

Text-to-TrajVis: Enabling Trajectory Data Visualizations from Natural Language Questions

Anonymous ACL submission

Abstract

This paper introduces the **Text-to-TrajVis** task, which aims to transform natural language questions into trajectory data visualizations, facilitating the development of natural language interfaces for trajectory visualization systems. As this is a novel task, there is currently no relevant dataset available in the community. To address this gap, we first devised a new visualization language called Trajectory Visualization Language (TVL) to facilitate querying trajectory data and generating visualizations. Building on this foundation, we further proposed a dataset construction method that integrates Large Language Models (LLMs) with human efforts to create high-quality data. Specifically, we devised a four-stage pipeline that begins with candidate extraction, proceeds through seed TVL generation and tree-based expansion, and concludes with LLM-driven question creation followed by human validation. This process results in the creation of the first large-scale Text-to-TrajVis dataset, named **TrajVL**, which contains 9,608 (*question, TVL*) pairs. We propose a framework called **TRCAT** for progressively converting natural language questions into TVLs. The framework incorporates TVL-RAG Chain Module and Area-Time Standardization Module, significantly enhancing the accuracy of LLMs in TVL generation. Based on the TrajVL dataset, we conduct a comprehensive evaluation of TRCAT’s performance across several mainstream LLMs (e.g., GPT, Qwen, LLaMA, and Gemma). Furthermore, we established a benchmarking system for this task, providing a foundation for future research in structured trajectory language generation.

1 Introduction

The growing availability of trajectory datasets such as GeoLife+ (Amiri et al., 2024b), T-Drive (Yuan et al., 2010), and WorldTrace (Zhu et al., 2024b) provides rich spatio-temporal data for research. Visualization plays a key role in applications such

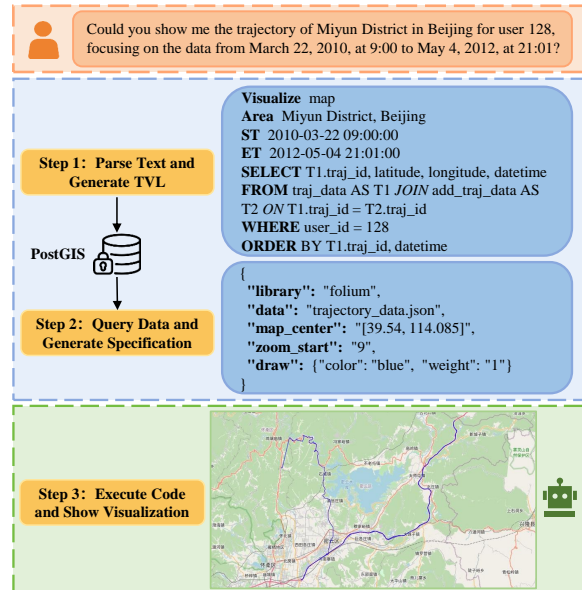


Figure 1: The process of transforming natural language questions into trajectory data visualizations.

as human mobility analysis (Toch et al., 2019; Zhu et al., 2024a), disease modeling (Amiri et al., 2024a; Kohn et al., 2023), and animal migration tracking (Konzack et al., 2019). However, existing query and visualization methods often demand specialized expertise, limiting non-expert users’ ability to extract insights from trajectory data.

Natural Language to Visualization (NL2VIS) technology focuses on the automatic translation of natural language questions (NLQs) into visualization specifications (Song et al., 2024, 2022). This capability enables non-expert users to more easily explore and visualize trajectory data, providing new opportunities to address the challenges discussed above. In this work, we introduce the Text-to-TrajVis task. Under this setting, users express trajectory visualization requests in natural language, and the Text-to-TrajVis system parses these queries to extract key semantic elements, including geographic areas, temporal constraints, and

064 query conditions. The NLQ is then transformed
065 into the Trajectory Visualization Language (TVL),
066 which serves as an intermediate representation for
067 automatically generating executable visualization
068 programs. These programs render the queried tra-
069 jectory data as either map-based visualizations or
070 statistical charts. Figure 1 illustrates an example
071 in which a natural language query is converted
072 into a map-based trajectory visualization, while the
073 transformation processes for other chart types are
074 detailed in Appendix A.1.

075 Large language models (LLMs) have shown re-
076 markable generalization capabilities in the field of
077 natural language processing. In addition to mak-
078 ing breakthroughs in traditional tasks such as text
079 generation and question-answering systems, LLMs
080 have enabled the construction of knowledge index-
081 ing libraries through Retrieval-Augmented Gen-
082 eration (RAG) (Lewis et al., 2020), significantly
083 improving the accuracy and domain adaptability of
084 generation tasks. Furthermore, with task decom-
085 position reasoning capabilities, LLMs can systemat-
086 ically break down complex tasks into well-defined
087 subtasks (Khot et al., 2022) and efficiently organize
088 and utilize retrieved knowledge. This structured
089 approach significantly enhances both the compre-
090 hension of intricate queries and the quality of the
091 final generated outputs. However, to the best of our
092 knowledge, there is currently no dedicated NL2VIS
093 dataset for trajectory data visualization that is suit-
094 able for benchmarking LLMs.

095 To address this gap, we introduce TrajVL, the
096 first benchmark dataset for the Text-to-TrajVis task.
097 We design a unified template and a systematic TVL
098 generation pipeline with 3 stages. First, seed TVLs
099 are created by instantiating the template with areas
100 and time ranges. Second, we construct a constraint
101 tree in which each leaf represents a specific query
102 constraint or combination, and progressively aug-
103 ment TVLs to increase query diversity. Finally, we
104 grouped 2–3 TVLs from the same geographic area
105 into single units, treating each as a multi-trajectory
106 visualization query. To annotate TVLs with natu-
107 ral language questions, we design two NLQ tem-
108 plates for general and complex queries involving
109 geographic and temporal expressions. NLQs are au-
110 tomatically generated using an LLM and manually
111 verified for semantic completeness and accuracy.
112 For each multi-trajectory group, a single NLQ is
113 assigned. The resulting benchmark contains 4,804
114 TVLs and 9,608 (*question, TVL*) pairs.

115 To enhance the accuracy of LLMs in generating

116 TVL, we proposed a framework named TRCAT,
117 which comprises two core modules: TVL-RAG
118 Chain module and Area-Time Standardization mod-
119 ule. The TVL-RAG Chain module incrementally
120 transforms natural language questions into struc-
121 tured TVL through a four-stage pipeline: (i) vi-
122 sualization type recognition and geographic area
123 analysis, (ii) SQL sketch construction with tem-
124 poral structure, (iii) sketch parameter filling, and
125 (iv) TVL structure integration and generation. The
126 core of this module is the RAG mechanism, which
127 retrieves semantically relevant examples based on
128 input queries, and the decomposition step. This de-
129 sign significantly improves the accuracy and struc-
130 tural consistency of the generated TVL. To further
131 improve robustness, TRCAT integrates the Area-
132 Time Standardization module. This module em-
133 ploys a multi-stage resolution strategy that com-
134 bines exact string matching and fuzzy semantic
135 matching with a gazetteer. Additionally, the mod-
136 ule leverages LLMs to parse diverse and ambiguous
137 temporal expressions. Collectively, these compo-
138 nents substantially improve the structural accuracy
139 of the generated TVLs.

140 We evaluated TRCAT against the baseline
141 method on TrajVL. Compared with the baseline
142 approach, our method improves the accuracy of
143 the best-performing model by 39.54%. Extensive
144 experiments demonstrate that our approach signifi-
145 cantly improves the performance of large language
146 models on the Text-TrajVis task.

147 In summary, our contributions are as follows:

- 148 • We are the first to introduce the Text-to-
149 TrajVis task, which aims to enable users to
150 generate trajectory data visualizations directly
151 from natural language questions.
- 152 • We construct TrajVL, the first benchmark
153 dataset specifically designed for the Text-to-
154 TrajVis task. It comprises 9,608 (*question,*
155 *TVL*) pairs.
- 156 • We propose a novel framework called TRCAT,
157 the first framework for the Text-to-TrajVis
158 task, which incrementally converts natural lan-
159 guage questions into TVLs.
- 160 • We conducted extensive experiments on
161 the proposed dataset, showing that TRCAT
162 achieves state-of-the-art performance, outper-
163 forming baseline implementation with main-
164 stream LLMs (e.g., GPT, Qwen, LLaMA, and
165 Gemma) on the Text-to-TrajVis task.

2 Related Work

2.1 Trajectory Data Visualization

Basic trajectory visualization typically represents trajectories by connecting sampled points into poly-lines overlaid on map canvases, with GeoJSON (Butler et al., 2016) serving as the standard data format. Based on this representation, Folium (Suma et al., 2024) builds on Leaflet (Crickard III, 2014) to render GeoJSON as interactive web layers, while GeoPandas (Jordahl et al., 2021) abstracts geospatial processing by integrating Shapely for topology and Cartopy for map projections. Deep learning techniques have been introduced to support higher-level semantic analysis. For instance, DeepHL (Maekawa et al., 2020) applies attention-based models to automatically label anomalous trajectory segments and enables near real-time transformation of large-scale raw data into interactive 3D visualizations. Despite these advances, most trajectory visualization tools remain technically complex and difficult for non-expert users to access.

