similarity between sentences is provided as follows:

$$\text{SenSynSim}(Q_a, Q_b) = \frac{\sum_{i=1}^n \text{HierDis}(r_{ai}, r_{bi})}{n} \quad (5)$$

$$\text{HierDis}(r_{ai}, r_{bi}) = 1 - \left| \frac{\text{Deep}(r_{ai}) - \text{Deep}(r_{bi})}{\text{Deep}(r_{ai}) + \text{Deep}(r_{bi})} \right|. \quad (6)$$

Here, $\text{Deep}(r_{ai})$ and $\text{Deep}(r_{bi})$ denote the depth of the respective roles in the semantic dependency parsing tree, using the root node as the reference point. The proximity of these depths to one another indicates the degree of similarity between the two questions.

In addition, to distinguish queries aimed at different databases, we have incorporated an analysis of sentence metadata, which includes sentence length and the frequency of common words. The similarity of common words between sentences denoted by $\text{ComWordSim}(Q_a, Q_b)$. Moreover, the sentence length similarity denoted as $\text{SenLenSim}(Q_a, Q_b)$

$$\text{ComWordSim}(Q_a, Q_b) = \frac{2 * \text{SameWord}}{\text{Len}(Q_a) + \text{Len}(Q_b)} \quad (7)$$

$$\text{SenLenSim}(Q_a, Q_b) = 1 - \left| \frac{\text{Len}(Q_a) - \text{Len}(Q_b)}{\text{Len}(Q_a) + \text{Len}(Q_b)} \right| \quad (8)$$

where SameWord represents the number of identical words present in questions Q_a and Q_b , while $\text{Len}(Q_a)$ and $\text{Len}(Q_b)$ denote the number of tokens in questions. To attain optimal results, it is imperative to adopt a comprehensive approach rather than focusing exclusively on a single aspect. Therefore, in the processing of sentence similarity, we treat the aforementioned factors as different feature items of the sentences. By integrating these feature items and assigning different weights according to their importance, we obtain the final similarity score between the sentences.

Finally, the final formula $\text{Similarity}(Q_a, Q_b)$ is derived by weighting the semantic and syntactic structures and the metadata. When querying different databases with the intention of expressing similar SQL query, the common vocabulary in the two queries tends to decrease, leading to a lower sensitivity to sentence metadata in the final metric, and consequently, a lesser weight is assigned to it. Conversely, the dependency parse tree and semantic analysis exhibit a higher sensitivity, and thus, are assigned a greater weight. The final calculation formula for $\text{Similarity}(Q_a, Q_b)$ is as follows:

$$\begin{aligned} \text{Similarity}(Q_a, Q_b) &= \gamma_1 \text{SenSemSim}(Q_a, Q_b) \\ &+ \gamma_2 \text{SenSynSim}(Q_a, Q_b) \\ &+ \gamma_3 \text{ComWordSim}(Q_a, Q_b) \\ &+ \gamma_4 \text{SenLenSim}(Q_a, Q_b) \end{aligned} \quad (9)$$

s.t. $\gamma_1 + \gamma_2 + \gamma_3 + \gamma_4 = 1.$

Employing our sentence similarity algorithms, we calculated the pairwise similarity among input sentences. This similarity metric was then applied to quantify sentence distances within the projection view of dimensionality reduction processes.

2) *Output SQL Result Side-by-Side Alignment*: The evaluation of T2S models typically involves a set of test dataset, comprising a series of query inputs, against which the model's predictions are compared to assess its performance against the

Algorithm 1: SQL Statement Comparison

Result: Two sets of the same token elements after elimination, p_{diff} and g_{diff}

```

1 Input: The string of model predict SQL  $p_{sql}$  and gold SQL  $g_{sql}$  provided by dataset;
2  $p_{sql}$  and  $g_{sql}$  are parsed into sets of tokens  $p_{tokens}$  and  $g_{tokens}$  respectively, based on SQL syntax. During the parsing process, table aliases are removed and SQL keywords are standardized to a uniform format;
3 /* Eliminate the same tokens between  $p_{tokens}$  and  $g_{tokens}$  */;
4 Set  $p_{diff} \leftarrow p_{tokens}$ ;
5 Set  $g_{diff} \leftarrow g_{tokens}$ ;
6 foreach  $p_{tokens}$   $p_t$  do
7   if  $p_t$  in  $g_{tokens}$  then
8      $p_{diff} \leftarrow$  remove  $p_t$  from  $p_{diff}$ ;
9      $g_{diff} \leftarrow$  remove  $p_t$  from  $g_{diff}$ ;
10  end
11 end
12 return ( $p_{diff}$ ,  $g_{diff}$ );

```

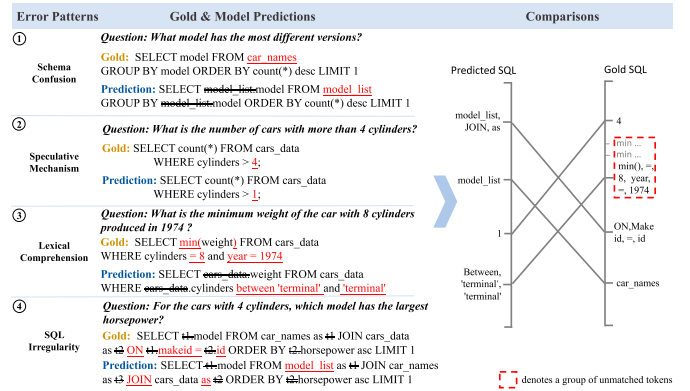


Fig. 4. Common error patterns in T2S models and principal methods of SQL alignment.

ground truth. Developers are required to conduct thorough analyses of extensive comparisons between predicted and actual SQL outputs when evaluating model performance. The complexity of SQL statement structures necessitates domain knowledge and an in-depth comprehension of the database context. Additionally, it is challenging to display an adequate number of SQL comparison pairs within the constraints of limited screen space.

Extensive research in visual analytics focuses on the alignment and comparison of textual data [37], such as sequence-aligned [38], aligned barcodes [39], and side-by-side [40]. We aim for developers to preserve the essential information in SQL statements during the analysis process, while minimizing redundant words, such as SQL keywords, that appear in both predictions and the ground truth. We used Algorithm 1, to preprocess the model's output of predict SQL and gold SQL. As shown in Fig. 4, considering the two SQL statements provided

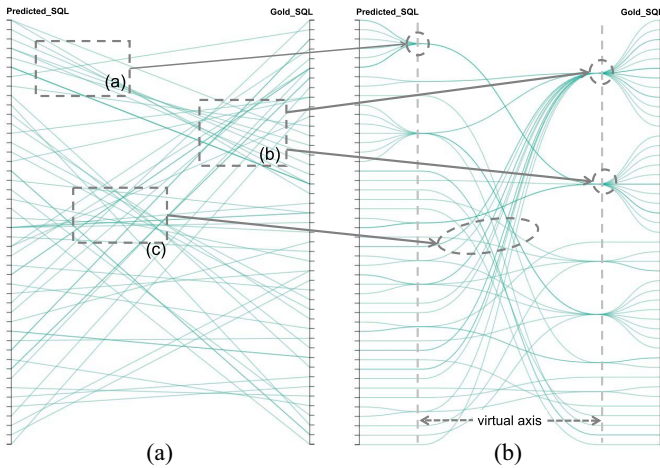


Fig. 5. Comparison of the different layouts for parallel coordinates plots (PCPs) depicting two dimensions from the comparison of predicted and gold SQL: (a) classic parallel coordinates plot; and (b) edge-bundling parallel coordinates plot. The edge-bundling technique effectively reduces visual clutter caused by line crossings and alleviates the visual cognitive load.

as input, this approach involved parsing each SQL statement into a set of tokens, while removing extraneous information such as table aliases and punctuation. We conducted a pairwise comparison between the generated tokens, eliminating any identical tokens and identifying discrepant ones to accentuate dissimilarities between the two token sets.

