# FedTDP: A Privacy-Preserving and Unified Framework for Trajectory Data Preparation via Federated Learning

Zhihao Zeng [1]   Ziquan Fang [1]   Wei Shao [1]   Lu Chen [1]   Yunjun Gao [1]

## Abstract

Trajectory data, which capture the movement patterns of people and vehicles over time and space, are crucial for applications like traffic optimization and urban planning. However, issues such as noise and incompleteness often compromise data quality, leading to inaccurate trajectory analyses and limiting the potential of these applications. While Trajectory Data Preparation (TDP) can enhance data quality, existing methods suffer from two key limitations: (i) they do not address data privacy concerns, particularly in federated settings where trajectory data sharing is prohibited, and (ii) they typically design task-specific models that lack generalizability across diverse TDP scenarios. To overcome these challenges, we propose FedTDP, a privacy-preserving and unified framework that leverages the capabilities of Large Language Models (LLMs) for TDP in federated environments. Specifically, we: (i) design a trajectory privacy autoencoder to secure data transmission and protect privacy, (ii) introduce a trajectory knowledge enhancer to improve model learning of TDP-related knowledge, enabling the development of TDP-oriented LLMs, and (iii) propose federated parallel optimization to enhance training efficiency by reducing data transmission and enabling parallel model training. Experiments on 6 real datasets and 10 mainstream TDP tasks demonstrate that FedTDP consistently outperforms 13 state-of-the-art baselines.

## 1. Introduction

Trajectory data are typically represented as sequences of spatio-temporal points that describe the movement of objects, such as people (Chib et al., 2024) and vehicles (Lin et al., 2023). The widespread use of GPS and location-based

services has led to the generation of vast amounts of trajectory data (Fang et al., 2023; Chang et al., 2024; Zheng, 2015), enabling various analytical applications, including route planning (Yuan et al., 2020), crowd clustering (Liang et al., 2024), and traffic prediction (Wang et al., 2021).

However, trajectory data often suffer from significant quality issues due to sensor malfunctions, limited equipment precision, and transmission interruptions, leading to inconsistent (Liu et al., 2024b), noisy (Fan et al., 2023), and missing values (Yuan et al., 2024a). For instance, GPS location estimates in Uber (Uber, 2018) can be inaccurate by over 50 meters in densely populated, highly built-up urban areas. Such low-quality data undermine the reliability of trajectory analyses, limiting their practical applications. To address these issues, **Trajectory Data Preparation (TDP)**—which includes pre-processing and data mining operations such as data imputation (Musleh & Mokbel, 2023), map matching (Liu et al., 2024b), trajectory-user linking (Chen et al., 2024a), anomaly detection (Gao et al., 2023), and trajectory recovery (Liu et al., 2024c)—has become essential for improving data quality prior to analysis and application. However, existing TDP methods face two key limitations that affect their privacy protection and generalizability.

**Limitations.** (i) *Previous studies have not addressed data privacy constraints.* According to government reports [1] and related studies (Zeng et al., 2024; Meng et al., 2021; Shi et al., 2023), trajectory data is often collected or stored across multiple stations or organizations. Consequently, a moving trajectory may span several geographic regions, with each region's data collected by its respective signal station. For example, Fig. 1 illustrates the trajectory data from GeoLife (Zheng et al., 2010), a real-world dataset collected in Beijing, which shows six distinct colored regions, each storing its trajectory data separately. Due to legal constraints (GDPR, 2016; Beijing, 2023; China, 2021), the exchange and sharing of trajectory data across regions are prohibited. However, existing studies assume centralized data, which increases the risk of privacy breaches. (ii) *All previous studies are single-task approaches.* Specifically, these models are tailored to a single TDP task, such as data
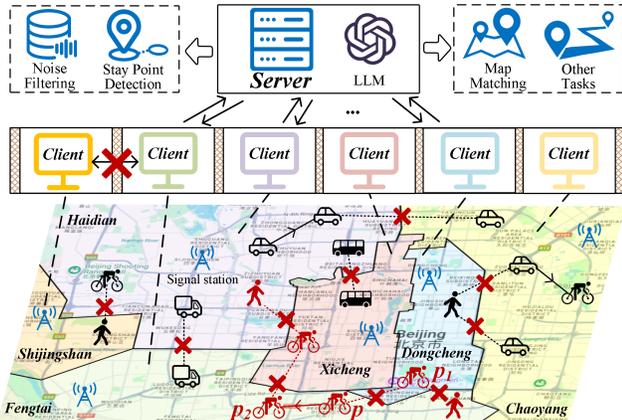
---

---

*Figure 1.* Trajectory Data Preparation in Federation (F-TDP)

imputation or anomaly detection. When addressing multiple TDP tasks, a new model must be trained for each task, resulting in high computational cost, extended training time, and limited generalizability.

Motivated by these limitations, we propose a **privacy-aware** and **generalized** framework for trajectory data preparation. (i) To protect data privacy, we introduce Federated Learning (FL)(Fanì et al., 2024; Zhang et al., 2024b), a privacy-preserving distributed learning paradigm. FL has been widely applied in domains such as urban computing (Wang et al., 2022) and transportation management (Yang et al., 2024) to address privacy concerns. For example, MobiSpaces (MS, 2022), a government-funded project by the European Union, collaborates with various transportation services to support Mobility-as-a-Service using FL. It provides a data governance platform for processing raw trajectory data from public transportation, traffic sensors, and maritime vessels for decentralized analysis. As shown in Fig. 1, FL enables multiple **Clients** (i.e., regions) to collaboratively train a model on a **Server** while keeping data decentralized, thereby preserving the privacy of each client (i.e., region). In this context, we focus on **Trajectory Data Preparation in Federated Environments (F-TDP)**. (ii) Inspired by the powerful capabilities of Large Language Models (LLMs) (Eltabakh et al., 2024; Liu et al., 2024a), particularly their success in multi-task learning, we aim to develop a **TDP-oriented LLM** to support various TDP tasks. Overall, our goal is to leverage LLMs to create a **privacy-preserving and unified framework FedTDP** for trajectory data preparation in federated environments. However, developing the FedTDP framework presents three key technical challenges that must be addressed.