2.2 Natural Language to Data Visualization

The field of NL2VIS has received growing attention with the advancement of visualization technologies (Cox et al., 2001; Gao et al., 2015; Liu et al., 2021; Narechania et al., 2020; Luo et al., 2021b,a). Numerous studies leverage deep neural networks to learn mappings between natural language questions and visualization specifications (Dibia and Demiralp, 2019; Song et al., 2022; Liu et al., 2023; Li et al., 2023b; Xiang et al., 2023). To support this research, several datasets have been introduced—most notably NVBench (Luo et al., 2021a), which serves as a standard benchmark, and Dial-NVBench (Song et al., 2024), which extends the task to interactive scenarios. Recently, LLM-based methods have become the dominant paradigm (Gan et al., 2021; Dong et al., 2023; Li et al., 2023a; Pourreza and Rafiei, 2023; Gao et al., 2024), benefiting from the synergy of semantic parsing and code generation. Pre-trained LLMs like Chat2vis (Maddigan and Susnjak, 2023) and ChartGPT (Tian et al., 2024) demonstrate effective prompting and task-specific fine-tuning strategies (Chung et al., 2024). Emerging multi-agent LLM systems further improve NL2VIS by distributing tasks across specialized agents. MatPlotAgent (Yang et al., 2024b) and NVAgent (Ouyang et al., 2025) enhance accuracy and interpretability through collaborative workflows, while PlotGen

```
VISUALIZE ::= map | bar | line | pie | scatter | heat map | contour plot
AREA area ∈ ({geographical_name})
ST st ∈ ({null, timestamp})
ET et ∈ ({null, timestamp}, et > st (if st, et ≠ null))
SELECT C | CC | CCC | C...C (C ::= column)
FROM T
TRANSFORM ::= avg | count | sum | group by | binning | max | min
ORDER BY ::= asc C | desc C
```

Figure 2: Structure of Trajectory Visualization Language.

(Goswami et al., 2025) and VisPath (Seo et al., 2025) integrate multimodal feedback to refine visual quality. However, the Text-to-TrajVis task remains largely unexplored. To address this gap, we propose the first systematic Text-to-TrajVis framework, accompanied by a dedicated dataset and a comprehensive evaluation on LLMs.

3 TrajVL Dataset

3.1 Trajectory Visualization Language

We designed a structured template for generating TVLs (Figure 2), which integrates four core components: (i) The VISUALIZE component defines the output visualization type and target rendering language (e.g., Python, Vega-Lite). (ii) Spatial dimensions are specified through the AREA parameter using geographical area names, while temporal dimensions employ a dual-anchor mechanism ([ST, ET]) to define closed intervals. (iii) The TRANSFORM operator implements trajectory-specific operations including temporal binning, spatial aggregation (GROUP BY), and basic statistics (AVG/COUNT/SUM/MAX/MIN), formally expressed as $T(x) \rightarrow x'$. (iv) The ORDER BY extension enables sorting based on transformed variables.

3.2 Dataset Construction

Data Preparation The trajectory data used to construct TrajVL dataset is sourced from the open-source GPS trajectory dataset released by the GeoLife project (Zheng et al., 2011). This dataset comprises 17,621 trajectories collected from 182 users, each trajectory as a time-ordered sequence of geospatial points with latitude, longitude, and altitude. The trajectories predominantly cover over 30 Chinese cities. For administrative boundary data of Chinese provinces at multiple levels, we employed Chinese area boundary data files provided by OpenStreetMap (Haklay and Weber, 2008) as the source. These boundary data were stored in

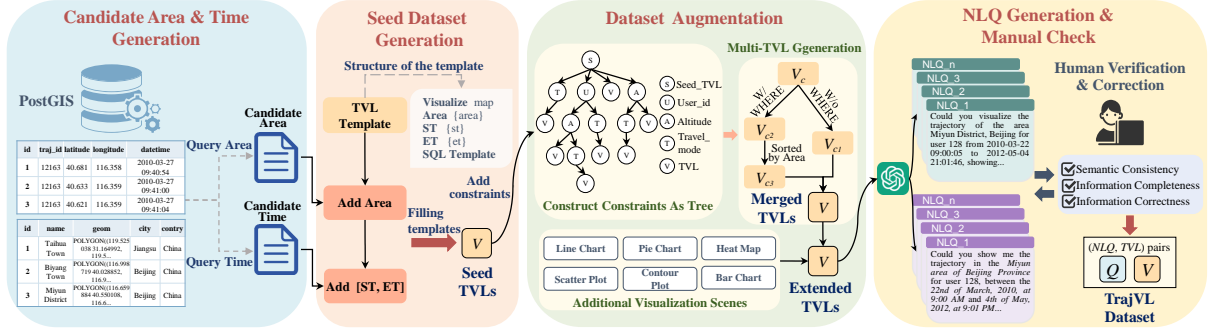


Figure 3: The pipeline of constructing the TrajVL dataset. The TrajVL dataset construction pipeline involves four main steps: (i) extraction of geographic areas and their associated temporal ranges from the trajectory database as candidates; (ii) generating representative TVLs as seed data by integrating geographic areas and temporal information into the structured template; (iii) enriching the dataset by expanding TVLs using a tree-structured augmentation approach, generating multi-trajectory queries and adding additional visualization scenarios; and (iv) generating multiple types of natural language questions using LLMs, with subsequent manual validation.

a PostGIS-enabled spatial database for efficient geospatial queries.

Seed Dataset Generation The TrajVL construction pipeline is shown in Figure 3. As the trajectory data were stored in database, we first extracted all geographic areas within the trajectory data to validate them before incorporating them into the template. These collected geographic areas constituted the area candidates $G^* = \{g_1, g_2, \dots, g_n\}$. Subsequently, we derived temporal ranges by identifying the earliest and latest timestamps of trajectories within each geographic area g_i . From each range, we selected 3-5 distinct time intervals to establish temporal candidates $T_i = \{\tau_1^{(i)}, \tau_2^{(i)}, \dots, \tau_k^{(i)}\}$. The seed TVLs are generated via combinatorial selection of validated geographic areas and temporal intervals. This process can be represented as:

$$V_{\text{seed}} = \{v_{ij} \mid g_i \in G^*, \tau_j^{(i)} \in T_i\} \quad (1)$$

Dataset Augmentation We designed a tree-structured constraint framework to enhance the expressiveness and diversity of TVLs. Beyond the fundamental latitude and longitude attributes in the primary table, the database contains auxiliary tables storing supplementary trajectory features. We systematically extracted and organized attributes including user ID, travel mode, and altitude to construct query constraints. These constraints form the hierarchy’s leaf nodes $\mathcal{L} = \{\ell_1, \ell_2, \dots, \ell_q\}$. By integrating these constraint nodes into TVL’s SQL query framework, we generated an enriched dataset. This process can be represented as:

$$V_{\text{aug}} = V_{\text{seed}} \cup \{v_{ij}^{\ell} = v_{ij} \wedge \ell \mid v_{ij} \in V_{\text{seed}}, \ell \in \mathcal{L}\} \quad (2)$$

Multi-Trajectory Generation Given a subset S , we proposed a composite TVL merging strategy $M(S)$ that combines multiple TVL queries into a single composite query containing 2-3 TVLs. This strategy follows the principle of merging only when a *WHERE* clause constraint w_i is present, ensuring that the merged natural language question retains the distinct semantics of each original trajectory. Formally, merging is performed only if all TVLs in S satisfy the following conditions: they share the same visualization type ν_i , correspond to the same geographic area $g_i = g^*$, and their corresponding SQL queries contain non-empty *WHERE* clauses. Conversely, any TVL with an empty *WHERE* clause is retained as an independent query and excluded from the merging process. This merging operation can be expressed as:

$$V_{\text{merge}} = \{M(S) \mid S \subseteq V_{\text{aug}}, \nu_i = \text{map}, g_i = g^*, w_i \neq \emptyset, \forall v_i \in S\} \cup \{v_j \in V_{\text{aug}} \mid w_j = \emptyset\} \quad (3)$$

We found that complementary trajectory features, such as travel mode and altitude, can be utilized to construct novel visual query paradigms. Consequently, variations in travel mode and altitude within trajectories were selected as primary analytical dimensions and combined them with geographic and temporal constraints to generate TVLs. This methodology was employed to generate queries for three common visualization types: pie charts, line charts, and bar charts, thereby expanding the dataset’s coverage. This dataset can be extended to support other visualization types, including scatter plots and heat maps. Formally, let V_{feat} denote the set of those TVLs. The final TVL

dataset is defined as the union of the previously merged set V_{merge} and this augmented set:

$$V_{\text{final}}^* = V_{\text{merge}} \cup V_{\text{feat}}. \quad (4)$$

Natural Language Questions Generation To reduce the cost and time overhead of manually constructing large-scale NLQs, we employ LLMs to automatically generate natural language expressions. Specifically, GPT-4o-mini is used with a temperature of 0.3 to balance generation diversity and consistency. Given the importance of geographic and temporal information in trajectory data, we adopt a sequential prompting strategy with two templates per TVL. These templates vary in area descriptions, temporal expressions, and sentence structures, enabling semantically diverse yet structurally faithful NLQs. For composite TVLs involving multiple trajectory subqueries, we design a specialized prompting strategy that explicitly extracts the geographic area, time range, and SQL filtering conditions of each subquery. These components are organized in parallel to preserve subquery independence and avoid predicate fusion, thereby ensuring that downstream systems can accurately reconstruct the original TVLs.

Manual Check To ensure the quality of NLQs generated from TVLs, we first categorized the dataset into three types based on query structure: single-trajectory, multi-trajectory, and other visualization types. A systematic human evaluation was subsequently conducted to assess the initial NLQ quality. The evaluation focused on three dimensions: semantic consistency, information completeness, and factual accuracy. Two independent annotation teams labeled separate data partitions, followed by cross-validation in which annotated sets were exchanged and re-evaluated to ensure consistency and objectivity. We employed targeted prompting strategies (see Appendix A.4) to repair faulty NLQs while preserving their semantic intent. A second round of manual verification further ensured correctness and clarity.

3.3 Dataset Analysis

The dataset includes four common visualization types: maps, bar charts, line charts, and pie charts, comprising a total of 4,804 visualization instances. Each instance is accompanied by two natural language questions, resulting in 9,608 (*question, TVL*) pairs (see Table 1). Maps account for 55.93% of the dataset, with single-trajectory queries representing 38.54% and multi-trajectory queries making up

Vis	TVL	TVL+	W/ Cons.	(NLQ, TVL)
Map	1,853	836	1,410	5,378
Bar	788	0	288	1,576
Line	727	0	727	1,454
Pie	600	0	268	1,200
Total	3,968	836	2,693	9,608

Table 1: Statistics of the TrajVL dataset. W/ Cons. denotes the TVLs with constraints, where TVL+ denotes the TVLs that support multi-trajectory query.