We visualize the data of SQL pairs using a side-by-side approach based on parallel coordinates shown in Fig. 5(a). Furthermore, we group SQL statements with common prefixes, allowing for class-level comparison of differences among SQL queries of this category. This approach is illustrated in Fig. 4.

3) *SQL Group Comparison Layout Optimization*: We aspire to enable users to interactively analyze a large number of SQL pairs within the constraints of limited screen space. To this end, we design views based on parallel coordinate visualization, which involve the comparison of SQL tokens in a side-by-side alignment, as shown in Fig. 5. However, with the increase in data volume, the resultant crossing of lines leads to visual clutter, posing challenges in pattern recognition. To address the issue, we employ a hierarchical density aggregation technique to group similar results and present them as word clouds organized by category. These approaches draw inspiration from the previous work on aggregation techniques [7] and visual designs [41]. We employ edge-bundling and categorization techniques to reduce visual clutter and cognitive load for users by sorting and grouping data on the axes.

In conventional parallel coordinate plots, lines are employed to connect the values between two axes. Nevertheless, this approach often engenders issues of visual clutter and excessive plotting, particularly as the quantity of data points and connections increases. Initially, following the previously discussed group method, we arrange SQL tokens belonging to the same group together for presentation. Items within a group will appear higher on the axis based on their quantity. Consequently,

as the volume of data on the axis escalates, the primary error patterns predicted by the model become increasingly concentrated toward the upper of the axis.

To provide smooth curves that are easy to follow with the eyes, we bundle cluster around each axis, after categorizing and sorting the data. The classic parallel coordinate plot employs lines to link the values between the two axes. However, this often leads to issues with visual clutter and overplotting, especially when the number of data points and links increases. Drawing inspiration from previous work [41], we add virtual bundling axes adjacent to each data axis in the parallel coordinate plot, as illustrated in Fig. 5. All data points belonging to the same category are bound to the corresponding position on the virtual axis through cubic Bézier curves. The y-coordinate of the points is calculated based on the median of all points in the same category. The points on the two aggregated virtual axes are also connected using Bézier curves, which helps to improve the visual grouping and reduce visual clutter. From Fig. 5, we can observe the differences between the classical parallel coordinate plot and the aggregated version.

B. Visualization

As shown in Fig. 3, MAVIDSQL consists of five visualization modules. The projection view (a) and the model performance statistics view (d) provide an overview of the input data and the model's prediction results. The control panel (b) provides users with options to select a dataset, choose specific models, and hardness levels for analysis. The SQL comparison view (c) and the raw data view (e) offer interactive analysis at different granularities, assisting users in identifying error patterns of interest. In this section, we provide a comprehensive overview of the visual representations and interactive capabilities of each module, along with a detailed explanation of our design considerations.

1) *Projection View*: To better illustrate the association between input NL question (R1) and help users analyze the impact of semantic patterns on prediction results (R2), we design the projection view to visualize the similarity of input sentence. In the projection view, we use the multidimensional scaling projection algorithm to map all the input NL questions. Each sentence is described by a series similarity distance, which is introduced in Section V-A1. The distance between these sentences is further mapped to 2-D space to obtain a quick overview of semantic association of input interrogatives. The projection view allows users to explore the similarity of combine question categories and semantic information and explore the impact of different categories of sentences on the input result. The visual encoding is used to differentiate between different question categories (e.g., blue for “how” questions, orange for “what” questions, red for “for which” questions, and green for “others” questions). It can be helpful for investigating whether the sentence with similar semantic structure share similar error patterns. In our implementation, we categorize the input question into seven groups according to wh-word in conjunction with expert opinion. When users lasso-select a set of sentences in a projection plane, the corresponding prediction

errors and accuracy will be displayed. This enables users to interactively explore the impact of nodes on model predictions.

Since we have extracted the semantic features of input question and model predict score, we are considering displaying more data information when projecting dimensionality reduction while avoiding the visual clutter caused by data superposition.

2) *SQL Comparison View*: The *SQL comparison view* is the key component of MAVIDSQL, which allowing users to compare and identify discrepancies between model-generated SQL and ground truth (R2, R3, R4). The system should reveal the difference between both predict and gold SQL queries. We work closely with the experts to incorporate their feedback and the core design consideration is to show scenarios of SQL prediction failures while retaining enough of the original SQL information. During the design process, we faced the challenge of presenting a large number of predicted and gold SQL pairs without aggravating the cognitive burden. Therefore, we followed both article and code-based prototyping approaches to refine this system according to user feedback. We summarize our design considerations and describe the details of each component as follows.

a) *SQL similarity component*: As introduced in Section V-A2, the SQL similarity component eliminates the same tokens from the predicted and ground truth SQL statements and defines the number of remaining words as similarity between them. As shown in Fig. 3(b1), this component provides an initial overview of SQL queries where the model has failed to accurately predict the expected output. It features a two-axis parallel coordinate component that showcases the similarity between the predicted and ground truth SQL. In this component's design, each polyline represents a specific number of remaining words after removing the common tokens. Users can explore the distribution of these errors by brushing on each coordinate to select a specific sequence within a certain range. The SQL-pairs comparison component then displays the corresponding detail SQL words for further class-level analysis.

b) *SQL-pairs comparison component*: To mitigate potential uncertainties associated with existing model evaluation strategies (e.g., false negative and worst case [42]), and to prevent visual clutter caused by displaying large amount of SQL statements in text form. We design a comprehensive SQL comparison component that allows users to easily compare the predicted and ground truth SQL statements at a class level. This view presents the data in a concise and informative manner, making it easy for users to interpret and analyze the results. Previous work has explored a class-level analysis of machine learning models to understand their behavior and identify common error patterns [31], [32]. Inspired by previous research, we use this strategy to investigate error patterns in T2S tasks. We employ an enhanced edge-bundling layout for interactive parallel coordinates to support a comparative analysis of predicted and gold SQL queries, present in Section V-A3. The SQL groups that have been aggregated will be presented in this view, supporting the exploration at the class level. This layout allows users to visually identify patterns and discrepancies in the data

side-by-side. Users can select which group of detail information they are interested in to be displayed in the raw data view.

As shown in Fig. 3(c), when user selects a group of predicted and ground truth SQL results from the SQL similarity component, the SQL comparison view will display more detailed SQL comparison information categorized by specific word feature metrics. The two axes in the SQL similarity component represent the remaining words in the predicted and ground truth SQL statements after eliminating the common tokens. The left axis represents the predicted SQL tokens and the right axis represents the ground truth SQL tokens. Both axes are categorized based on the common initial tokens and sorted by their quantity. Each tick on axis represents a specific state after eliminating the common tokens. The tokens are mapped on both sides of the axis, as shown in Fig. 3(c1) and 3(c2). Given the wide variety of forms in which predicted SQL statements can occur, the associated tokens may be scattered. To avoid the potential loss of important cases, we display the original text directly on the axis. However, as gold SQL statements are typically structured uniformly, we present them in the form of word clouds to facilitate class level pattern recognition and improve user comprehension. The predict and gold SQL pairs between the two axes are connected using an improved version of the cubic Bézier curves.