*Challenge 1: How to safeguard trajectory data privacy in the FedTDP framework?* TDP tasks often necessitate considering the data context (Gao et al., 2023; Musleh & Mokbel, 2023; Liu et al., 2024b), which involves the exchange and sharing of data and demands collaborative processing across clients (i.e., cross-client TDP), raising privacy concerns. As shown in Fig. 1, if the data $p$ is missing, the

Fengtai region needs to utilize the context of the missing data (i.e., $p_1$ and $p_2$) for data imputation (Liu et al., 2024b; Xu et al., 2023). However, due to data privacy constraints, the Fengtai region cannot access $p_1$ from the Dongcheng region. Ensuring the privacy of trajectory data thus constitutes the first challenge the FedTDP framework must address.

*Challenge 2: How to develop a TDP-oriented LLM in the FedTDP framework?* Existing LLMs are primarily designed for text data (Fan et al., 2024; Li et al., 2024a). However, trajectory data exhibits unique spatio-temporal features (Musleh & Mokbel, 2023; Fang et al., 2023), such as temporal regularity and spatial dependency, which differ significantly from text data and are not inherently understood by LLMs. Furthermore, the pre-training of existing LLMs relies largely on publicly available unsupervised corpora (Eltabakh et al., 2024; Gu et al., 2023), which capture only general textual knowledge. In contrast, TDP tasks involve intricate spatio-temporal relationships and patterns (Liu et al., 2024b; Hu et al., 2024b), knowledge that is not adequately represented in these corpora. As a result, LLMs often perform poorly on TDP tasks. Effectively training a TDP-oriented LLM thus represents the second challenge that FedTDP must overcome.

*Challenge 3: How to improve the training efficiency of the FedTDP framework?* Due to the limited computational resources and storage capacities of clients, directly deploying and training LLMs locally on clients is infeasible (Sun et al., 2024; Zhang et al., 2024a; Hou et al., 2024). As a result, LLMs are typically hosted on servers, requiring clients (i.e., regions) to transmit their local data to the server for TDP processing. This introduces storage burdens on the server and wastes computational resources on the client. Additionally, LLMs often contain a large number of parameters, and even with techniques like Parameter-Efficient Fine-Tuning (PEFT) (Hu et al., 2022), training a TDP-oriented LLM remains highly time- and resource-intensive. Therefore, enhancing training efficiency represents the third challenge that the FedTDP framework must address.

**Contributions.** To address the challenges outlined above, we introduce the **Small Language Model (SLM)**, a compact version of the server's LLM, which is deployed on each client for local TDP. This approach leverages the computational resources of clients and reduces the server's workload, enabling distributed computation within the FedTDP framework. To address *Challenge 1*, we propose the **Trajectory Privacy AutoEncoder (TPA)**, which encodes trajectory data into spatio-temporal embeddings for transmission, rather than sending the raw data. This ensures data privacy while preserving the spatio-temporal correlations essential for TDP tasks. Additionally, we develop a decentralized secret-sharing method to safeguard against trajectory data recovery or inference through embedding (Song & Raghunathan, 2020; Chen et al., 2024b; Huang et al., 2024) and

gradient (Wang et al., 2024; Zheng et al., 2024; Zhu et al., 2023) inversion attacks. To address ***Challenge 2***, we design the **Trajectory Knowledge Enhancer (TKE)**, which helps both the SLM and LLM understand trajectory data and learn the specific knowledge required for TDP tasks. This enhances the model's ability to learn TDP-related patterns while reducing the number of parameters. To tackle ***Challenge 3***, we introduce **Federated Parallel Optimization (FPO)** to improve training efficiency. Specifically, FPO decomposes the federated training between server and clients through split learning, employs alternating optimization to minimize data transmission, and accelerates training via parallel execution. Finally, experiments on 6 real-world datasets demonstrate that the proposed FedTDP framework outperforms 13 baselines, achieving a performance improvement ranging from 4.84% to 45.22% across 10 mainstream TDP tasks. All code and data are available at `https://anonymous.4open.science/r/FedTDP`.

## 2. Preliminary

The frequently used notations is shown in **Appendix B**.

**Definition 1 (Spatio-Temporal Point)**. *A spatio-temporal point is represented as $p = \langle l, t \rangle$, where $l = (lon, lat)$ is a tuple of longitude and latitude location coordinates, and $t$ refers to the observed time associated with the point.*

**Definition 2 (Trajectory)**. *A trajectory comprises chronological spatio-temporal points, denoted as $T = \{p_1, p_2, \ldots\}$, which is typically represents the movement of a user. In addition, a trajectory can be segmented into multiple sub-trajectories, denoted as $T = \{ST^{(1)}, ST^{(2)}, \ldots\}$.*

**Definition 3 (Data Silo)**. *A data silo $S$ has its own collected trajectory dataset $D$. In the federated environment, a data silo $S$ is represented as a client $C$, typically a regional data storage platform or institution, responsible for the collection and management of trajectory data within that region. Specifically, a trajectory $T = \{p_1, p_2, \ldots\}$ is segmented into distinct sub-trajectories based on the geographic locations, denoted as $T = \{ST^{(C_1)}, ST^{(C_2)}, \ldots\}$, where sub-trajectory $ST^{(C_i)}$ is stored in client $C_i$.*

**Federated Architecture.** A federation usually consists of a central server and several clients. A trajectory may be segmented into multiple sub-trajectories, which are stored across different clients based on the spatial location of the spatio-temporal points. The server is responsible for aggregating the results from the clients to output the final result.