17.39%, thus reflecting the predominance of maps in trajectory data visualization. In addition to maps, the dataset includes visual queries such as “Display the percentage of travel by each mode in Beijing over the past year” and “Display the change in altitude of the user’s trajectory over a specified period”. To capture such diversity, we expanded the number of bar charts, line charts, and pie charts, which account for 16.39%, 15.12%, and 12.48% of the dataset, respectively. Our dataset covers 9 provinces and includes 638 regions, with the majority situated in Hebei and Beijing.

4 TRCAT Framework

4.1 Pipeline of TRCAT

For a given NLQ, TRCAT encodes the query into an embedding vector using the *all-mpnet-base-v2* model and computes cosine similarities with pre-computed embeddings of all stored NLQs. The top- K most relevant NLQs and their corresponding TVLs are retrieved to construct structured prior knowledge. This structured knowledge, along with the database schema, is subsequently passed to the LLM to generate a preliminary TVL. The Area-Time Standardization module further refines geographic area-related fields via a hybrid strategy that combines exact and fuzzy matching with a curated geographic gazetteer. It also employs LLMs to extract and normalize temporal expressions. The standardized geographic area name and temporal constraint are reintegrated to form the final TVL. Figure 4 illustrates the TRCAT pipeline.

TVL-RAG Chain This module incorporates a RAG mechanism. For each input query, the system retrieves a set of semantically similar annotated examples from the knowledge base. Rather than directly reusing these examples, the TVL-RAG Chain module extracts structured knowledge aligned with the TVL format from the retrieved items, including visualization type, geographic

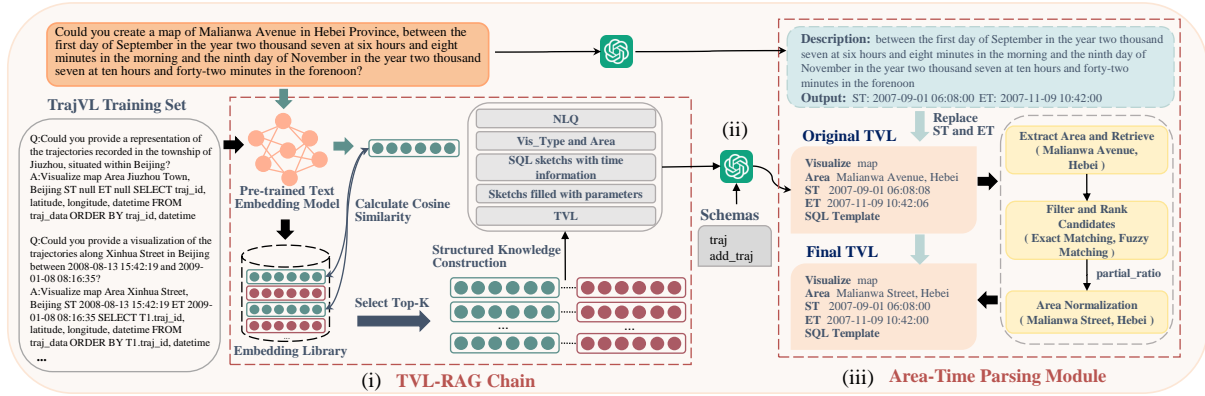


Figure 4: The working pipeline of our approach, which includes three steps: (i) Input NLQ into the retriever to obtain top- K NLQs and their associated TVLs, and generate structured knowledge based on these instances; (ii) Input structured knowledge and database schema into the LLM to generate TVL; (iii) Extract geographic area and time from the generated TVL, resolve them using Area-Time Standardization module to produce the final output.

area, abstract SQL sketches with temporal structure, parameter-filled sketches, and the complete TVL. These components are subsequently organized into a dynamic sequence of intermediate decomposition prompts, which serve as contextual demonstrations for the language model (Appendix A.5 provides complete implementation details). By leveraging these structurally consistent examples, the TVL-RAG Chain module guides LLMs to generate logically coherent and semantically faithful outputs that align with the user’s query intent.

Area-Time Standardization Module To ensure accurate interpretation of geographic area names and temporal expressions in natural language queries, we introduce an Area-Time Standardization module integrated into the TVL-RAG Chain. This module operates after TVL generation to normalize geographic and temporal fields in the structured output. Specifically, the geographic area name is first extracted from the generated TVL. It is subsequently passed to a dedicated toponym resolver, which employs both exact string matching and fuzzy semantic similarity techniques against a curated geographic gazetteer. The extracted candidate toponyms are further filtered and ranked to identify the most plausible standardized toponym. The module also supports extracting temporal descriptions from NLQs and converting them into standardized start and end timestamps. It can identify and normalize time-related expressions, including implicit and ambiguous temporal ones. The standardized geographic and temporal information is then reintegrated into the final TVL. This post-processing step improves the robustness and executability of TVLs, particularly for vague or non-

standard geographic and temporal expressions.

5 Experimental Setup

Baseline To comprehensively compare the performance of various LLMs on the TVL generation task, we selected four representative models: GPT-4o-mini (Achiam et al., 2023), LLaMA3-8B (Grattafiori et al., 2024), Qwen2.5-7B (Yang et al., 2024a), and Gemma-7B (Team et al., 2024). Among these, GPT-4o-mini is a proprietary model, while LLaMA3-8B, Qwen2.5-7B, and Gemma-7B are open-source models. All baseline models were evaluated using their respective optimal few-shot prompting strategies (Brown et al., 2020). Appendix B.1 provides specifications of the models.

Metrics To evaluate LLM performance on TVL generation, we use three common NL2VIS metrics (Luo et al., 2021a): *Vis Accuracy* (*Vis.Acc*), *Axis Accuracy* (*Axis.Acc*), and *Data Accuracy* (*Datas.Acc*). We also introduce *TVL Accuracy* (*TVL.Acc*) to evaluate TVL quality. Considering the geographic and temporal nature of trajectory data, Data Accuracy is further subdivided into Area Accuracy and Time Accuracy to evaluate geographic area and temporal dimensions. For the generated multi-trajectory TVL queries, as changes in TVL order may lead to misjudgments in exact matching methods, we adopt *F1-score* (*TVL.F1*) as the primary evaluation metric to reflect the model’s performance in multi-trajectory semantic generation tasks. Detailed metric definitions are provided in Appendix B.2.

Dataset Splitting TrajVL is split into 1,746 TVLs for training, 1,500 for validation, and 1,558 for testing. The validation set is used for few-shot parameter tuning. The test set is constructed via

Test Set	Model	Vis.Acc	Axis.Acc	Data.Acc			TVL.Acc	TVL.F1
				Area	Time	SQL		
TrajVL _{normal}	Gemma-7B	77.21%	76.12%	74.52%	76.06%	41.59%	39.41%	44.24%
	LLaMA3-8B	94.22%	93.65%	92.04%	92.30%	71.12%	68.87%	72.32%
	Qwen2.5-7B	95.70%	94.09%	94.29%	94.29%	70.80%	69.19%	72.74%
	GPT-4o-mini	97.05%	95.25%	95.70%	94.74%	77.34%	74.52%	77.32%
	Gemma-7B + Ours	86.65%	85.17%	86.20%	85.82%	67.01%	64.63%	70.85%
	LLaMA3-8B + Ours	98.10%	97.37%	97.75%	97.43%	81.51%	79.08%	81.41%
	Qwen2.5-7B + Ours	94.87%	93.13%	94.42%	94.16%	79.14%	76.77%	79.99%
	GPT-4o-mini + Ours	98.14%	95.83%	97.82%	97.50%	83.76%	81.39%	83.41%
TrajVL _{complex}	Gemma-7B	77.27%	71.31%	17.91%	61.68%	37.87%	7.25%	8.26%
	LLaMA3-8B	88.90%	88.32%	24.26%	78.75%	63.41%	15.21%	16.06%
	Qwen2.5-7B	94.35%	92.49%	26.12%	83.57%	69.13%	16.56%	18.23%
	GPT-4o-mini	95.89%	94.03%	26.83%	86.01%	74.07%	18.55%	17.83%
	Gemma-7B + Ours	83.76%	82.03%	57.64%	75.87%	64.57%	38.51%	42.59%
	LLaMA3-8B + Ours	97.43%	96.73%	76.32%	88.51%	80.36%	54.75%	57.08%
	Qwen2.5-7B + Ours	96.85%	94.87%	67.33%	87.87%	78.69%	48.27%	50.83%
	GPT-4o-mini + Ours	97.95%	95.51%	70.67%	89.28%	83.25%	53.59%	55.16%

Table 2: Performance comparison of various LLMs under different settings on the TrajVL dataset.

stratified sampling over visualization types and constraint patterns to ensure semantic and structural diversity while reducing bias toward specific query forms. To assess generalization over geographic entities, all geographic area names appearing in the test set are excluded from the training set. For each test TVL, two semantically distinct NLQs are generated, resulting in two evaluation subsets: a standard test set (TrajVL_{normal}) and a Spatiotemporal complexity set (TrajVL_{complex}), with the latter emphasizing queries involving challenging geographic and temporal expressions.

Implementation Details We employed three open-source large language models (LLaMA3-8B (Grattafiori et al., 2024), Qwen2.5-7B (Yang et al., 2024a), and Gemma-7B (Team et al., 2024)) as well as one proprietary model (GPT-4o-mini (Achiam et al., 2023)). During the inference stage, the temperature parameter for all LLMs was set to 0.1 and the max_length was set to 2048 tokens.

6 Results and Discussion

6.1 Performance Comparison

To comprehensively evaluate the performance of four LLMs and our proposed method on the TrajVL dataset, we conducted experiments on two test sets: TrajVL_{normal} and TrajVL_{complex}. The results are presented in Table 2.