In general, the SQL comparison view aims to present the distribution of predicted results at the class level. To address the visual clutter and reduce cognitive load during parallel coordinate plotting, we aggregate word clouds using a combination of classification and edge bundling based on axis number ordering. This technique effectively reduces visual clutter and enhances the user experience.

3) *Model Performance Statistics View*: To facilitate evaluation and calculation, each SQL statement within a structure is assigned to an individual object instance, as depicted in Fig. 6. We first compare the ground truth and predicted components on a module-by-module basis. This allows users to obtain the accuracy of each module of the predicted SQL [Fig. 6(a)]. Furthermore, since all databases have executable SQLite files, we can also measure the EX of the T2S model. To do so, we execute the generated SQL statements in the SQLite environment and compare the results of the gold and predicted SQL runs [Fig. 6(b)].

We design the model performance statistics view to help users intuitively perceive the distribution of the model performance, as shown in Fig. 3(d). The model performance statistics view includes two components: a multiset combination matrix plot and a stack bar chart, inspired by Upset [43], LineUp [44], and VSumVis [45]. For evaluation of the model performance, ESM is the primary metric, which is a multidimensional metric that predicts various components of the SQL statements such as select, where, group, order, IUEN, and keywords. The goal is to classify and analyze the model performance based on these metrics to identify potential error patterns in the model. Therefore, we employ the combination matrix from UpSet [43] to display different categories in the model performance statistics view.

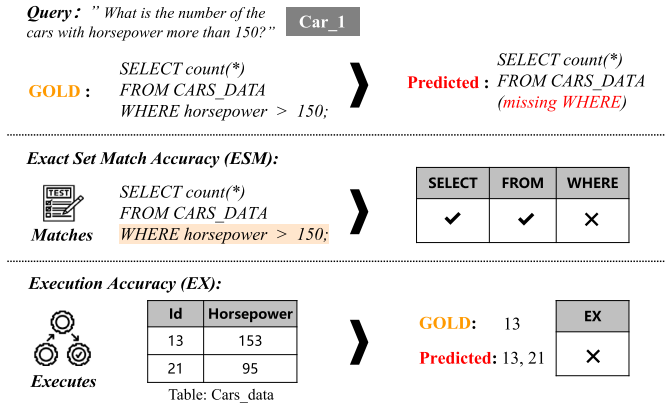


Fig. 6. Examples of evaluation metrics for T2S tasks, illustrating how model predictions are evaluated through exact set match accuracy and execution accuracy metrics. Evaluating errors in different modules can assist users in analyzing the semantic understanding issues of the model across different SQL slots. In addition, EX can prevent misclassification caused by SQL variations in different forms.

Moreover, queries in the dataset are often categorized by their level of difficulty. Examining the model predictions for different levels of difficulty can greatly aid in identifying and improving the performance of the model. We display the model's performance on different difficulty levels through four bar charts.

4) *Raw Data View*: As shown in Fig. 3(e), the raw data view is linked interactively with the SQL comparison view, enabling users to access detailed information by conducting an instance-level exploration. In the process of designing the system, it becomes apparent that users would benefit from intuitive guidance at the instance-level when exploring raw data, without being overwhelmed by cognitive load. In raw data view, we provide a comprehensive set of raw data on the model's input and output, including the NL questions, the SQL statements generated by the model and dataset ground truth. All the raw data are represented in original text format, annotated with database numbers and corresponding data index. Researchers can combine SQL comparison view and model performance statistics view to analyze the results of the model. Additional offline data can be exported, allowing users to debug and improve their models.

C. User Interactions

The MAVIDSQL provides a set of interactions, which facilitate the interplay of different views and enable multilevel exploration with details on demand.

1) *Lassoing for Input Exploration*: First, users can view the overall performance of the model, or obtain a preliminary overview, through the projection view and model performance statistics view. To enhance scalable exploration in the projection view, users may employ the lasso tool to select concentrate on particular instances of interest. Then, the detailed information will be displayed in the SQL similarity component and the SQL comparison view.

2) *Filtering for Dataset Selection*: To filter datasets and text complexity according to user interests, users have the capability to select and view these filters within the control panel.

3) *Brushing for Feature Analysis*: Users can employ the brushing feature on SQL similarity and SQL comparison components to assess the unmatched SQL tokens. This process facilitates the discovery of distinct model prediction patterns. Furthermore, the SQL predictions brushed in the SQL comparison view will be displayed in the raw data view, assisting users in detailed analysis.

4) *Clicking for Detailed Access*: Our system offers click-responsive features for viewing detailed information. For instance, users can click on the SQL comparison component to generate a word cloud visualization, highlighting a group of SQL keywords present in the gold SQL after performing a comparison.

VI. EVALUATION

In this section, we demonstrate the effectiveness of MAVIDSQL in facilitating the understanding and diagnosis of T2S tasks at both the global-class and instance levels. Two case studies and expert interviews are conducted with three domain experts (*E1*, *E2*, and *E3*), within the Spider dataset [46] and test suite data [42]. *E1* is an NLP researcher with a decade of experience in NLP. Furthermore, *E1* possesses expertise in applying and designing deep learning models for the purpose of converting unstructured text into structured formats. Both *E2* and *E3* are senior Ph.D. students in computer science. *E2* is engaged in the visualization of NLP tasks, while *E3* primary focus on data visual analytics. Notably, none of these experts are coauthors in our research. The two cases are discovered by *E1* and *E2* during the system exploration. Furthermore, comprehensive feedback from all the experts is also collected and summarized.

A. Datasets and Benchmark

In our case study, we employ Spider [46] dataset to establish analytical tasks. It is recognized as one of the most widely used single-turn datasets in the T2S community. The datasets is a large-scale, cross-domain semantic parsing dataset containing 10 181 NL questions and 5693 unique SQL queries across 200 databases, spanning 138 different domains.

The SQL queries are divided into four levels based on their difficulty: easy, medium, hard, and extra hard. This categorization allows for a more comprehensive evaluation of model performance across various types of queries. This categorization allows for a more comprehensive evaluation of model performance across various types of queries. As illustrated in Fig. 6, each Spider instance consists of: *database id*, *question*, *query*, and *gold sql*. The model's input-output performance is assessed using the 21 leaderboard submissions for the Spider dataset as proposed by Zhong et al. [42] to ensure the validity of our results.

B. Model Metrics and Error Patterns

1) *Evaluation Metrics of T2S Tasks*: Typically, T2S parsers are evaluated by comparing the generated SQL queries to the ground truth. Specifically, there are two main types of evaluation metrics used for evaluating the T2S tasks: ESM and

EX [46]. Additionally, the official evaluation metric for the Spider dataset is the test suites accuracy [42], which is an expansion of the ESM.

ESM is determined by comparing whether the sets of SQL clause match exactly between the ground-truth SQL query and the prediction. Both the prediction and the ground truth are parsed into bags of several subcomponents, including SELECT, WHERE, GROUP BY, ORDER BY, KEYWORDS, and all SQL keywords without column names and operators. The evaluation script compares the sets of subcomponents in each SQL component side by side to determine if they match between the ground-truth and predicted SQL queries. EX has been introduced to consider the accurate execution of a predicted SQL query on a specific instance of a database. Considering the possible variance in SQL structure predictions across different models, the developers introduce that these queries must be executed directly within the SQL database.

Moreover, existing *T2S* datasets are often divided into classes based on the difficulty level of the queries, but most evaluation metrics are only used for global-level analysis, making it difficult to conduct a more detailed exploration of the performance of *T2S* tasks. To overcome these limitations, we propose a global-class level exploration that combines our designed projection view, model performance statistics view, and SQL similarity view in MAVIDSQL. These views allow for a more fine-grained analysis of the performance of *T2S* models and facilitate the discovery of patterns and errors in the model's predictions.