**Problem Formulation (F-TDP).** Given the server's LLM $\theta_{LLM}$ and the trajectory dataset $\mathcal{D} = \{D_1, D_2, \ldots\} \rightarrow \{T_1, T_2, \ldots\}$ of all clients $\mathcal{C} = \{C_1, C_2, \ldots\}$, where client $C_i$ holds dataset $D_i$, F-TDP is to employ $\theta_{LLM}$ on $\mathcal{D}$ for performing various trajectory data preparation tasks, where collected trajectories $D_i$ of client $C_i$ cannot be shared and exchanged to the server and other clients, as shown below:

$$F\text{-}TDP(\mathcal{D}) = \theta_{LLM}(T_i), T_i = \{ST_i^{(C_1)}, ST_i^{(C_2)}, \ldots\}, \quad (1)$$

where $\theta_{LLM}(T_i)$ is the result of $\theta_{LLM}$ on the trajectory $T_i$, with different forms of output depending on the TDP task, such as the cleaned trajectory, points, or classification result.

## 3. Trajectory Data Preparation Task

We demonstrate all major types of TDP tasks, with the rough processing shown in **Appendix B**, as detailed below.

**T-1: Anomaly Detection (*AD*).** It aims to detect trajectories that deviate significantly from typical movement behaviors. These anomalies could result from unusual user behavior, errors in data collection, or potential malicious activities.

**T-2: Trajectory Imputation (*TI*).** It aims to reconstruct a complete trajectory by estimating the missing points based on available spatio-temporal points. This often occurs when GPS signals are lost or data collection is interrupted.

**T-3: Noise Filtering (*NF*).** It aims to identifying and removing irrelevant spatio-temporal points that deviate from a trajectory. These noisy points can result from GPS inaccuracies, signal interference, or sensor malfunctions.

**T-4: Stay Point Detection (*SPD*).** It aims to identify locations where a moving object remains within an area for a certain period of time. A stay point typically represents a place of interest, such as a rest stop, home, or office.

**T-5: Map Matching (*MM*).** It aims to map the spatio-temporal point to the most probable segment in the road network. This is often the case when there is a deviation in the collected GPS position.

**T-6: Trajectory-User Link (*TUL*).** It aims to link an anonymous trajectory with its corresponding user. These trajectories are often collected ignoring the user's identity.

**T-7: Travel Mode Identification (*TMI*).** It aims to identify the travel mode based on the moving pattern of trajectory, which is walking, biking, taking the bus, or driving a car.

**T-8: Trajectory Simplification (*TSim*).** It aims to reduce the number of spatio-temporal points in a trajectory while preserving its essential shape and features.

**T-9: Trajectory Segmentation (*TSeg*).** It aims to divide a trajectory into meaningful segments based on specific criteria such as stay points or travel modes.

**T-10: Trajectory Recovery (*TR*).** It aims to reconstruct a complete trajectory from partially observed or incomplete spatio-temporal points. This often occurs when some parts of the trajectory are missing or unobserved.

## 4. Our Approach

Fig. 2 shows an overview of the FedTDP framework, which involves a server and multiple clients. FedTDP consists of
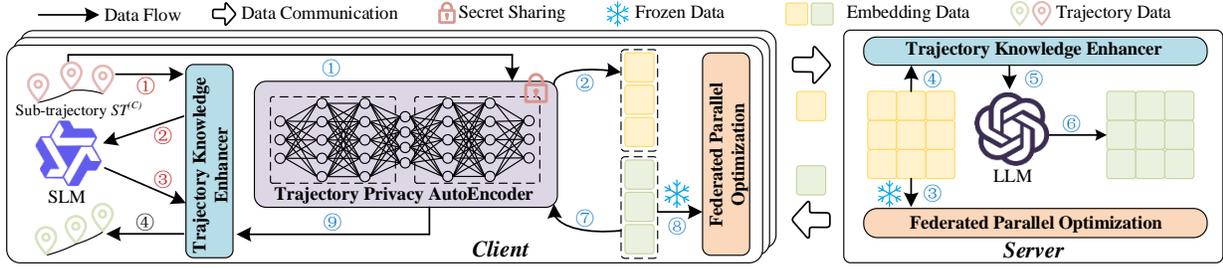
*Figure 2.* The Overview of Our Framework

four modules, i.e., Trajectory Privacy AutoEncoder (TPA), Trajectory Knowledge Enhancer (TKE), and Federated Parallel Optimization (FPO). To enable distributed computing of the FedTDP framework, we introduce the Small Language Model (SLM), a small-scale version of the server's Large Language Model (LLM), which is deployed on each client for local TDP, to leverage clients' computational resources and reduce the server's workload. Specifically, if the data context of TDP on the client's sub-trajectory $ST^{(C)}$ does not involve data from other clients for joint processing, the locally deployed SLM is used for local TDP, or $ST^{(\bar{C})}$ must be uploaded to the server and use the LLM for cross-client TDP. The overall process of FedTDP is as following. For the local TDP, TKE generates the TDP prompt as input for SLM (①–②). Next, TKE enhances the TDP knowledge with results outputted by the client's SLM to get final result (③–④). For the cross-client TDP, TPA first encodes trajectory data and transmits the encoded embeddings to the server (①–②). Then, FPO freezes the data transmitted from clients (③) and TKE generates the TDP prompt as input for LLM (④–⑤). Next, TPD decodes results outputted by the server's LLM (⑥–⑦). Finally, FPO freezes the data transmitted from the server (⑧) and TKE enhances the TDP knowledge with decoded results to get final results (⑨–④).