On TrajVL_{normal}, the baseline GPT-4o-mini achieved the highest TVL accuracy of 74.52%. Among open-source models, Qwen2.5 slightly outperformed LLaMA3, with accuracies of 69.19%

and 68.87%, respectively, while Gemma performed substantially worse at 39.41%. Performance degraded markedly on the more challenging TrajVL_{complex}, which contains complex spatial and temporal expressions. TVL accuracy dropped by 55.97 percentage points for GPT-4o-mini, 53.66 for LLaMA3, and 32.16 for Gemma, highlighting the increased difficulty introduced by ambiguous geographic and temporal descriptions.

Applying the proposed TRCAT framework consistently improved performance across all models. On TrajVL_{normal}, GPT-4o-mini with TRCAT reached 81.39% TVL accuracy, a 6.87 percentage point gain over the baseline, while all open-source models also showed clear improvements. Notably, TRCAT enabled LLaMA3 to surpass Qwen2.5, and boosted Gemma to 64.63%, although it remained behind other models.

The gains from TRCAT were more pronounced on TrajVL_{complex}. GPT-4o-mini improved by 35.04 percentage points to 53.59% TVL accuracy. LLaMA3, Qwen2.5, and Gemma achieved improvements of 39.54, 31.71, and 31.26 percentage points, respectively. Under this setting, LLaMA3 outperformed Qwen2.5 by 6.48 percentage points and even surpassed GPT-4o-mini. This suggests that TRCAT effectively enhances LLaMA3’s ability to interpret implicit multi-trajectory structures in NLQs, leading to more accurate generation of multiple TVLs. Overall, these results demonstrate TRCAT’s robustness in handling vague and complex geographic and temporal expressions.

Beyond TVL accuracy, TRCAT outperformed

all baselines on Vis.Acc, Axis.Acc, and Data.Acc, indicating strong robustness in identifying visualization types and layout components. Axis.Acc was consistently higher on TrajVL_{normal} than on TrajVL_{complex}, due to the strong dependence of map-based visualizations on precise spatial descriptions. In addition, TRCAT improved SQL.Acc by leveraging the TVL-RAG Chain, which decomposes generation into sketch construction and parameter filling, enabling LLMs to better capture implicit SQL structures in NLQs.

6.2 Parameter Study

To examine how the number of retrieved exemplars affects the performance of TRCAT and baselines, we conducted experiments using TVL.Acc as the evaluation metric. Under the baseline setting on TrajVL_{complex}, GPT-4o-mini, Qwen2.5, LLaMA3, and Gemma achieved their best accuracies with 6-shot or 7-shot configuration. Under the TRCAT setting, GPT-4o-mini and LLaMA3 performed best with 5-shot or 6-shot configuration, Qwen2.5 and Gemma again with 3-shot configuration. Increasing the number of exemplars beyond these optimal values led to performance degradation, likely due to excessive context length hindering reasoning. All reported results correspond to each model’s best-performing few-shot setting. Accuracy curves are provided in Appendix B.3.

6.3 Ablation Study

We conduct an ablation study to assess the contribution of each component in TRCAT. The full framework with all modules enabled is first evaluated, followed by configurations with individual components removed: (i) without the TVL-RAG Chain and Area-Time Standardization modules (w/o TRC&AT); (ii) without the TVL-RAG Chain module (w/o TRC); and (iii) without the Area-Time Standardization module (w/o AT). All configurations use the same number of demonstration examples as the full TRCAT to ensure fair comparison.

The results, presented in Table 3, validate the importance of both modules in the proposed framework. We observe that the TVL-RAG Chain module plays a critical role in the Text-to-TrajVis task, as it guides LLMs to align their output structure with the TVL format, thereby enhancing the accuracy of the generated TVLs. It also enhances the LLM’s ability to interpret the implicit multi-trajectory semantics embedded in the NLQ, improving multi-TVL generation accuracy. Further-

Model	Setting	Area	TVL	TVL+
Gemma	TRCAT	57.64%	38.51%	22.14%
	- w/o TRC&AT	17.91%	7.25%	4.80%
	- w/o TRC	41.34%	17.59%	6.27%
	- w/o AT	32.80%	20.47%	15.50%
LLaMA3	TRCAT	76.32%	54.75%	53.14%
	- w/o TRC&AT	24.26%	15.21%	16.97%
	- w/o TRC	54.81%	34.85%	25.83%
	- w/o AT	54.04%	37.87%	43.91%
Qwen2.5	TRCAT	67.33%	48.27%	42.07%
	- w/o TRC&AT	26.12%	16.56%	18.08%
	- w/o TRC	43.26%	27.54%	20.3%
	- w/o AT	32.86%	23.36%	28.41%
GPT-4o-mini	TRCAT	70.67%	53.59%	55.35%
	- w/o TRC&AT	26.83%	18.55%	26.57%
	- w/o TRC	59.69%	40.82%	39.85%
	- w/o AT	40.31%	31.19%	40.96%

Table 3: Ablation study results on TrajVL_{complex}.

more, we observe that the decomposition step enables LLMs to better identify the underlying SQL structure within NLQ, thus improving the accuracy of SQL generation. The Area-Time Standardization module is essential for normalizing geographic and temporal entities. It standardizes geographic area names that LLMs may fail to recognize and converts complex temporal expressions into precise timestamps, thereby improving the accuracy of area and time parameters in TVLs.

7 Conclusion

We introduce TrajVL, the first benchmark dataset for the Text-to-TrajVis task, and design TVL, a visualization language for trajectory data querying and visualization. TrajVL is constructed through template-based seed TVL generation, tree-structured constraint augmentation, and multi-trajectory query grouping. Using LLMs for NLQ generation with human verification, we obtained 9,608 (*question*, *TVL*) pairs. To enhance TVL generation, we propose TRCAT, a framework integrating TVL-RAG Chain module and Area-Time Standardization module, which enhances accuracy and structural consistency. Experimental results demonstrate that TRCAT-enhanced LLMs outperform baseline models, especially for complex geographic and temporal expressions. Future work will aim to improve end-to-end TVL generation under complex geographic and temporal constraints, further advancing trajectory visualization research.

620 Limitations

621 We propose a dataset generation framework that
622 integrates LLMs with human intervention, and con-
623 struct the first large-scale Text-to-TrajVis bench-
624 mark dataset (dubbed TrajVL) to address the
625 scarcity of benchmark datasets in this domain.
626 However, the trajectory data used to construct the
627 dataset primarily derives from the GeoLife Project
628 Dataset. Although the dataset was effectively ex-
629 panded through a systematic approach, the insuf-
630 ficient volume of auxiliary information associated
631 with the trajectory data limited the complexity of
632 constraint trees. Thus, future research should incor-
633 porate more comprehensive auxiliary information
634 related to trajectories and construct more complex
635 constraint trees to expand the TrajVL dataset. Addi-
636 tionally, more diverse and complex spatiotemporal
637 descriptions of real-world trajectories should be
638 incorporated into the NLQ generation process to
639 enhance the dataset’s difficulty.

640 References

641 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
642 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
643 Diogo Almeida, Janko Altenschmidt, Sam Altman,
644 Shyamal Anadkat, and 1 others. 2023. *Gpt-4 techni-
645 cal report*. *arXiv preprint arXiv:2303.08774*.

646 Hossein Amiri, Ruochen Kong, and Andreas Züfle.
647 2024a. *Urban anomalies: A simulated human mobil-
648 ity dataset with injected anomalies*. In *Proceedings
649 of the 1st ACM SIGSPATIAL International Workshop
650 on Geospatial Anomaly Detection*, pages 1–11.

651 Hossein Amiri, Richard Yang, and Andreas Züfle.
652 2024b. *Geolife+: Large-scale simulated trajectory
653 datasets calibrated to the geolife dataset*. In *Pro-
654 ceedings of the 7th ACM SIGSPATIAL International
655 Workshop on GeoSpatial Simulation*, pages 25–28.

656 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
657 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
658 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
659 Askeel, and 1 others. 2020. *Language models are
660 few-shot learners*. *Advances in neural information
661 processing systems*, 33:1877–1901.

662 Howard Butler, Martin Daly, Allan Doyle, Sean Gillies,
663 Stefan Hagen, and Tim Schaub. 2016. *The geojson
664 format*. Technical report.

665 Hyung Won Chung, Le Hou, Shayne Longpre, Barret
666 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
667 Wang, Mostafa Dehghani, Siddhartha Brahma, and
668 1 others. 2024. *Scaling instruction-finetuned lan-
669 guage models*. *Journal of Machine Learning Re-
670 search*, 25(70):1–53.