Global-class level evaluation metrics provide a rapid and high-level overview of model performance. In contrast, instance-level analysis enables users to gain a more comprehensive understanding of the model's performance by identifying specific errors and patterns. In our system, we present global-class level model metrics as an overview in the statistical view, allowing users to understand the model's overall performance distribution. Subsequently, through interaction, users can examine more detailed, instance-level analysis. By leveraging instance-level analysis, users can explore individual instances in detail, compare the predicted and gold SQL statements, and evaluate the model's performance on specific subtasks, ultimately leading to a more comprehensive understanding of the model's strengths and limitations. Additionally, users can also uncover additional instances based on existing error patterns.

2) *Error Patterns of T2S Tasks*: Through collaborative exploration within our experts, we have characterized the common error patterns that MAVIDSQL can identify.

a) *P1: Schema confusion*: In the *T2S* tasks, it is imperative for the model to effectively comprehend the intent of the questions. This involves discerning the primary element of the query and creating a schema that links it to the corresponding database [17]. Subsequently, it is crucial to precisely interpret the semantic query's intentions, such as discerning whether it is seeking a maximum or minimum value or determining which column to group by. As illustrated in Fig. 4(1), while the model accurately predicted subsequent filtering and aggregation operations, confusion regarding the table structure led

to the selection of an incorrect database table, resulting in a prediction failure.

b) *P2: Lexical comprehension*: When the input query contains vocabulary or expressions not encountered during the model's training, there is a potential for misunderstanding the meaning of these words, which might lead to incorrect interpretation or selection of the query's intent [1]. As demonstrated in Fig. 4(2), the model misinterpreted the intent of the question, failing to predict the "MIN" operation.

c) *P3: Speculative mechanism*: Due to the characteristics of the Spider dataset's evaluation mechanism, test suite EX focuses on assessing the SQL execution and the alignment of various SQL components. During the evaluation process, it does not examine the values [42]. During the prediction process, several models resort to utilizing default value inputs to populate the SQL queries. As shown in Fig. 4(3), although the SQL structure and column names are correctly forecasted, there is a mismatch with the actual raw data in the SQL.

d) *P4: SQL irregularity*: Several models produce prediction results that do not conform to SQL syntax rules. For instance, as shown in Fig. 4(4), although the models accurately identify the "JOIN" module, they fail to specify the conditions for the "JOIN." This discrepancy leads to a divergence in EX between the predicted SQL and the gold SQL.

C. Case 1: Analyzing Error Patterns From Model Prediction

The first case conducted by *E1* involved an exploration of the database "car_1" within the Spider dataset. This was carried out in conjunction with our system, and a comparative analysis was made against the workflow of the original analytical model. When *E1* load in a *T2S* model and accesses MAVIDSQL, the system extracts the relevant model prediction data and add it to the existing dataset. Then, the system recalculates the feature extraction results and update the visualization views. As shown in Fig. 3(d), *E1* acquire a comprehensive understanding of the model's performance while also obtaining information on varying degrees of hardness. *E1* analyzed the projection views and model performance statics views to gain an overview from a global level to the class level. During the exploration, *E1* observed that in the "car_1" database, the predictions yielded satisfactory results for easy and medium tasks but underperformed in handling hard-level predictions.

E1 further investigated the projection view, focusing on exploring specific questions of interest. The projection view clusters semantically similar sentences together based on the SQL similarity measure, facilitating easy navigation and exploration of the input questions. This combination of interactive visualization techniques enables *E1* to gain insight into how well the model performs for different question types. As depicted in Fig. 7(1), *E1* employ the lasso tool to select several areas of query sentences, where sentences of the same category tend to cluster together. The corresponding model prediction performance is displayed in the SQL Similarity Component, Fig. 7(2). *E1* found that the predicted axis values predominantly cluster at or below five, whereas the actual gold axis values are mainly distributed at one, eight, and above. *E1* speculated

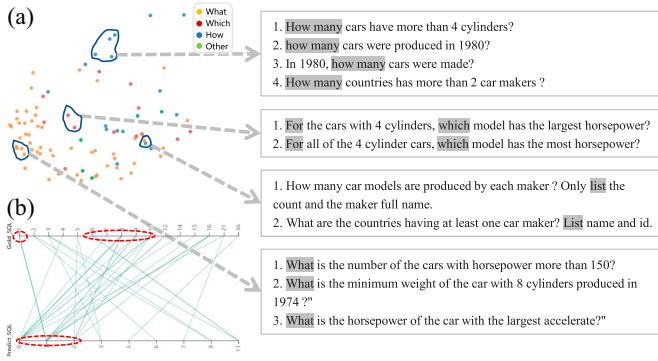


Fig. 7. Exploring at the global-class level combining projection view and SQL similarity component: (a) examples of a user exploring input sentences freely, it is common for questions in the same category to exhibit correlations in their similarity metrics; and (b) potential error patterns in initial user inference can help users compare the differences between the predicted and ground truth SQL statements.

the model's reduced similarity score, when compared to the gold SQL statement, may stem from inadequate prediction of certain words or the omission of key components in the SQL statement. *E1* further filtered the results in the SQL similarity view and performed instance-level comparisons by combining specific SQL comparison components and raw data views to verify false evidence.

As shown in Fig. 3(c2), *E1* observed in the SQL comparison view that the gold SQL contains a larger number of SQL keywords. The word cloud visualization demonstrates that SQL keywords, such as “AS,” “JOIN,” “ON,” and “WHERE” constitute a substantial proportion of the visualization. *E1* integrating experience with analysis of raw data, observed that while the model accurately predicted certain database values, some predictions were erroneous. This inaccuracy was primarily due to a disregard for SQL syntax rules, specifically the failure to properly match the 'ON' condition during SQL joins (*P4*). Additionally, table names, such as “car names” and “cars data” and column names, such as “model” and “makeid” also occupy a considerable proportion in the word cloud. *E1* further explored this category of SQL statements by utilizing the Raw data view. The specific details related to the problem, including the predicted and ground truth SQL statements, along with the corresponding NL questions, are displayed. As depicted in Fig. 3(e1), *E1* also noted that the model evaluation does not examine the values, leading to the replacement of the WHERE clause with “1” instead of leaving it unstandardized (*P3*). This is exemplified in the query “AMC Hornet Sportabout (sw).”

Furthermore, while the predicted SQL statement correctly identifies the table and column names, the lack of uniformity in the form of the predicted and gold SQL statements, including the absence of keywords such as “AS” and “ON,” results in many SQL keywords remaining after comparison. Upon inspecting the original evaluation results for this example, *E1* found that the match accuracy was *True* while the EX was *False*. *E1* indicated that this inactive instance-level analysis can help distinguish and validate the reasons for inconsistencies

between match and EX. In Fig. 3(e2), *E1* observed that it is demonstrated that the model has misunderstood the correspondence between “car maker” and the corresponding column name. “American Motor Company” should correspond to the “full name” column in the “car makers” table, but the model mistakenly predicted the “maker” column (*P2*). Further analysis revealed that cases with fewer unmatched tokens in the Gold SQL generally involve errors in value predictions or incorrect table and column names. *E1* noted that since the model metrics calculation does not check the values in SQL, some models often overlook this aspect, typically substituting with default values. This issue is a subsequent challenge in the *T2S* task.