### 4.1. Trajectory Privacy AutoEncoder

FedTDP proposes a Trajectory Privacy AutoEncoder (TPA) to protect trajectory data privacy while maintaining spatio-temporal correlations. Specifically, TPA employs an encoder-decoder architecture that encodes trajectory data $T = \{p_1, p_2, \dots\}$ into embeddings $E = \{e_1, e_2, \dots\}$, where each spatio-temporal point $p_i$ is independently encoded as $e_i = \theta_{Enc}(p_i)$. Then, these clients' embeddings are transmitted to the server for aggregation $\mathcal{E} = \bigcup_{i=1}^{|\mathcal{C}|} E_i$, preserving both intra-client and inter-client spatio-temporal dependencies, which helps the LLM to capture spatio-temporal relationships in the trajectory data. Next, the server splits and distributes results $\tilde{\mathcal{E}} = \{\tilde{e_1}, \tilde{e_2}, \dots\}$ outputted by the LLM to clients, where the decoder reconstructs the estimated trajectory $\tilde{T} = \{\tilde{p_1}, \tilde{p_2}, \dots\}$ through $\tilde{p_i} = Dec(\tilde{e_i})$. Here, TPA is implemented as a three-layer Multi-Layer Perception (MLP) (Rosenblatt, 1958) with GELU (Hendrycks & Gimpel, 2016) activation, where the embedding and hidden dimensions are set to 32 and 256, respectively.

However, merely using embeddings for transmission cannot safeguard data privacy completely in FL, as attackers can recover the raw data by embedding (Song & Raghunathan, 2020; Chen et al., 2024b; Huang et al., 2024) and gradient (Wang et al., 2024; Zheng et al., 2024; Zhu et al., 2023) inversion attacks during TPA model aggregation. Specifically, traditional FL model aggregation, which exchange client gradients and aggregated parameters, are vulnerable to these attacks. While homomorphic encryption (Rivest et al., 1978) and differential privacy (Dwork et al., 2006) offer solutions, they introduce computational overhead or degrade model accuracy. In contrast, we propose a decentralized aggregation approach based on secret sharing (Shamir, 1979), achieving secure TPA aggregation without compromising efficiency or accuracy. Initially, each client pair $(C_i, C_j)$ generates a shared secret key $sk_{i,j} = sk_{j,i}$ stored locally, respectively. Then, the TPA parameters are partitioned into $|\mathcal{C}|$ parameter blocks $\{P^{(1)}, P^{(2)}, \dots\}$. For aggregation, the client $C_i$ masks its parameter block using secret keys $\{sk_{i,0}, sk_{i,1}, \dots\}$ determined with the other clients to mask parameter blocks, adding $sk_{i,j}$ if $i > j$ or subtracting it if $i < j$:

$$\tilde{P}_i^{(k)} = P_i^{(k)} + \sum_{j=1 \& j \neq i}^{|\mathcal{C}|} a_{i,j} * sk_{i,j}, \quad a_{i,j} = \begin{cases} 1, & i < j \\ -1, & i > j \end{cases}, \quad (2)$$

where client $C_i$ holds parameter block $P_i^{(k)}$, which is aggregated by client $C_k$, and $\tilde{P}_i^{(k)}$ is the mask parameter block.

**Theorem 4.1.** *Given the masked parameter block $\{\tilde{P}_1^{(k)}, \tilde{P}_2^{(k)}, \dots\}$ from all clients, the result of aggregating them is equal to the result of aggregating raw parameter blocks $\{P_1^{(k)}, P_2^{(k)}, \dots\}$ for all clients directly shown below:*

$$\sum_{i=1}^{|\mathcal{C}|} \tilde{P}_i^{(k)} = \sum_{i=1}^{|\mathcal{C}|} P_i^{(k)} \quad (3)$$

*Proof.* The proofs are provided in **Appendix C.1**. □

According to Theorem 4.1, the client $C_k$ can obtain the aggregation result $\overline{P}^{(k)}$ of the parameter block $P^{(k)}$ by aggregating the mask parameter blocks transmitted from clients:

$$\overline{P}^{(k)} = \frac{1}{|\mathcal{C}|} \sum_{i=1}^{|\mathcal{C}|} \tilde{P}_i^{(k)} = \frac{1}{|\mathcal{C}|} \sum_{i=1}^{|\mathcal{C}|} P_i^{(k)} \quad (4)$$

Finally, the aggregated parameter block is broadcasted to clients for TPA model updates.

## 4.2. Trajectory Knowledge Enhancer

To develop a TDP-oriented LLM, FedTDP designs Trajectory Knowledge Enhancer (TKE) consisting of trajectory prompt engineering, trajectory offsite-tuning, LoRA sparse-tuning, and bidirectional knowledge learning, to enhance the model learning abilities on TDP knowledge while reducing the number of training parameters.

**i) Trajectory Prompt Engineering** To help the SLM and LLM understand trajectory data and learn TDP knowledge, TKE designs a trajectory instruction paradigm to generate the TDP prompt, defined as (***Task***, ***Data***, ***Information***, ***Format***).

***Task***. It is the textual instruction consisting of the task name and the task description, as listed in the Section 3.

***Data***. It is the input trajectory data, either as trajectory data $T = \{p_1, p_2, \ldots\}$ to the SLM for local TDP or embeddings $E = \{e_1, e_2, \ldots\}$ to the LLM for cross-client TDP.