Kenneth Cox, Rebecca E Grinter, Stacie L Hibino, 671
Lalita Jategaonkar Jagadeesan, and David Mantilla. 672
2001. *A multi-modal natural language interface to an 673
information visualization environment*. *International 674
Journal of Speech Technology*, 4:297–314. 675

Paul Crickard III. 2014. *Leaflet. js essentials*. Packt 676
Publishing Ltd. 677

Victor Dibia and Çağatay Demiralp. 2019. *Data2vis: 678
Automatic generation of data visualizations us- 679
ing sequence-to-sequence recurrent neural networks*. 680
IEEE computer graphics and applications, 39(5):33– 681
46. 682

Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, 683
Yunjun Gao, Jinshu Lin, Dongfang Lou, and 1 others. 684
2023. *C3: Zero-shot text-to-sql with chatgpt*. *arXiv 685
preprint arXiv:2307.07306*. 686

Yujian Gan, Xinyun Chen, Jinxia Xie, Matthew Purver, 687
John R Woodward, John Drake, and Qiaofu Zhang. 688
2021. *Natural sql: Making sql easier to infer from 689
natural language specifications*. In *Findings of the 690
Association for Computational Linguistics: EMNLP 691
2021*, pages 2030–2042. 692

Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, 693
Yichen Qian, Bolin Ding, and Jingren Zhou. 2024. 694
*Text-to-sql empowered by large language models: A 695
benchmark evaluation*. *Proceedings of the VLDB 696
Endowment*, 17(5):1132–1145. 697

Tong Gao, Mira Dontcheva, Eytan Adar, Zhicheng Liu, 698
and Karrie G Karahalios. 2015. *Datatone: Managing 699
ambiguity in natural language interfaces for data vi- 700
sualization*. In *Proceedings of the 28th annual acm 701
symposium on user interface software & technology*, 702
pages 489–500. 703

Kanika Goswami, Puneet Mathur, Ryan Rossi, and 704
Franck Dernoncourt. 2025. *Plotgen: Multi-agent 705
llm-based scientific data visualization via multimodal 706
retrieval feedback*. In *Companion Proceedings of the 707
ACM on Web Conference 2025*, pages 1672–1676. 708

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, 709
Abhinav Pandey, Abhishek Kadian, Ahmad Al- 710
Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, 711
Alex Vaughan, and 1 others. 2024. *The llama 3 herd 712
of models*. *arXiv e-prints*, pages arXiv–2407. 713

Mordechai Haklay and Patrick Weber. 2008. *Open- 714
streetmap: User-generated street maps*. *IEEE Perva- 715
sive computing*, 7(4):12–18. 716

Kelsey Jordahl, Joris Van den Bossche, Jacob Wasser- 717
man, James McBride, Jeffrey Gerard, Jeff Tratner, 718
Matthew Perry, and Carson Farmer. 2021. *geopan- 719
das/geopandas: v0. 5.0*. *Zenodo*. 720

Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao 721
Fu, Kyle Richardson, Peter Clark, and Ashish Sab- 722
harwal. 2022. *Decomposed prompting: A modular 723
approach for solving complex tasks*. In *The Eleventh 724
International Conference on Learning Representa- 725
tions*. 726

838 Eran Toch, Boaz Lerner, Eyal Ben-Zion, and Irad Ben-
839 Gal. 2019. [Analyzing large-scale human mobility](#)
840 [data: a survey of machine learning methods and](#)
841 [applications](#). *Knowledge and Information Systems*,
842 58:501–523.

843 Yanzheng Xiang, Qian-Wen Zhang, Xu Zhang, Ze-
844 jie Liu, Yunbo Cao, and Deyu Zhou. 2023. [G3r:](#)
845 [A graph-guided generate-and-rerank framework for](#)
846 [complex and cross-domain text-to-sql generation](#). In
847 *Findings of the Association for Computational Lin-*
848 *guistics: ACL 2023*, pages 338–352.

849 An Yang, Baosong Yang, Beichen Zhang, Binyuan
850 Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Day-
851 iheng Liu, Fei Huang, Haoran Wei, and 1 others.
852 2024a. Qwen2.5 technical report. *arXiv preprint*
853 *arXiv:2412.15115*.

854 Zhiyu Yang, Zihan Zhou, Shuo Wang, Xin Cong,
855 Xu Han, Yukun Yan, Zhenghao Liu, Zhixing Tan,
856 Pengyuan Liu, Dong Yu, and 1 others. 2024b. [Mat-](#)
857 [plotagent: Method and evaluation for llm-based agen-](#)
858 [tic scientific data visualization](#). In *Findings of the*
859 *Association for Computational Linguistics ACL 2024*,
860 pages 11789–11804.

861 Jing Yuan, Yu Zheng, Chengyang Zhang, Wenlei Xie,
862 Xing Xie, Guangzhong Sun, and Yan Huang. 2010.
863 [T-drive: driving directions based on taxi trajectories](#).
864 In *ACM SIGSPATIAL International Workshop on Ad-*
865 *vances in Geographic Information Systems*.

866 Yu Zheng, Hao Fu, Xing Xie, Wei Ying Ma, and Quan-
867 nan Li. 2011. [Geolife gps trajectory dataset - user](#)
868 [guide](#).

869 Feng Zhu, Chen Chang, Zhiheng Li, Boqi Li, and Li Li.
870 2024a. [A generic optimization-based enhancement](#)
871 [method for trajectory data: Two plus one](#). *Accident;*
872 *analysis and prevention*, 200:107532.

873 Yuanshao Zhu, James Jianqiao Yu, Xiangyu Zhao, Xue-
874 tao Wei, and Yuxuan Liang. 2024b. [Unitraj: Learn-](#)
875 [ing a universal trajectory foundation model from](#)
876 [billion-scale worldwide traces](#). *CoRR*.

A Implementation Details

A.1 Visualization Generation Examples

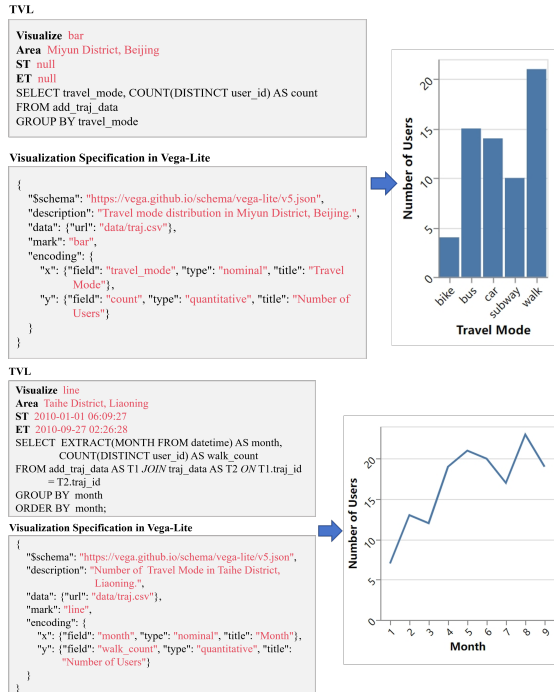


Figure 5: Examples of the visualization specification in Vega-Lite, and their corresponding visualizations.

The proposed Text-to-TrajVis system implements an automatic mapping mechanism from natural language questions to TVLs. As shown in Figure 5, the generated TVL contains three core components: (i) visualization type statement: specifying the chart presentation form; (ii) geographic and temporal constraint parameters: constructing geographic and temporal filtering conditions; and (iii) SQL query statement: realizing semantic fusion of trajectory attributes and geographic and temporal coordinates through multi-table association. For standardized chart types, including but not limited to bar charts, line charts, and pie charts, the system employs Vega-Lite, a state-of-the-art declarative visualization specification framework. By leveraging the widely adopted JSON syntax, Vega-Lite allows for the precise specification of visual coding rules, such as mapping data attributes to visual channels (e.g., position, color, size) and defining the overall chart structure. The proposed architecture ensures end-to-end interpretability from semantic parsing to visual presentation through formal language transformation mechanisms.

TVL

```
Visualize map
Area Miyun District, Beijing
ST 2010-03-22 09:00:05
ET 2012-05-04 21:01:46
SELECT T1.traj_id, latitude, longitude, datetime
FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id
WHERE user_id = 128
ORDER BY T1.traj_id, datetime
```

Synthesized SQL

```
SELECT T1.traj_id, latitude, longitude, datetime
FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id
WHERE ST_Within(ST_GeomFromText('POINT(' || longitude || ' ' || latitude || ')', 4326), (SELECT geom FROM boundaries WHERE city='Beijing' AND name='Miyun District'))='true'
AND datetime BETWEEN '2010-03-22 09:00:05' AND '2012-05-04 21:01:46'
AND user_id = 128
ORDER BY T1.traj_id, datetime
```

Figure 6: An example of synthesizing SQL based on TVL.

A.2 Synthesizing SQL from TVL

Given the TVL specification, the system achieves translation to SQL by parsing the logical structure inherent in TVL. During this conversion, geographic and temporal constraints declared within the TVL are dynamically mapped to the WHERE clause of the resulting SQL query, specifically encompassing the core parameters: Area, ST (Start Time), and ET (End Time). For instance, when TVL specifies that the value of Area is "Miyun District, Beijing", ST is "2010-03-22 09:00:00", ET is "2012-05-04 21:01:00", and the basic SQL structure is "WHERE ST_Within(ST_GeomFromText('POINT(' || longitude || ' ' || latitude || ')', 4326), (SELECT geom FROM boundaries WHERE city='Beijing' AND name='Miyun District'))='true' AND datetime BETWEEN '2010-03-22 09:00:05' AND '2012-05-04 21:01:46' ". Using these parameters, the system is able to dynamically generate the SQL required to query the data, as shown in Figure 6.

A.3 Prompts for Generating NLQs

Prompt for generating multi-trajectory queries

Please write a single NLQ whose intent covers all of the above TVLs.