E1 observed that our system significantly enhances the efficiency of analyzing model predictions through preliminary interactive exploration of extensive model prediction data. Traditional methods of analyzing model predictions involve directly comparing extensive textual information, including each token of the model's erroneously predicted SQL statements with the original data, which depends heavily on the developers' proficiency in SQL statements for accurate problem analysis. Our interactive tool greatly improves the efficiency of data analysis.

D. Case 2: Discovering Additional Cases From Existing Patterns

In this case, *E2* combined visualization and interactive techniques to explore instance-level analysis. Guided by an overview from the global-level exploration, the approach involves comparing groups of predicted and gold SQL statements to identify differences and similarities, thereby quickly discovering error patterns in SQL predictions.

E2 noted that instance-level analysis facilitates a more efficient and targeted approach to locating known errors and incorporating more cases into their model, thereby allowing users to conduct a fine-grained diagnosis of their models. In Fig. 8, an example of robustness issue mentioned in BRIDGE [1] is illustrated. Specifically, the model's inability to incorporate significant information from the utterance, even when the underlying logic is relatively straightforward, suggests that the model may have learned spurious correlations during training. In the example from BRIDGE, the model includes the “Horsepower” field in the “SELECT” clause, while the question specifically asks for “the model of the car.” *E2* hopes to uncover additional cases through instance-level analysis using the MAVIDSQL system. First, *E2* locate the “model,” “car names,” “as,” “join” cluster in the SQL comparison view, and then conduct a detailed analysis by utilizing the raw data view. For instance, in the second and third case where the model predicts “SELECT Cylinders” while the question asks for “which model.” *E2* suggest that improved modeling of compositionality in NL may help reduce such errors. This could involve modeling the span structure of the language, as well as constructing interpretable grounding with the database schema (*P1*).

E2 summarized that the SQL comparison view facilitates rapid validation of the hypotheses presented in the global-class

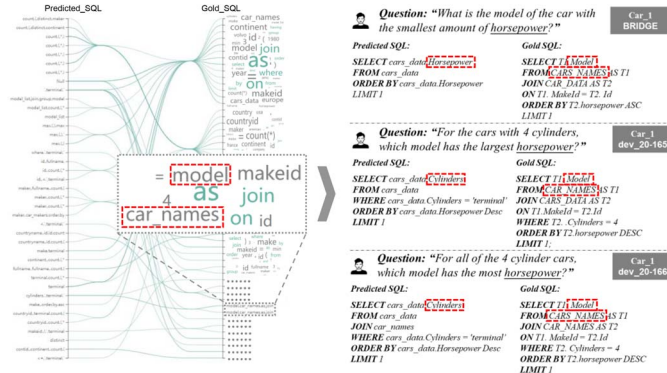


Fig. 8. Building upon the established robustness error analysis, instance-level exploration can be used to identify additional relevant cases and diagnosis models accordingly.

level through the projection and SQL similarity views. This approach provides users with an in-depth understanding of the model's strengths and limitations and can inform further refinements to the model.

Furthermore, *E2* hopes to identify issues such as lexical misunderstandings of the model and SQL irregular through this system. Specifically, he pinpoints problems within the SQL comparison view, such as finding operators namely "MIN," or instances containing numerous JOINS and aliases. This is to be followed by conducting a fine-grained analysis in combination with the raw data view.

E. User Study

We conduct a user study to further evaluate the effectiveness and usability of MAVIDSQL in facilitating the interpretation and diagnosis of T2S tasks. Moreover, we also aim to further evaluate whether the visual design can effectively improve the identification and discovery of error patterns.

1) *Participants*: we invited 12 participants (P1–P12, seven females and five males, age 22–29, $\mu = 25$, $\sigma = 2.38$) to perform the user evaluation. All 12 participants possess varied background knowledge and are currently pursuing further studies at the graduate level. P1–P5 are primarily focused on visual analysis, while P6–P8 have experience in developing T2S models. P9–P11 specialize in NLP, and participant P12 concentrates on computer vision. All participants have experience in designing and debugging deep learning models. None of the 12 participants are coauthors of this article.

2) *Procedures*: We first provided a brief introduction of the research background, including the research motivation and the exploration workflow. Subsequently, we collected demographic data from the participants and obtained their informed consent for recording their activities and outcomes for subsequent analysis. Then we introduced the system's features and demonstrated the case study mentioned above. We have designed five tasks for participants, guiding them in exploring the distribution patterns of T2S predictions and explaining the reasons for model failure. After participants concluded their exploration, we invited them to complete a postinterview questionnaire featuring five-point Likert scale questions, ranging from 1 (strongly disagree) to 5 (strongly agree), to collect their

TABLE I
USER STUDY ON EFFECTIVENESS (Q1–Q4), VISUAL DESIGN (Q5–Q9), AND USABILITY (Q10–Q13) OF MAVIDSQL. SCORE (MEAN \pm STD) FOR EACH QUESTION ARE REPORTED

	Question	Score	Score Distribution
Effectiveness			
Q1	The system can help me understand the overview of the dataset and the prediction results of Text-to-SQL tasks.	4.50 \pm 0.49	5 7
Q2	The system can help me identify the error patterns in Text-to-SQL tasks.	4.00 \pm 0.56	2 8 2
Q3	The system can help me explore model prediction at multiple levels.	3.92 \pm 0.76	4 5 3
Q4	The system can help me analyze the predictions more efficiently compared to the traditional method of directly comparing extensive textual information.	4.17 \pm 0.76	2 6 4
Visual Design			
Q5	The Projection View facilitates my comprehension of the inherent distribution characteristics in the input queries.	3.58 \pm 0.76	1 4 6 1
Q6	The Statistics View facilitates my understanding of the distribution characteristics within model prediction outcomes.	4.50 \pm 0.65	1 4 7
Q7	The Control Panel and SQL Comparison View facilitate my exploration and discovery of potential predictive error patterns.	4.00 \pm 0.71	3 6 3
Q8	The Raw Data View facilitates my validation of specific error instances.	4.41 \pm 0.64	1 5 6
Q9	The overall system is intuitive and easy to understand.	4.33 \pm 0.62	1 6 5
Usability			
Q10	It is easy to learn and use the system.	3.91 \pm 0.76	4 5 3
Q11	I think it is useful to use this system to explore the prediction of Text-to-SQL tasks.	4.00 \pm 0.82	1 8 3
Q12	I would use this system to diagnose errors of Text-to-SQL tasks in the future.	4.00 \pm 0.82	4 4 4
Q13	I would like to recommend this system to others who are working on explore Text-to-SQL tasks.	4.67 \pm 0.47	4 8

1 (Strongly Disagree) 2 3 4 5 (Strongly Agree)

Note: Score Distribution maps out the scoring outcomes for each participant.

feedback on MAVIDSQL. As illustrated in Table I, the questionnaire primarily evaluates the effectiveness, visual design, and usability of MAVIDSQL. The results and feedback have been thoughtfully summarized.