***Information***. It is the optional trajectory context (e.g., road network, weather) from public sources such as Open-StreetMap (OSM, 2024) and weather services (OWM, 2024), to enhance the model ability to perform TDP tasks.

***Format***. It is the task-specific output format, such as classification results for TDP tasks including AD, TUL, and TMI, trajectories for TDP tasks including TI, NF, TSim, TSeg, MM, and TR, and spatio-temporal points for the SPD task.

A few examples of TDP tasks using the trajectory prompt engineering are shown in **Appendix C.2**.

**ii) Trajectory Offsite-Tuning.** To enhance the SLM's learning capability, TKE employs the LLM to assist it in learning trajectory knowledge by trajectory offsite-tuning. Specifically, inspired by the offsite-tuning (Xiao et al., 2023), we divide the LLM into two components, denoted as $\theta_{LLM} = [\mathcal{A}, \mathcal{F}]$, where adapter $\mathcal{A}$ is the last few layers of the LLM to specialize general features for specific tasks, enabling task-specific feature mapping and decision-making, and foundation $\mathcal{F}$ is the remaining layers excluding $\mathcal{A}$ to extract and encode general data patterns, transforming raw inputs into meaningful representations.

Initially, it dispatches the server's adapter $\mathcal{A}$ to the client as the final few layers to be integrated into the client's SLM. Consequently, the SLM is composed of two components, denoted as $\theta_{SLM} = [\mathcal{A}, \mathcal{F}']$, where $\mathcal{F}'$ is the foundation of the SLM. Subsequently, the SLM employs LoRA to reduce the number of parameters in the adapter and then transmits the fine-tuning adapter to the server for aggregation and updates. Note that, rather than directly transferring the trained LLM's adapter to SLM, it utilizes and trains it to augment the SLM's learning capacity during training.

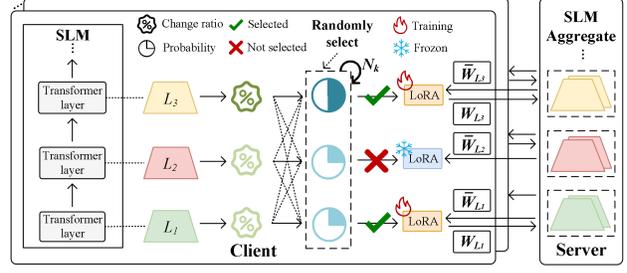**ii) LoRA Sparse-Tuning.** To reduce the number of training



*Figure 3.* LoRA Sparse-Tuning

parameters, TKE proposes LoRA sparse-tuning, as shown in Fig. 3. According to works on sparsity (Alistarh et al., 2018; Dai et al., 2022; Zhang et al., 2023a), more significantly varying parameters have a greater contribution to model convergence. Therefore, we only choose the layer in the SLM where the LoRA parameter change rate is the top $m$ for training. Specifically, the client calculates the ratio of the LoRA parameters change rate of each layer to the global LoRA parameters change rate of all layers ("ratio" for short):

$$R^{(r)}(L_i) = \frac{CR^{(r)}(L_i)}{\sum_{j=1}^{N} CR^{(r)}(L_j)}, \tag{5}$$

where $N$ is the number of layers and $CR^{(r)}(L_i)$ is the LoRA parameters change rate of layer $L_i$ in training round $r$:

$$CR^{(r)}(L_i) = \left| \frac{L_i^{(r)} - L_i^{(r-1)}}{L_i^{(r-1)}} \right| \tag{6}$$

Then, we randomly select $N_m$ layers to participate in the next round of training, where $N_m = \lfloor m * N \rfloor$.

**Theorem 4.2.** *Given the ratio $R^{(r)}(L_i)$ of layer $L_i$ in training round $r$ and the number of layers $N_m$ to be trained, the probability $Pr^{(r+1)}(L_i, N_m)$ of layer $L_i$ to train in the next round $r + 1$ is shown below:*

$$Pr^{(r+1)}(L_i, N_m) = R^{(r)}(L_i) + \sum_{j_1=1}^{N} \frac{R^{(r)}(L_i) * R^{(r)}(L_{j_1})}{1 - R^{(r)}(L_{j_1})} + \ldots$$

$$+ \sum_{j_1=1}^{N} \cdots \sum_{j_{N_m}=1}^{N} \frac{R^{(r)}(L_i) * R^{(r)}(L_{j_1}) * \ldots * R^{(r)}(L_{j_{N_m}})}{1 - R^{(r)}(L_{j_1}) - \ldots - R^{(r)}(L_{j_{N_m}})},$$

$$j_1 \neq \ldots \neq j_{N_m} \neq i, \tag{7}$$

*Proof.* The proofs are provided in **Appendix C.3**. □

According to Theorem 4.2, it chooses the training layers based on their probability at each training round. Finally, the client uploads the LoRA parameters of the trained layers to the server for aggregation, and the server assigns different weights to the parameters based on the number of clients involved in training on these layers, as shown below:

$$\overline{W}_{L_i}^{(r)} = \frac{(|\mathcal{C}| - |\mathcal{C}'|) * \frac{\sum_{j=1}^{|\mathcal{C}'|} n_j * W_{L_i,j}^{(r)}}{\sum_{j=1}^{|\mathcal{C}'|} n_j} + \overline{W}_{L_i}^{(r-1)}}{|\mathcal{C}| - |\mathcal{C}'| + 1}, \tag{8}$$

where $W_{L_i,j}^{(r)}$ is the LoRA parameters of layer $L_i$ in training round $r$ sent by client $j$, $\overline{W}_{L_i}^{(r)}$ is the aggregated LoRA parameters of layer $L_i$ in training round $r$, $|\mathcal{C}'|$ is the number of clients that have trained layer $L_i$ in training round $r$, and $n_j$ is the number of data used to train in client $C_j$.
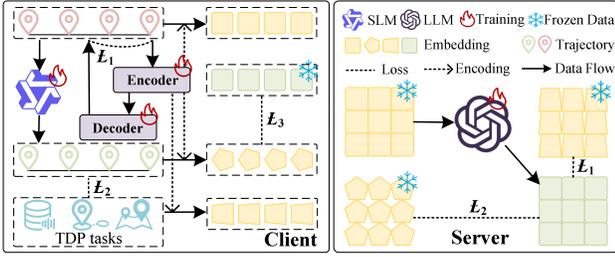
*Figure 4.* Federated Parallel Optimization