Special Rules

1. The resulting NLQ must be fluent, match the user's tone and consistent with the real world.
2. The NLQ is a fluent user query in one sentence, don't appear to be listed in separate entries.
3. The natural language question must contain the intent of all TVLs.
4. The NLQ must accurately preserve all distinct constraints from each TVL, the

939 same constraints can be merged.
 940 5. The goal is for a downstream system to
 941 reconstruct all original TVLs from a single
 942 structured NLQ, without merging or
 943 confusing their individual logic.
 944 6. The geographic area name and time in the
 945 natural language question need to be
 946 consistent with the TVL and do not need to
 947 be described differently.
 948
 949 Here are some examples:
 950 TVL 1: Visualize map Area Baiyun Street,
 951 Shandong ST 2008-08-21 12:42:33 ET
 952 2009-09-12 01:50:33 SELECT T1.traj_id,
 953 latitude, longitude, datetime FROM
 954 traj_data AS T1 JOIN add_traj_data AS T2 ON
 955 T1.traj_id = T2.traj_id WHERE user_id =
 956 140 ORDER BY T1.traj_id, datetime;
 957 TVL 2: Visualize map Area Baiyun Street,
 958 Shandong ST 2008-08-21 12:42:33 ET
 959 2009-09-12 01:50:33 SELECT T1.traj_id,
 960 latitude, longitude, datetime FROM
 961 traj_data AS T1 JOIN add_traj_data AS T2 ON
 962 T1.traj_id = T2.traj_id WHERE user_id =
 963 180 AND altitude BETWEEN -9725.41 AND
 964 102482.8 ORDER BY T1.traj_id, datetime;
 965 Generated NLQ:
 966 Please show me the trajectory data for the
 967 area Baiyun Street, Shandong from
 968 2008-08-21 12:42:33 to 2009-09-12 01:50:33,
 969 including trajectories of user 140, and
 970 those of user 180 with altitude between
 971 -9725.41 and 102482.8?
 972
 973 Now please generate the NLQ based on the
 974 given TVLs:
 975
 976 TVL: {tv1}
 977 Generated NLQ:

Prompt for generating normal data

978 Please generate a NLQ for the given TVL.
 979
 980
 981 ## Special Rules
 982 1. The resulting NLQ must be fluent, match
 983 the user's tone and consistent with the
 984 real world.
 985 2. The natural language question must
 986 contain the intent of all TVLs.
 987 3. The geographic area name and time in the
 988 natural language question need to be
 989 consistent with the TVL and do not need to
 990 be described differently.
 991
 992 Now please generate the NLQ based on the
 993 given TVL:
 994 TVL: {tv1}
 995 Generated NLQ:

Prompt for different descriptions of geographic area and time

996 You are a helpful assistant to complicate
 997 the description of area name and time
 998 strictly according to the user's
 999 requirements.
 1000
 1001 ## Requirements:
 1002 1. First, you must ensure that you do not

change the description of any information
 in the base NLQ other than the area name
 and time.
 2. Complicate the area name descriptions in
 a way that is different from the base NLQ,
 but make sure they still represent the
 same area name.
 3. If there is time information in the base
 NLQ, complicate the time descriptions in a
 way that is different from the base NLQ,
 but ensure that they still represent the
 same time.
 4. To ensure that you do not modify the
 information in the base NLQ, it is
 recommended that you first extract the area
 name in the base NLQ, complicate the area
 name, and then replace the area name in the
 base NLQ.
 Base NLQ:
 {nlq}
 Modified NLQ:

A.4 Prompts for Correcting NLQs

Prompt for correcting semantic redundancy

1027 There is a problem with the description of
 the NLQ in the example data, the generated
 description contains information that is
 not relevant to the TVL. Please examine the
 original NLQ carefully, simplify that NLQ
 and save the information describing only
 the individual fields in the TVL, correct
 it, and then regenerate a new NLQ to ensure
 that the information described is not
 redundant (Note: the answer just outputs
 each NLQ without additional text).
 Original NLQ:
 {nlq}
 Revised NLQ:

Prompt for correcting information missing

1044 There is a problem with the description of
 the NLQ in the example data, the
 description information generated based on
 the TVL is incomplete with missing
 information. Please regenerate a new NLQ
 based on the original NLQ and make sure
 that the description information is
 complete with all the information in the
 TVL (Note: the answer just outputs each NLQ
 without additional text).
 Original NLQ:
 {nlq}
 Revised NLQ:

Prompt for correcting information error

1059 There is a problem with the description of
 the NLQ in the example data, the generated
 description information is inconsistent
 with the TVL in fields such as time, area,
 and SQL. Please examine the original NLQ
 carefully, find the inconsistent
 description, fix it, and regenerate a new

1067	NLQ to ensure that the information in the	the final TVL.	1133
1068	description is consistent with the	Final TVL:	1134
1069	information in the TVL (Note: the answer	TVL1:	1135
1070	just outputs each NLQ without additional	Visualize map Area Jingshan Street, Beijing	1136
1071	text).	ST 2008-04-13 09:45:34 ET 2010-01-04	1137
1072		05:44:45 SELECT T1.traj_id, latitude,	1138
1073	Original NLQ:	longitude, datetime FROM traj_data AS T1	1139
1074	{nlq}	JOIN add_traj_data AS T2 ON T1.traj_id = T2.	1140
1075	Revised NLQ:	traj_id WHERE user_id = 17 AND travel_mode	1141
		= 'walk' ORDER BY T1.traj_id, datetime;	1142
		TVL2:	1143
		Visualize map Area Jingshan Street, Beijing	1144
		ST 2008-04-13 09:45:34 ET 2010-01-04	1145
		05:44:45 SELECT T1.traj_id, latitude,	1146
		longitude, datetime FROM traj_data AS T1	1147
		JOIN add_traj_data AS T2 ON T1.traj_id = T2.	1148
		traj_id WHERE user_id = 2 AND altitude	1149
		BETWEEN -12875.79 AND 5657.95 ORDER BY T1.	1150
		traj_id, datetime;	1151
1076	A.5 Implementation Details of TVL-RAG		
1077	Chain		
1078	## Question		
1079	Please provide the trajectory data for the		
1080	area Jingshan Street, Beijing from		
1081	2008-04-13 09:45:34 to 2010-01-04 05:44:45,		
1082	including trajectories of user 17 on foot		
1083	and those of user 2 within an altitude		
1084	range between -12875.79 and 5657.95.		
1085			
1086	Decompose the task into sub tasks,		
1087	considering [Background] [Special Rules] [
1088	Constraints], and generate the TVL after		
1089	thinking step by step:		
1090			
1091	**Sub task 1:** Generate visualize type and		
1092	area name		
1093	By parsing the Question, we generated the		
1094	visualize type and standardized area name.		
1095	Visualize type: map, area name: Jingshan		
1096	Street, Beijing		
1097			
1098	**Sub task 2:** Construct sketches		
1099	By parsing the question, we found that 3		
1100	specific sketches with temporal structure		
1101	need to be constructed for the Question.		
1102	Sketch 1:		
1103	ST _ ET _ SELECT _ , _ , _ , _ FROM _ JOIN		
1104	_ ON _ WHERE _ ORDER BY _		
1105	Sketch 2:		
1106	ST _ ET _ SELECT _ , _ , _ , _ FROM _ JOIN		
1107	_ ON _ WHERE _ ORDER BY _		
1108			
1109	**Sub task 3:** Fill sketches		
1110	By parsing the question, we extracted		
1111	relevant information and filled them in the		
1112	sketches to generate 3 complete SQL		
1113	queries with temporal information.		
1114	SQL 1:		
1115	ST 2008-04-13 09:45:34 ET 2010-01-04		
1116	05:44:45 SELECT T1.traj_id, latitude,		
1117	longitude, datetime FROM traj_data AS T1		
1118	JOIN add_traj_data AS T2 ON T1.traj_id = T2.		
1119	traj_id WHERE user_id = 17 AND travel_mode		
1120	= 'walk' ORDER BY T1.traj_id, datetime;		
1121	SQL 2:		
1122	ST 2008-04-13 09:45:34 ET 2010-01-04		
1123	05:44:45 SELECT T1.traj_id, latitude,		
1124	longitude, datetime FROM traj_data AS T1		
1125	JOIN add_traj_data AS T2 ON T1.traj_id = T2.		
1126	traj_id WHERE user_id = 2 AND altitude		
1127	BETWEEN -12875.79 AND 5657.95 ORDER BY T1.		
1128	traj_id, datetime;		
1129			
1130	**Sub task 4:** Combine into Final TVL		
1131	Add visualize type and area name to the SQL		
1132	queries with time information to generate		
		For a given NLQ, the TVL-RAG Chain first re-	1152
		trieves the top-k most similar NLQs from an em-	1153
		bedding library and uses each retrieved NLQ as	1154
		the Question in an in-context example. Since ev-	1155
		ery NLQ is paired with a TVL, we decompose the	1156
		corresponding TVL into structured components:	1157
		(i) visualization type and geographic area name,	1158
		(ii) SQL sketches with temporal structure, (iii)	1159
		information-filled sketches, and (iv) the final TVLs.	1160
		These structured elements guide the LLM to inter-	1161
		nalize the target TVL format and generate correct	1162
		outputs step by step. All constructed examples are	1163
		then embedded into the prompt and fed to the LLM.	1164
		B Experimental Details	1165
		B.1 Baselines	1166
		We evaluated TRCAT and baseline methods on	1167
		LLMs using the TrajVL dataset, with the baseline	1168
		methods and LLMs listed as follows:	1169
		• Few-shot LLM: The few-shot prompting is a	1170
		core mechanism in In-Context Learning, us-	1171
		ing a limited set of examples to guide LLMs in	1172
		performing tasks within specialized domains.	1173
		• LLaMA3-8B: LLaMA3-8B, an 8-billion-	1174
		parameter open-weight language model devel-	1175
		oped by Meta, is optimized for high efficiency	1176
		while maintaining robust performance across	1177
		diverse NLP tasks. Pretrained on a large-scale	1178
		multilingual corpus, it demonstrates strong	1179
		reasoning and text-generation capabilities, po-	1180
		sitioning it as a versatile choice for applica-	1181
		tions requiring a balance between model size	1182
		and inference speed.	1183
		• Qwen2.5-7B: Qwen2.5-7B, a 7-billion-	1184
		parameter large language model developed	1185

by Alibaba Group, is engineered to deliver competitive performance in both general and domain-specific tasks. Trained on diverse multilingual and multimodal corpora, it excels in natural language understanding, generation, and code-related tasks, while maintaining computational efficiency.

- **Gemma-7B**: Gemma-7B, an 7-billion-parameter large language model developed by Google’s Gemma family, a series of lightweight, state-of-the-art open-source models built on the same research and technologies used to develop the Gemini models. As text-to-text decoder-based large language models, Gemma models are available in English, with both open-weight pre-trained and instruction-tuned variants. They are well-suited for a variety of text generation tasks, including question answering, summarization, and reasoning. Thanks to their relatively small size, Gemma models are also optimized for deployment in resource-constrained environments.
- **GPT-4o-mini**: GPT-4o-mini, a smaller version in OpenAI’s GPT-4o series, leverages advanced architectures and training methodologies for strong performance on diverse NLP tasks. Designed for efficiency and responsiveness, it offers a balance between high-quality generation and reduced latency, making it suitable for applications requiring both sophisticated language understanding and rapid response times.