3) *Results*: Drawing upon the user ratings from the questionnaire and feedback obtained during the interviews, we evaluate the effectiveness, visual design and usability of MAVIDSQL.

a) *Effectiveness*: Most participants indicated that the system facilitated their understanding of the dataset's overview and the prediction outcomes for T2S tasks. P6 commented, "I can easily identify the primary error patterns in MAVIDSQL." As shown in Table I, participants agreed that MAVIDSQL can help users analyze the predictions more efficiently compared to the traditional method of directly comparing extensive textual information. P5 commented that the visualization is clear and insightful, allowing model developers to analyze model predictions systematically and hierarchically with MAVIDSQL, without the need to directly handle extensive textual data. This approach facilitates a more rapid analysis of model results. Moreover, several participants expressed a preference for detailed comparisons of raw data, finding that it aligned with their perception and enabled them to validate hypotheses based on practical experience.

b) *Visual design*: Based on the results of the evaluation of visual design, it is observed that the most of the

participants agreed that the overall visual design of MAVID-SQL is intuitive and easy to understand. Sometimes participants lacked familiarity with the database schema. For instance, P5 found the depiction of the projection view unclear, advocating for additional demonstration of model input queries to enhance user intuition. The participants agreed that the statistics view is the most intuitive and easy to understand, as shown in Table I. Regarding the control panel and SQL comparison view, participants believe that aggregated comparisons could help users understand the distribution of model prediction outcomes, though the learning curve is a bit steep.

c) *Usability*: The participants thought that the workflow of our system is intuitive and the interface is user-friendly. They also appreciated the system's filtering and interaction features, which facilitated a more flexible exploration of prediction instances of interest. P6 suggested displaying more information within the raw data view to support users in diagnosing model errors more effectively. Last, all the participants expressed their willingness to recommend the system to others who are working on explore T2S tasks in the future. P11 stated, "Exploring the predict and gold SQL pairs in the system was immersive, and I enjoyed the process."

F. Expert Interviews

To further evaluate the effectiveness and practicality of MAVIDSQL, we obtained feedback through individual interviews with the previously mentioned three domain experts (*E1*, *E2*, and *E3*). Prior to the interviews, none of the experts had experience with the system. We first presented an overview of the system's background and design. Subsequently, we requested the experts to engage with MAVIDSQL, exploring the model's prediction on two distinct databases. Following a 40-min exploration, we gathered their feedback regarding the system's workflow, design, application scenarios, and suggestions for improvements.

1) *System Workflow*: All the experts confirmed the effectiveness of the system workflow of MAVIDSQL in providing explanations for T2S tasks. They indicated that their typical approach to model evaluation involves relying on performance matrices and conducting instance-level analyses individually. However, this approach often lacks comprehensive details and does not facilitate in-depth support. Analyzing individual instances of model predictions is a tedious process, consuming substantial time and requiring manual summarization based on domain knowledge. Our system enhances this approach by providing both global-class level and instance-level explanations, thereby facilitating a comprehensive and systematic understanding of model predictions and identifying points where the model is prone to failure. *E1* and *E3* praised that the interactive exploration tools (i.e., lasso and brushing) are impressive and useful for analyzing T2S tasks and discovering error patterns. *E2* mentioned that if he finds the model predictions ineffective in identifying database table and column names, he might consider optimizing schema linking to enhance the model's association with the database structure. *E3* added that the aggregation in the SQL comparison group helps to generalize the model

error patterns. *E1* summarized that the system assisted users in discovering interesting insights into the models. For example, she was surprised to find that certain models unexpectedly used default values to fill in the SQL value section.

2) *Visual Design and Interactions*: Overall, the experts concurred that the visualizations are both efficacious and comprehensible, with fluent interactions. The SQL comparison view is highly favored by experts for providing a quick overview of large-scale SQL pairs. The design of the model performance statistics view is also highly regarded by the experts. *E3* really liked the SQL similarity component for providing an overview of the overall distribution of model predictions. *E1* thought the scatter plot was very intuitive, and the interactions such as lasso and brush are really helpful for the exploration of a large amount of data. Moreover, she valued how the raw data view restored the original data, enhancing the credibility of the exploration results. Nevertheless, *E1* and *E2* mentioned that in the raw data view, they still need to spend several time linking the content in the questions with the corresponding SQL. They suggested that color-coding entity fields would greatly facilitate their review of these contents.

3) *Improvements*: The experts offered constructive suggestions for improvements. During exploration in the raw data view, *E2* pointed out: "Currently, I still need to spend some time connecting the content in the question to its corresponding SQL. If entities could be mapped out with color coding, this would greatly facilitate my review of these contents." *E3* requested that in the projection view, mapping more information about queries could be implemented, as this might aid in discovering additional patterns. For a more thorough review, *E1* suggested that the system could include multimodel comparisons to help uncover more error patterns. At the same time, during the exploration by *E1* and *E2*, it was observed that some significant model errors are caused by the SQL structure. They recommended that the system should support filtering high-frequency errors.

VII. DISCUSSIONS AND FUTURE WORK

In this work, we presented a technique for understanding and exploring T2S tasks using MAVIDSQL. We preliminarily validated the effectiveness of the visual interface, MAVIDSQL, through two use scenarios and an expert interview.

MAVIDSQL can be applied to analyze various kinds of T2S models and different datasets. As a model-agnostic visual analytics workflow, the system has been designed to facilitate the discovery and diagnosis of model errors both globally and at the instance level. By focusing on the model's inputs and outputs, users can gain insight into its decision-making process. Additionally, the system's adaptable framework can be applied to various machine learning models, regardless of their specific algorithms or techniques. Despite the remarkable performance achieved by LLMs in various NL tasks, such as semantic parsing and generation, further fine-tuning is required in specific vertical domains. Moreover, the high training cost of large models hinders their flexibility in deployment. Therefore, there is still research value in exploring the development of

customizable deep learning language models for specific domain tasks, which can be quickly deployed and responsively updated. However, it may not be suitable for prior classification exploration of T2S prediction results as it is designed to facilitate free-form exploration tasks. Moreover, the system cannot directly support and apply datasets that involve multiturn conversations.

Our scalability analysis on the Spider dataset revealed that MAVIDSQL is capable of analyzing up to thousands of questions and their corresponding SQL statements. We reduce visual clutter in our system by binding SQL comparison data with the same category using a method based on sorting and category aggregation sampling. In SQL comparison view, we present a method for grouping SQL tokens with common prefixes. In the SQL comparison view, the more item a SQL group contains, the higher it is positioned on the axes. As the volume of predicted SQL data escalates, error patterns are likely to become more concentrated, with principal inaccuracies manifesting more distinctly within the views. However, there may still be issues with excessive visual clutter and cognitive burden in the word cloud and SQL similarity parallel coordinate plot visualization, despite our efforts to bind data with the same category using sorting and category aggregation.

In our work, additional error patterns not fully delineated, such as logic errors, robustness issues, and missing common-sense, as mentioned in previous works [1], [47]. If the tokens in the similar parts of the SQL structure are incorrect, they can be explored through our system. Conversely, our system still exhibits evident deficiencies in cases of significant differences in SQL structure. As mentioned earlier, the current version of MAVIDSQL only supports analysis of single-turn dialogues. When facing multiturn dialogues and complex conversations, our system may exhibit potential insufficiencies in adaptability to dynamic data and handling of complex queries. Multiturn dialogues are prevalent in many real-world applications, making the development of methods for analyzing and visualizing such dialogues an important direction for future work.

While our system has shown promising results in detecting and diagnosing errors, including patterns P1–P4 mentioned in Section VI-B2, it is imperative to recognize that there are still exist error patterns yet to be uncovered. Our system employs algorithms for SQL result alignment and SQL group comparison to identify similarities between predicted and Gold SQL structures despite differences in specific tokens. This approach aids in error pattern categorization at both class and instance levels, enhancing model refinement efforts. However, with complex query situations, such as nested queries and intricate multitable queries, the potential structural divergence between model predictions and the ground truth in SQL may be uncontrollable. This issue complicates the process of identifying error patterns through SQL comparison. In future work, semantic parsing algorithms need to be introduced to further identify the syntactical structure of SQL, thereby enhancing its capability to analyze complex model predictions. Therefore, future work should explore techniques for identifying errors in the ground truth labels and incorporating these errors into the analysis process.