**iii) Bidirectional Knowledge Learning.** To improve the learning capabilities of the SLM and LLM, TKE develops bidirectional knowledge learning to enhance their TDP knowledge. Specifically, in order for the SLM to learn useful TDP knowledge in the complex output space of the LLM, it aligns the SLM's output with LLM's high frequency output using the inverse Kullback–Leibler (KL) divergence (Kullback & Leibler, 1951), shown as follows:

$$\min_{\theta_{SLM}} D_{KL}(P_{\theta_{SLM}}||P_{\theta_{LLM}}) = \sum_T P_{\theta_{SLM}}(T) \log(\frac{P_{\theta_{SLM}}(T)}{P_{\theta_{LLM}}(T)}) \quad (9)$$

where $P_{\theta_{SLM}}$ and $P_{\theta_{LLM}}$ are the output distribution of the SLM and LLM, respectively. Besides, since the SLM can access raw trajectory data, it aligns the LLM's output with SLM's overall output using the forward KL divergence, which enables the LLM to learn the whole knowledge of the SLM in the raw trajectory data, shown as follows:

$$\min_{\theta_{LLM}} D_{KL}(P_{\theta_{SLM}}||P_{\theta_{LLM}}) = \sum_T P_{\theta_{SLM}}(T) \log(\frac{P_{\theta_{SLM}}(T)}{P_{\theta_{LLM}}(T)}) \quad (10)$$

### 4.3. Federated Parallel Optimization

To improve training efficiency, FedTDP introduces Federated Parallel Optimization (FPO), which utilizes split learning, alternating optimization, and parallel training, to reduce data transmission and enhance the training parallelism. The overall process of the FPO module is shown in Fig. 4.

First, to enable the simultaneous training of the client and server, it employs split learning (Fu et al., 2024) to decompose the federated training process into client and server training. Specifically, the client is responsible for the training of the TPA model (i.e., the encoder and decoder) and SLM, while the server manages the training of the LLM.

Besides, to reduce data transmission, it utilizes alternating optimization (Meng et al., 2021) to freeze the data required by the client and server, respectively. During training, the server freezes the embeddings uploaded by the client for the LLM training, while the client freezes the results outputted by the server's LLM for the TPA model and SLM training.

Finally, to enhance the training parallelism, it uses parallel training to optimize several objectives in parallel. Specifically, the client focuses on three optimization objectives: (i) minimizing the reconstruction loss $\mathcal{L}_1$ of the TPA model, (ii) reducing the inverse KL loss $\mathcal{L}_2$ between SLM and LLM outputs, and (iii) minimizing the loss $\mathcal{L}_3$ between SLM outputs and labels. On the other hand, the server has two

*Table 1.* The Evaluated Trajectory Data Preparation Tasks

| Type | Task | Dataset |
|---|---|---|
| **Seen** **(seen in training)** | Anomaly Detection (**AD**) | GeoLife |
| | Trajectory Imputation (**TI**) | |
| | Map Matching (**MM**) | |
| | Trajectory-User Link (**TUL**) | |
| | Travel Mode Identification (**TMI**) | |
| | Trajectory Simplification (**TSim**) | |
| | Trajectory Recovery (**TR**) | |
| **Unseen** **(unseen in training)** | Anomaly Detection | Porto |
| | Trajectory Imputation | |
| | Noise Filtering (**NF**) | T-Drive |
| | Stay Point Detection (**SPD**) | |
| | Trajectory Simplification | |
| | Trajectory Recovery | |
| | Map Matching | Tencent |
| | Trajectory-User Link | Gowalla |
| | Travel Mode Identification | SHL |
| | Trajectory Segmentation (**TSeg**) | |

optimization objectives: (i) minimizing the forward KL loss $\mathcal{L}_1$ between LLM and SLM outputs, and (ii) reducing the loss $\mathcal{L}_2$ between LLM outputs and labels.

The overall training process of the FedTDP framework is shown in **Appendix C.4**.

## 5. Experiment

We use the Research Questions (RQs) to guide experiments.
*RQ1: How does FedTDP perform compared to state-of-the-art (SOTA) approaches in various TDP tasks?*
*RQ2: How does each module (i.e., TPA, TKE, and TPO) affect the performance of the FedTDP framework?*
*RQ3: How generalized is the FedTDP framework when using different numbers of seen TPD tasks for training?*
*RQ4: How does the FedTDP framework perform in different LLM and SLM model bases?*
*RQ5: How efficient is the FedTDP framework in terms of runtime and communication size?*
*RQ6: Is FedTDP sensitive to the hyperparameter settings?*