B.2 Metrics

We use the same metrics as NVBench, including Vis Accuracy, Axis Accuracy, and Data Accuracy, to assess LLM performance in TVL generation. Since geographic and temporal information is important for trajectory data, SQL statements for trajectory queries are synthesized based on geographic area, time and basic SQL query structure. Consequently, Data Accuracy is decomposed into Area Accuracy, Time Accuracy, SQL Accuracy and TVL Accuracy. For the generated multi-trajectory TVL queries, we adopt F1-score as the primary evaluation metric.

The evaluation metrics are formally defined as follows:

- **Vis Accuracy (Vis.Acc)**: This metric evaluates the accuracy of the Visualize component

of the TVL generated by the model, measuring the model’s ability to identify the type of visualization. It is calculated as:

$$Acc_{type} = \frac{N_{type}}{N}$$

Where N_{type} represents the count of exact match visualization types, N represents the total number of visualization types in the test set.

- **Axis Accuracy (Axis.Acc)**: This metric evaluates the performance of the model in recognizing vis components. For map visualizations, it evaluates the accuracy of the Area generated by the model, since proper visualization requires centering of the specified geographic area for map visualizations. For the x-axis and y-axis components of other types of visualizations such as bar, line, pie, etc., we measure the Select component of the vis query. It is calculated as:

$$Acc_{comp} = \frac{N_{comp}}{N}$$

Where N_{comp} represents the count of exact matches to the Area or Select component, N represents the total number of queries in the test set.

- **Area Accuracy (Area.Acc)**: This metric evaluates the accuracy of the Area component of the TVLs generated by the model to reveal the model’s performance in recognizing geographical area names. It is calculated as:

$$Acc_{area} = \frac{N_{area}}{N}$$

Where N_{area} represents the count of exact matching Area components, N represents the total number of queries in the test set.

- **Time Accuracy (Time.Acc)**: This metric evaluates the accuracy of the Time component of the TVLs generated by the model to reveal the model’s performance in recognizing temporal information. It is calculated as:

$$Acc_{time} = \frac{N_{time}}{N}$$

Where N_{time} represents the count of exact matching Time components, N represents the total number of queries in the test set.

- **SQL Accuracy (SQL.Acc):** This metric evaluates the accuracy of the model in generating critical SQL by parsing and comparing SQL query structures in TVL. Specifically, it eliminates interference from syntactically equivalent but ordering differences through normalization (e.g., alphabetizing logical conditions in WHERE clauses), and recursively compares SQL structures through parse trees. It is calculated as:

$$Acc_{sql} = \frac{N_{sql}}{N}$$

Where N_{sql} represents the count of exact matching SQL queries, N represents the total number of SQL queries in the test set.

- **TVL Accuracy (TVL.Acc):** This metric evaluates the accuracy of the model in generating a complete TVL. Based on the structure of the TVL, a model is considered to be able to generate the TVL correctly when it correctly generates Visualize, Area, Time and SQL. It is calculated as:

$$Acc_{tvl} = \frac{N_{tvl}}{N}$$

Where N_{tvl} represents the count of exact match TVLs, N represents the total number of TVLs in the test set.

- **F1 Score (TVL.F1):** This metric addresses the challenge of evaluating multi-trajectory TVL queries, where changes in TVL order can cause misjudgments in traditional exact matching methods. We adopt F1-score as the primary evaluation metric to comprehensively reflect the model’s performance in multi-trajectory semantic generation tasks. It is calculated as:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$

- **True Positive (TP):** The number of generated TVLs that correctly match the target TVLs.
- **False Positive (FP):** The number of generated TVLs that do not match the target TVLs.
- **False Negative (FN):** The number of target TVLs that are not matched by the generated TVLs.

B.3 Experimental Details

Under the baseline setting on TrajVL_{normal}, GPT-4o-mini, Qwen2.5, and LLaMA3 achieved their best accuracies with 6-shot or 7-shot configuration, whereas Gemma peaked at 3-shot configuration. On the TrajVL_{complex}, all four LLMs achieved their best accuracies with 6-shot or 7-shot configuration. Under the TRCAT setting on TrajVL_{normal}, GPT-4o-mini and LLaMA3 performed best with 5-shot or 6-shot configuration, Qwen2.5 and Gemma again with 2-shot configuration. On the TrajVL_{complex}, GPT-4o-mini and LLaMA3 performed best with 5-shot or 6-shot configuration, Qwen2.5 and Gemma again with 3-shot configuration. When the number of retrieval examples exceeds the optimal value for each LLM, excessively long context lengths may have a negative impact on the model’s performance.

B.4 Experimental Prompt for Few-shot Configuration

Given a [Database schema] and a [Question], generate valid TVL sentences.

```
## Background
### TVL Structure:
Visualize [visualize] Area [area] ST [start time] ET [end time] SELECT [COLUMNS] FROM [TABLES] [JOIN] [WHERE] [GROUP BY] [ORDER BY]
```

```
### Key Components:
1. Visualization: The type of visual chart that users expect.
2. Area: Extracting geographic area information.
3. Time: Extracting time information. The time format is: XXXX-XX-XX XX:XX:XX. (ST: Start Time, ET: End Time. If no time is specified, set to "null")
4. SQL Components: SELECT, FROM, JOIN, WHERE, GROUP BY, ORDER BY.
You must consider which part in the SQL components is necessary, which is unnecessary.
```

When generating TVL, we should always consider special rules and constraints:

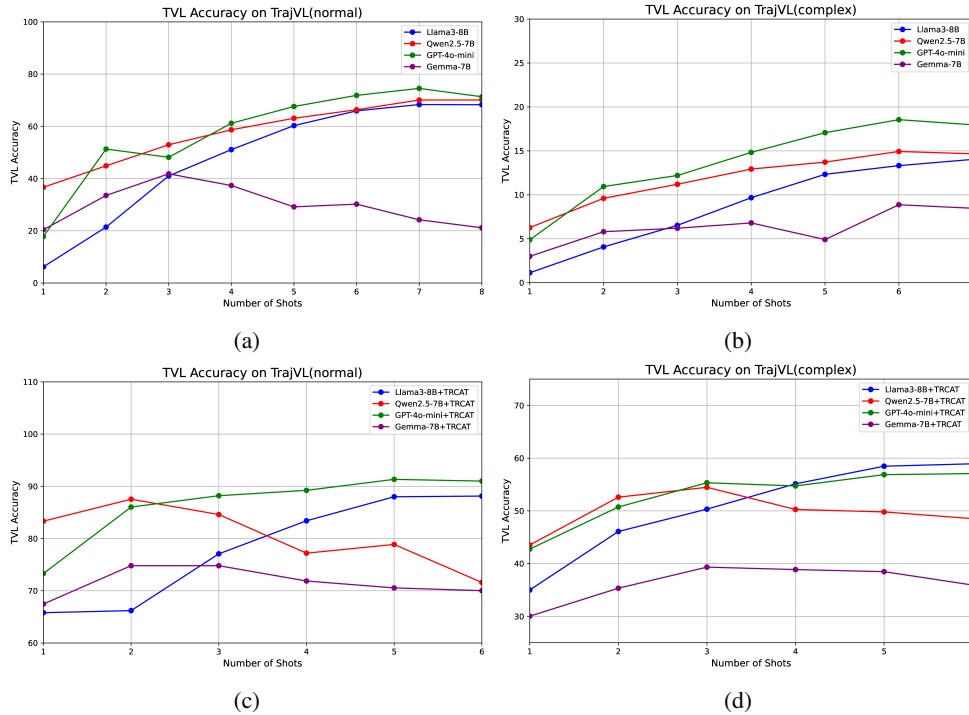


Figure 7: Parameter study of Baseline on TrajVL. The vertical axis denotes TVL.Acc, and the horizontal axis denotes the number of retrieved exemplars.

```

1371 ## Special Rules
1372 a. For simple visualizations (bar, line,
1373 pie):
1374     - SELECT exactly TWO columns, X-axis and
1375       Y-axis.
1376 b. Aggregate Functions:
1377     - Use COUNT for counting occurrences
1378     - Use SUM only for numeric columns
1379     - When in doubt, prefer COUNT over SUM
1380     - When extract the month from datetime,
1381       use EXTRACT
1382
1383 ## Constraints
1384     - In SELECT <column>, make sure there are
1385       at least two selected.
1386     - Use only table names and column names
1387       from the given database schema.
1388     - Ensure GROUP BY precedes ORDER BY for
1389       distinct values.
1390
1391 Here are some examples:
1392 ## Database schema
1393 ### Table: traj_data
1394 [
1395     (id, Value examples: [20038172, 1, 2, 3].
1396     And this is an id type column),
1397     (traj_id, Value examples: [1522, 12917,
1398     14059, 8149]). And this is an id type
1399     column),
1400     (latitude, Value examples: [39.97539,
1401     39.97535, 39.97525, 39.97515].),
1402     (longitude, Value examples: [116.32769,
1403     116.327375, 116.327325, 116.45166].),
1404     (datetime, Value examples: ['2011-08-30
1405     00:51:53']).)
1406 ]
1407 ### Table: add_traj_data

```

```

[
(id, Value examples: [20038172, 1, 2, 3].
And this is an id type column),
(traj_id, Value examples: [1522, 12917,
14059, 8149]). And this is an id type
column),
(altitude, Value examples: [0.0, 164.0,
128.0, 157.5].),
(travel_mode, Value examples: ['walk', '
bus', 'bike', 'car', 'train', 'subway'].),
(user_id, Value examples: [128, 41, 17,
68]). And this is an id type column)
]
## Question
{nl_query}
## Final TVL:
{TVL}

Now here is a new question:
## Database Schema
{database_schema}

## Question
{query}

Now, please generate the **Final TVL** for
the database schema and question.
(stop answering after you give the final
ttl)

```

B.5 Experimental Prompt for TRCAT Framework

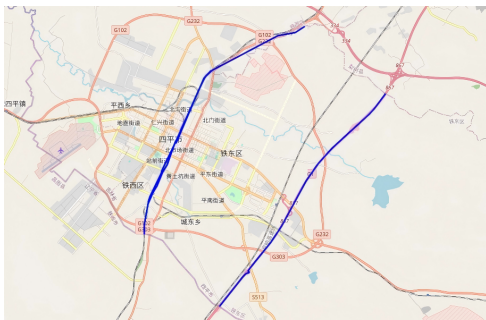
```

Given a [Database schema] and a [Question],
generate valid TVL sentences.