In addition, experts have offered numerous valuable insights regarding visual analysis and interactive exploration which subsequently assisted in the refinement of our work. To improve the efficiency of error detection during system interaction, the introduction of error model retrieval is recommended. To enhance user understanding and alleviate over-plotting issues in PCPs, incorporating additional visual mappings in SQL comparison views is suggested. In exploring and analyzing diagnostic results of the model, it is recommended to include additional semantic hierarchical information. This can be achieved by introducing association analysis methods that are closely aligned with SQL semantics and strengthening the token link between input queries and output SQL.

Finally, we aim to improve the customization options for the projection view and SQL comparison view, allowing users to define their own metrics, including personalized similarity metrics. This will provide better insights for researchers and practitioners to comprehend and interpret these models.

VIII. CONCLUSION

In this article, we present MAVIDSQL, a visual analytics system to help model developers and users understand and diagnose T2S model prediction results. MAVIDSQL comprises four visualization components: The projection view allows for the display of multiple 2-D projections of the input sentence according to similarity summarized from different perspectives enabling users to extract potential semantic relationships. The SQL comparison view allows users to compare and analyze the differences between the model-generated SQL and the ground truth at a class-level. The model performance statistics view displays the evaluation results of the models, categorized by the performance metric ESM. The raw data view provides users with access to the raw data, allowing for detailed analysis and exploration. All four visualization components are linked together to support users in analyzing T2S tasks from multiple perspectives simultaneously and identifying common error patterns in T2S prediction results. Two case studies demonstrate the effectiveness and usability of our system T2S. While LLMs such as T5 [48] and GPT-4 [49] have shown promising performance in various tasks, the feasibility of using lightweight language models for specific NL subdomains, such as T2S, has not been fully explored. Future research in this direction can provide insights into cost-effective and rapid deployment options for such applications.

REFERENCES

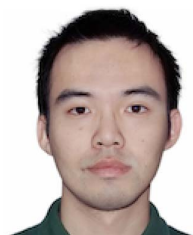
- [1] X. V. Lin, R. Socher, and C. Xiong, "Bridging textual and tabular data for cross-domain text-to-SQL semantic parsing," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2020, pp. 4870–4888.
- [2] K.-M. George and G. Koutrika, "A survey on deep learning approaches for text-to-SQL," *VLDB J.*, vol. 32, no. 4, pp. 905–936, 2023.
- [3] T. Xie et al., "UnifiedSKG: Unifying and multi-tasking structured knowledge grounding with text-to-text language models," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2022, pp. 602–631.
- [4] F. Hohman, M. Kahng, R. Pienta, and D. H. Chau, "Visual analytics in deep learning: An interrogative survey for the next frontiers," *IEEE Trans. Vis. Comput. Graph.*, vol. 25, no. 8, pp. 2674–2693, Aug. 2019.
- [5] H. Strobel, S. Gehrmann, H. Pfister, and A. M. Rush, "LSTMVis: A tool for visual analysis of hidden state dynamics in recurrent neural networks," 2017, *arXiv:1606.07461*.