### 5.1. Experimental Settings

**Tasks and Datasets.** We evaluate the framework by the 10 mainstream tasks and 6 datasets in Table 1, which are widely studied in TDP communities (Liu et al., 2024b; Musleh & Mokbel, 2023; Hu et al., 2024a). The seen task is the training task and the unseen task is the testing task not shown in training. For the seen task, we conduct experiments using the GeoLife (Zheng et al., 2010) dataset, which was collected from April 2007 to August 2012. It contains various quality issues such as positional inaccuracies, data noise, and lower precision, which makes it suitable for various tasks. For the unseen task, the following datasets are used:

- **Porto** (Porto, 2015). It was collected in Porto from July 2013 to June 2014 with 442 taxis, which contains quality issues such as anomalies and missing data.
- **T-Drive** (Yuan et al., 2010). It was collected in Beijing in February 2008 with 10,357 trajectories, which contains quality issues such as noisy and incomplete points.
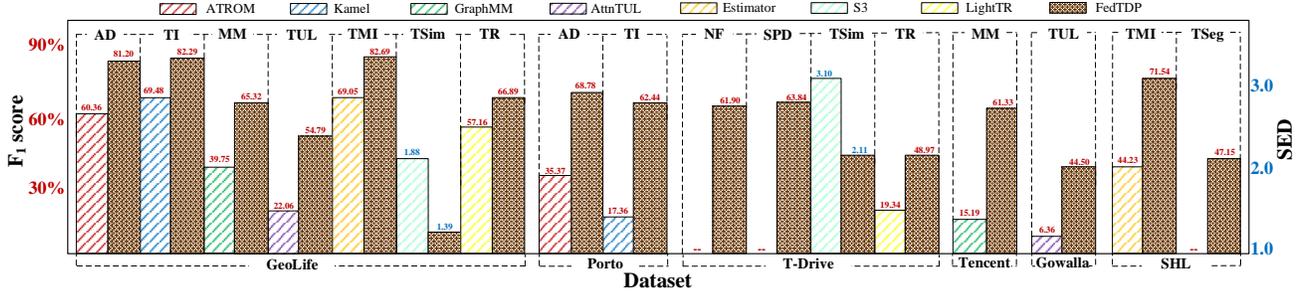- **Tencent** (Liu et al., 2024b). It was collected in Beijing

*Figure 5.* The Performance Comparison between FedTDP and None-LLM Trajectory Data Preparation Methods on Different Datasets
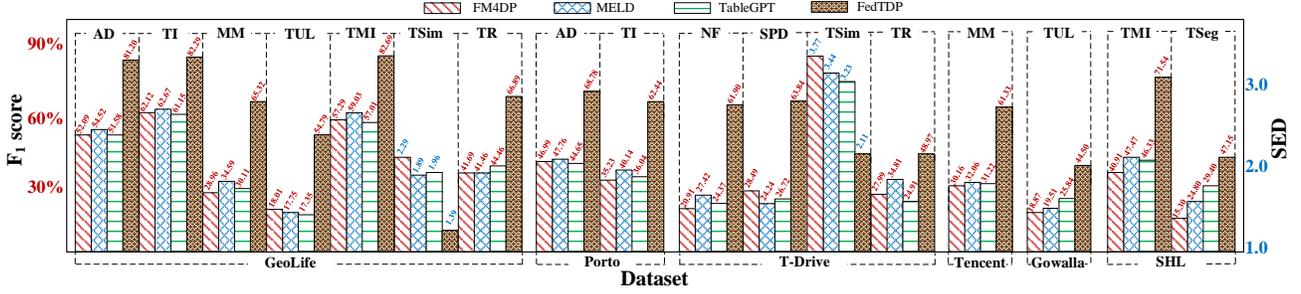


*Figure 6.* The Performance Comparison between FedTDP and LLM-based Table Data Preparation Methods on Different Datasets

city for 3 months, which contains quality issues such as inaccurate points due to the low sampling rate.

- **Gowalla** (Cho et al., 2011). It was collected in the social network from January to June in 2010 with 6,442,890 check-in data from 10,336 users.

- **SHL** (SHL, 2017). It was collected by the University of Sussex over 7 months in 2017 from 3 users, which contains various travel and movement modes.

We provide more details of these datasets in **Appendix D.1**.

**Baselines.** To answer ***RQ1***, (i) we compare FedTDP with **none-LLM methods** including ATROM (Gao et al., 2023) (for anomaly detection task), Kamel (Musleh & Mokbel, 2023) (for trajectory imputation task), GraphMM (Liu et al., 2024b) (for map matching task), AttnTUL (Chen et al., 2024a) (for trajectory-user link task), Estimator (Hu et al., 2024a) (for travel mode identification task), S3 (Fang et al., 2023) (for trajectory simplification task), and LightTR (Liu et al., 2024c) (for trajectory recovery task), which are the leading approaches in their respective research tasks; (ii) we compare with three SOTA **LLM-based table data preparation methods**, namely FM4DP (Narayan et al., 2022), MELD (Yan et al., 2024), and TableGPT (Li et al., 2024b); and (iii) we evaluate three SOTA **LLM-based spatio-temporal data analysis methods**, including PromptGAT (Da et al., 2024), UniST (Yuan et al., 2024b), and UrbanGPT (Li et al., 2024c). To answer ***RQ4***, we choose widely-used model bases for LLMs (Llama-8B (Touvron et al., 2023), GPT3-7B (Brown et al., 2020), and Qwen-7B (Bai et al., 2023)) and SLMs (GPT3-Small-125M (Brown et al., 2020), GPT2-Small-117M (Radford et al., 2019), and T5-Small-60M (Raffel et al., 2020)). More

details of these baselines are provided in **Appendix D.2**.