```

1442	## Background	116.327375, 116.327325, 116.45166.),	1512
1443	### TVL Structure:	(datetime, Value examples: ['2011-08-30	1513
1444	Visualize [visualize] Area [area] ST [start	00:51:53'].)	1514
1445	time] ET [end time] SELECT [COLUMNS] FROM]	1515
1446	[TABLES] [JOIN] [WHERE] [GROUP BY] [ORDER	### Table: add_traj_data	1516
1447	BY]	[1517
1448		(id, Value examples: [20038172, 1, 2, 3].	1518
1449	### Key Components:	And this is an id type column),	1519
1450	1. Visualization: The type of visual chart	(traj_id, Value examples: [1522, 12917,	1520
1451	that users expect.	14059, 8149]. And this is an id type	1521
1452	2. Area: Extracting geographic area	column),	1522
1453	information.	(altitude, Value examples: [0.0, 164.0,	1523
1454	3. Time: Extracting time information. The	128.0, 157.5].),	1524
1455	time format is: XXXX-XX-XX XX:XX:XX. (ST:	(travel_mode, Value examples: ['walk', '	1525
1456	Start Time, ET: End Time. If no time is	bus', 'bike', 'car', 'train', 'subway'].),	1526
1457	specified, set to "null")	(user_id, Value examples: [128, 41, 17,	1527
1458	4. SQL Components: SELECT, FROM, JOIN,	68]. And this is an id type column)	1528
1459	WHERE, GROUP BY, ORDER BY.]	1529
1460	You must consider which part in the SQL	## Question	1530
1461	components is necessary, which is	{result['Matched NL Query'][i]}	1531
1462	unnecessary.		1532
1463			1533
1464	When generating TVL, we should always	Decompose the task into four sub-tasks,	1534
1465	consider special rules and constraints:	considering [Background] [Special Rules] [1535
1466	## Special Rules	Constraints], and generate the TVL after	1536
1467	a. For simple visualizations (bar, line,	thinking step by step:	1537
1468	pie):		1538
1469	- SELECT exactly TWO columns, X-axis and	**Sub task 1:** Generate visualize type and	1539
1470	Y-axis.	area name	1540
1471	b. Aggregate Functions:	By parsing the Question, we generated the	1541
1472	- Use COUNT for counting occurrences	visualize type and standardized area name.	1542
1473	- Use SUM only for numeric columns	{result['structure_fields'][i]}	1543
1474	- When in doubt, prefer COUNT over SUM		1544
1475	- When extract the month from datetime,	**Sub task 2:** Construct sketches	1545
1476	use EXTRACT	By parsing the question, we found that {	1546
1477		result['tv_num'][i]} specific sketches	1547
1478	## Constraints	with temporal structure need to be	1548
1479	- In SELECT <column>, make sure there are	constructed for the Question.	1549
1480	at least two selected.	{result['TVL Sketch'][i]}	1550
1481	- Use only table names and column names		1551
1482	from the given database schema.	**Sub task 3:** Fill sketches	1552
1483	- Ensure GROUP BY precedes ORDER BY for	By parsing the question, we extracted	1553
1484	distinct values.	relevant information and filled them in the	1554
1485		sketches to generate {result['tv_num'][i]}	1555
1486	Now we could think step by step:	complete SQL queries with temporal	1556
1487	1. First choose visualize type and area	information.	1557
1488	name.	{result['SQL'][i]}	1558
1489	2. Second construct a sketch with temporal		1559
1490	structure for the natural language question.	**Sub task 4:** Combine into Final TVL	1560
1491		Add visualize type and area name to the SQL	1561
1492	3. Third extract the relevant information	queries with time information to generate	1562
1493	and fill in the sketched to generate	the final TVL.	1563
1494	complete SQL queries with temporal	Final TVL:	1564
1495	information.	{result['TVL'][i]}	1565
1496	4. Fourth add visualize type and area name		1566
1497	to the SQL queries with time information to	Now here is a new question:	1567
1498	generate the final TVL.	## Database Schema	1568
1499		{database_schema}	1569
1500	Here are some examples:		1570
1501	## Database Schema	## Question	1571
1502	### Table: traj_data	{query}	1572
1503	[1573
1504	(id, Value examples: [20038172, 1, 2, 3].	Now, please generate the **Final TVL** for	1574
1505	And this is an id type column),	the database schema and question after	1575
1506	(traj_id, Value examples: [1522, 12917,	thinking step by step.	1576
1507	14059, 8149]. And this is an id type	(stop answering after you give the final	1577
1508	column),	tvL)	1578
1509	(latitude, Value examples: [39.97539,		
1510	39.97535, 39.97525, 39.97515].),		
1511	(longitude, Value examples: [116.32769,		

NLQ	Show the trajectory data for the Tiedong District, Liaoning Province from 5:10:50 AM on February 28, 2009, to 10:37:04 AM on April 21, 2010, including all trajectories with an altitude between -7984.86 and 54056.95, and the trajectories of user 71 with an altitude between -16987.76 and 22297.35 from 4:42:47 AM on September 10, 2009, to 8:18:57 AM on January 2, 2012.
Target TVL	TVL1: Visualize map Area Tiedong District, Liaoning ST 2009-02-28 05:10:50 ET 2010-04-21 10:37:04 SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE altitude BETWEEN -7984.86 AND 54056.95 ORDER BY T1.traj_id, datetime TVL2: Visualize map Area Tiedong District, Liaoning ST 2009-09-10 04:42:47 ET 2012-01-02 08:18:57 SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE user_id = 71 AND altitude BETWEEN -16987.76 AND 22297.35 ORDER BY T1.traj_id, datetime
Few-shot	TVL1: Visualize map Area Tiedong District, Liaoning ST 2009-02-28 05:10:50 ET 2010-04-21 10:37:04 SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE (altitude BETWEEN -7984.86 AND 54056.95 OR (user_id = 71 AND altitude BETWEEN -16987.76 AND 22297.35 AND datetime BETWEEN '2009-09-10 04:42:47' AND '2012-01-02 08:18:57')) ORDER BY T1.traj_id, datetime →Figure 4a
TRCAT	TVL1: Visualize map Area Tiedong District, Liaoning ST 2009-02-28 05:10:50 ET 2010-04-21 10:37:04 SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE altitude BETWEEN -7984.86 AND 54056.95 ORDER BY T1.traj_id, datetime TVL2: Visualize map Area Tiedong District, Liaoning ST 2009-09-10 04:42:47 ET 2012-01-02 08:18:57 SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE user_id = 71 AND altitude BETWEEN -16987.76 AND 22297.35 ORDER BY T1.traj_id, datetime →Figure 4b



(a) ✘



(b) ✔

Table 4: Case study on multi-trajectory query. TVLs generated under few-shot and TRCAT configurations, along with their corresponding map visualizations (Errors are marked in red).

C Case Study

Table 4 presents a case study on multi-trajectory queries using GPT-4o-mini under different configurations. As shown, under the few-shot prompting configuration, the model fails to correctly identify the implicit multi-trajectory query structure within the NLQ, generating only a single TVL instance.

Although this TVL instance includes all trajectory query conditions and visualizes two trajectories, it does not distinguish between them, making it difficult for users to distinguish individual trajectories (Figure 4a). This is because even if a user submits a single query, the corresponding results may contain multiple trajectory data points. For example,

1586
1587
1588
1589
1590
1591
1592

NLQ	Provide a visualization of the trajectory taken on the thoroughfare known as Mingshui Avenue in the province of Shandong, filtering for those with an altitude between -7109.45 and 102496.41.
Target TVL	TVL1: Visualize map Area Mingshui Street, Shandong ST null ET null SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE altitude BETWEEN -7109.45 AND 102496.41 ORDER BY T1.traj_id, datetime
Few-shot	TVL1: Visualize map Area Mingshui Avenue, Shandong ST null ET null SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE altitude IS NOT NULL ORDER BY T1.traj_id, datetime → Figure 5a
TRCAT	TVL1: Visualize map Area Mingshui Street, Shandong ST null ET null SELECT T1.traj_id, latitude, longitude, datetime FROM traj_data AS T1 JOIN add_traj_data AS T2 ON T1.traj_id = T2.traj_id WHERE altitude BETWEEN -7109.45 AND 102496.41 ORDER BY T1.traj_id, datetime → 5b



(a) X



(b) ✓

Table 5: Case study. TVLs generated by few-shot and TRCAT configuration, together with their corresponding map visualizations (Errors are marked with red colors).

when a query covers a specific time period, the system may return multiple independent trajectories recorded for the user within that period. In contrast, TRCAT not only reproduces the target TVL structure but also accurately infers the multi-trajectory query structure from the NLQ, generating multiple TVL instances as needed. Benefiting from this structure, the map visualization can distinguish between trajectories using attributes such as color and query ID, as illustrated in Figure 4b.

Table 5 presents a case study of trajectory query using GPT-4o-mini under different configurations. As shown, the few-shot configuration produced incorrect Area keyword and erroneous data query constraints, resulting in the absence of the corresponding trajectory in Figure 5a. In contrast, GPT-4o-mini under TRCAT configuration not only generates TVL with the correct structure but also produces accurate area name and data constraints, leading to precise visualization as illustrated in Figure 5b.