- [6] H. Strobelt, S. Gehrmann, M. Behrisch, A. Perer, H. Pfister, and A. M. Rush, "Seq2Seq-Vis: A visual debugging tool for sequence-to-sequence models," 2018, *arXiv:1804.09299*.
- [7] Y. Ming et al., "Understanding hidden memories of recurrent neural networks," in *Proc. IEEE Conf. Vis. Analytics Sci. Technol.*, 2017, pp. 13–24.
- [8] S. Zhu, G. Sun, Q. Jiang, M. Zha, and R. Liang, "A survey on automatic infographics and visualization recommendations," *Vis. Inform.*, vol. 4, no. 3, pp. 24–40, 2020.
- [9] Q. Jiang, G. Sun, Y. Dong, and R. Liang, "DT2VIS: A focus+context answer generation system to facilitate visual exploration of tabular data," *IEEE Comput. Graph. Appl.*, vol. 41, no. 5, pp. 45–56, Sep./Oct. 2021.
- [10] Q. Jiang et al., "Qutaber: Task-based exploratory data analysis with enriched context awareness," *J. Visualization*, early access, 2024.
- [11] V. Zhong, C. Xiong, and R. Socher, "Seq2SQL: Generating structured queries from natural language using reinforcement learning," 2017, *arXiv:1709.00103*.
- [12] X. Xu, C. Liu, and D. Song, "SQLNet: Generating structured queries from natural language without reinforcement learning," 2017, *arXiv:1711.04436*.
- [13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics*, 2019, pp. 4171–4186.
- [14] Y. Liu et al., "RoBERTa: A robustly optimized BERT pretraining approach," 2019, *arXiv:1907.11692*.
- [15] T. Yu et al., "GraPPa: Grammar-augmented pre-training for table semantic parsing," 2021, *arXiv:2009.13845*.
- [16] P. Yin, G. Neubig, W.-t. Yih, and S. Riedel, "TaBERT: Pretraining for joint understanding of textual and tabular data," in *Proc. Annu. Meeting Assoc. Comput. Linguistics*, 2020, pp. 8413–8426.
- [17] B. Bogin, J. Berant, and M. Gardner, "Representing schema structure with graph neural networks for text-to-SQL parsing," in *Proc. Annu. Meeting Assoc. Comput. Linguistics*, 2019, pp. 4560–4565.
- [18] B. Bogin, M. Gardner, and J. Berant, "Global reasoning over database structures for text-to-SQL parsing," in *Proc. Conf. Empirical Methods Natural Lang. Process. Int. Joint Conf. Natural Lang. Process.*, 2019, pp. 3659–3664.
- [19] B. Wang, R. Shin, X. Liu, O. Polozov, and M. Richardson, "RAT-SQL: Relation-aware schema encoding and linking for text-to-SQL parsers," in *Proc. Annu. Meeting Assoc. Comput. Linguistics*, 2020, pp. 7567–7578.
- [20] W. Hwang, J. Yim, S. Park, and M. Seo, "A comprehensive exploration on WikiSQL with table-aware word contextualization," 2019, *arXiv:1902.01069*.
- [21] D. Choi, M. C. Shin, E. Kim, and D. R. Shin, "RYANSQL: Recursively applying sketch-based slot fillings for complex text-to-SQL in cross-domain databases," *Comput. Linguistics*, vol. 47, no. 2, pp. 309–332, 2021.
- [22] T. Yu et al., "SyntaxSQLNet: Syntax tree networks for complex and cross-domain text-to-SQL task," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2018, pp. 1653–1663.
- [23] J. Guo et al., "Towards complex text-to-SQL in cross-domain database with intermediate representation," 2019, *arXiv:1905.08205*.
- [24] T. Scholak, N. Schucher, and D. Bahdanau, "PICARD: Parsing incrementally for constrained auto-regressive decoding from language models," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2021, pp. 9895–9901.
- [25] A. Liu, X. Hu, L. Wen, and P. S. Yu, "A comprehensive evaluation of ChatGPT's zero-shot text-to-SQL capability," 2023, *arXiv:2303.13547*.
- [26] W. Liu, J.-L. Qiu, W.-L. Zheng, and B.-L. Lu, "Comparing recognition performance and robustness of multimodal deep learning models for multimodal emotion recognition," *IEEE Trans. Cogn. Develop. Syst.*, vol. 14, no. 2, pp. 715–729, Jun. 2022.
- [27] X. Li, Y. Hou, P. Wang, Z. Gao, M. Xu, and W. Li, "Trear: Transformer-based RGB-D egocentric action recognition," *IEEE Trans. Cogn. Develop. Syst.*, vol. 14, no. 1, pp. 246–252, Mar. 2022.
- [28] F. Cheng et al., "VBridge: Connecting the dots between features and data to explain healthcare models," *IEEE Trans. Vis. Comput. Graph.*, vol. 28, no. 1, pp. 378–388, Jan. 2022.
- [29] S. Liu, Z. Li, T. Li, V. Srikumar, V. Pascucci, and P.-T. Bremer, "NLIZE: A perturbation-driven visual interrogation tool for analyzing and interpreting natural language inference models," *IEEE Trans. Vis. Comput. Graph.*, vol. 25, no. 1, pp. 651–660, Jan. 2019.
- [30] D. Ren, S. Amershi, B. Lee, J. Suh, and J. D. Williams, "Squares: Supporting interactive performance analysis for multiclass classifiers," *IEEE Trans. Vis. Comput. Graph.*, vol. 23, no. 1, pp. 61–70, Jan. 2017.
- [31] J. Zhang, Y. Wang, P. Molino, L. Li, and D. S. Ebert, "Manifold: A model-agnostic framework for interpretation and diagnosis of machine learning models," *IEEE Trans. Vis. Comput. Graph.*, vol. 25, no. 1, pp. 364–373, Jan. 2019.
- [32] A. Hinterreiter et al., "ConfusionFlow: A model-agnostic visualization for temporal analysis of classifier confusion," *IEEE Trans. Vis. Comput. Graph.*, vol. 28, no. 2, pp. 1222–1236, Feb. 2022.
- [33] Z. Jin, X. Wang, F. Cheng, C. Sun, Q. Liu, and H. Qu, "Shortcut-Lens: A visual analytics approach for exploring shortcuts in natural language understanding dataset," *IEEE Trans. Vis. Comput. Graph.*, early access, 2023.
- [34] C. Chen et al., "OoDAnalyzer: Interactive analysis of out-of-distribution samples," *IEEE Trans. Vis. Comput. Graph.*, vol. 27, no. 7, pp. 3335–3349, Jul. 2021.
- [35] T. Nora Raju, P. A. Rahana, R. Moncy, S. Ajay, and S. K. Nambiar, "Sentence similarity—A state of art approaches," in *Proc. Int. Conf. Comput., Commun., Secur. Intell. Syst.*, 2022, pp. 1–6.
- [36] H. He and J. D. Choi, "The stem cell hypothesis: Dilemma behind multi-task learning with transformer encoders," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2021, pp. 5555–5577.
- [37] T. Yousef and S. Janicke, "A survey of text alignment visualization," *IEEE Trans. Vis. Comput. Graph.*, vol. 27, no. 2, pp. 1149–1159, Feb. 2021.
- [38] S. Jänicke, A. Geßner, G. Franzini, M. Terras, S. Mahony, and G. Scheuermann, "TRAViz: A visualization for variant graphs," *Digit. Scholarship Humanities*, vol. 30, no. 1, pp. 83–99, 2015.
- [39] B. Gipp, N. Meuschke, and J. Beel, "Comparative evaluation of text- and citation-based plagiarism detection approaches using guttenplag," in *Proc. Annu. Int. ACM/IEEE Joint Conf. Digit. Libraries*, 2011, pp. 255–258.
- [40] S. Jänicke and D. J. Wrisley, "Interactive visual alignment of medieval text versions," in *Proc. IEEE Conf. Vis. Analytics Sci. Technol.*, 2017, pp. 127–138.
- [41] G. Palmas, M. Bachynskyi, A. Oulasvirta, H. P. Seidel, and T. Weinkauff, "An edge-bundling layout for interactive parallel coordinates," in *Proc. IEEE Pacific Visualization Symp.*, 2014, pp. 57–64.
- [42] R. Zhong, T. Yu, and D. Klein, "Semantic evaluation for text-to-SQL with distilled test suites," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2020, pp. 396–411.
- [43] A. Lex, N. Gehlenborg, H. Strobelt, R. Vuilleumot, and H. Pfister, "UpSet: Visualization of intersecting sets," *IEEE Trans. Vis. Comput. Graph.*, vol. 20, no. 12, pp. 1983–1992, Dec. 2014.
- [44] S. Gratzl, A. Lex, N. Gehlenborg, H. Pfister, and M. Streit, "LineUp: Visual analysis of multi-attribute rankings," *IEEE Trans. Vis. Comput. Graph.*, vol. 19, no. 12, pp. 2277–2286, Dec. 2013.
- [45] G. Sun et al., "VSumVis: Interactive visual understanding and diagnosis of video summarization model," *ACM Trans. Intell. Syst. Technol.*, vol. 12, no. 4, pp. 1–28, 2021.
- [46] T. Yu et al., "Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2018, pp. 3911–3921.
- [47] B. Hui et al., "S²SQL: Injecting syntax to question-schema interaction graph encoder for text-to-SQL parsers," in *Proc. Assoc. Comput. Linguistics*, 2022, pp. 1254–1262.
- [48] C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, no. 1, pp. 5485–5551, 2020.
- [49] S. Bubeck et al., "Sparks of artificial general intelligence: Early experiments with GPT-4," 2023, *arXiv:2303.12712*.



Jingwei Tang received the B.E. degree in communication engineering from Zhejiang University of Technology, Hangzhou, China, in 2017, where he is currently working toward the Ph.D. degree in computer science and technology with Zhejiang University of Technology.

His research interests include data mining, visual analytics of network, and information visualization.



Guodao Sun received the B.Sc. degree in computer science and technology, in 2010 and the Ph.D. degree in control science and engineering from Zhejiang University of Technology, Hangzhou, China, in 2016.

He is a Professor with the College of Computer Science and Technology, Zhejiang University of Technology. His research interests include urban visualization, visual analytics of social media, and information visualization.



Baofeng Chang received the B.Sc. degree in automation from the College of Information Engineering in 2017 and the Ph.D. degree in computer science and technology from the College of Computer Science and Technology, Zhejiang University of Technology, Hangzhou, China, in 2023.

He is currently a Researcher with Zhejiang Airport Innovation Institute, Hangzhou, China.

His research interests include dynamic network visualization, traffic information visualization, and temporal sequence visualization.



Jiahui Chen received the B.E. degree in software engineering from the College of Computer Science and Technology, Zhejiang University of Technology, Hangzhou, China, in 2022. He is currently working toward the M.S. degree in computer technology with Zhejiang University of Technology, Hangzhou, China.

His research interests including data mining and information visualization.



Haixia Wang (Member, IEEE) received the Ph.D. degree in computer engineering from Nanyang Technological University, Singapore, in 2012.

She is currently a Professor with Zhejiang University of Technology, Hangzhou, China. Her research interests include image processing and pattern recognition.



Gefei Zhang received the B.Ed. in education technology from the College of Educational Science and Technology, Zhejiang University of Technology, Hangzhou, China. She is currently working toward the Ph.D. degree from the College of Computer Science and Technology.

Her research interests include information visualization and educational data mining.



Ronghua Liang (Senior Member, IEEE) received the Ph.D. degree in computer science from Zhejiang University, Hangzhou, China, in 2003.

He worked as a Research Fellow with the University of Bedfordshire, Luton, U.K., from 2004 to 2005, and as a Visiting Scholar with the University of California, Davis, Davis, CA, USA, from 2010 to 2011. He is currently a Professor with Zhejiang University of Technology, Hangzhou, China. His research interests include visual analytics and computer vision.