**Implementations.** Synchronized Euclidean Distance (SED) is used for the trajectory simplification task while $F_1$ score are used for other tasks. The lower the SED and the higher the $F_1$ score, the better the performance. Besides, we use the running time and communication size to evaluate the efficiency. All baselines run under their optimal settings. Besides, FedTDP can protect data privacy with the TPA module, while other baselines do not consider data privacy in F-TDP. To solve F-TDP, one alternative approach for baselines is to employ differential privacy (Dwork et al., 2006). Specifically, clients apply differential privacy to perturb local trajectory data before transmitting them to the server. Therefore, to safeguard data privacy and ensure fairness in experiments, we extend baselines combined with this optional approach to solve the F-TDP problem. Moreover, all experiments are conducted in the federation with 9 nodes, one as a server and the other 8 nodes as clients, each equipped with two Intel Xeon CPU E5-2650 12-core processors, two GeForce RTX 3090, and 100 MB/s internet.

### 5.2. Overall Performance (*RQ1*)

We present the overall performance comparison between the FedTDP framework with various SOTA baselines across different datasets and tasks. First, as shown in Fig. 5 (where the dash "− −" denotes tasks that are not supported), FedTDP demonstrates the best performance and robust generalization with an improvement of at least **18.38%** compared with none-LLM TDP methods, highlighting its effectiveness in addressing the F-TDP problem. Besides, the superior performance achieved by FedTDP on unseen TDP tasks also
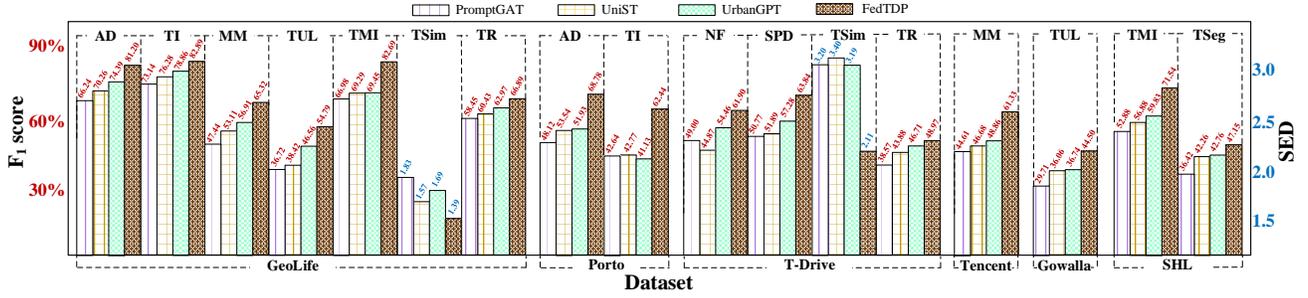
*Figure 7.* The Performance Comparison between FedTDP and LLM-based Spatio-temporal Data Analysis Methods on Different Datasets
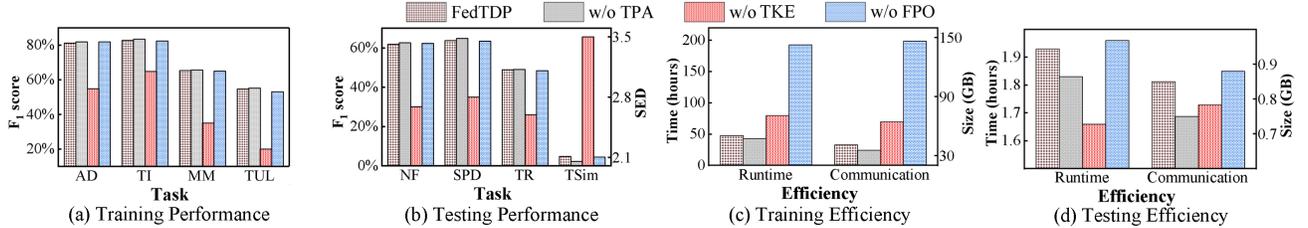


*Figure 8.* The Ablation Study

demonstrates its strong generalization. Second, as shown in Fig. 6, compared with SOTA LLM-based table data preparation, FedTDP achieve the best performance across different datasets and tasks with an improvement of at least **32.26%**. Third, as shown in Fig. 7, compared with SOTA LLM-based spatio-temporal data analysis methods, FedTDP also shows the best performance with an improvement of **4.84% to 45.22%**. We attribute these improvements to the developed TDP-oriented LLM and SLM in the FedTDP framework.

### 5.3. Ablation Study (*RQ2*)

We evaluate the effectiveness of each module in the FedTDP framework by systematically removing one at a time, with the following configurations: FedTDP without Trajectory Privacy AutoEncoder (w/o TPA), without Trajectory Knowledge Enhancer (w/o TKE), and without Federated Parallel Optimization (w/o FPO). The results are shown in Fig. 8.

First, the performance of FedTDP is slightly degraded compared to w/o TPA as TPA can not fully capture the spatio-temporal information of the trajectory data, leading to a marginal performance decline when using TPA. Besides, FedTDP has a slight increase in runtime and communication costs because TPA transmits higher-dimensional embedding data instead of three-dimensional spatio-temporal points, introducing greater communication size and runtime when TPA is employed. However, to safeguard data privacy, the use of TPA in the FedTDP framework is essential.

Second, the performance of FedTDP improves dramatically compared to w/o TKE, with at least **27.52%** improvement. This is because TKE enhances the model learning abilities on TDP knowledge to develop the TDP-orient LLM and SLM. Additionally, FedTDP has lower runtime and communication costs during training, since the TKE module can reduce the number of parameters that need to be trained and

transmitted, which speeds up the model training.

Finally, the performance of FedTDP does not change significantly compared to w/o FPT, but its training runtime and communication overhead are significantly reduced with almost 4 times less. This reduction is because FPO can reduce data transmission and improve training efficiency.

### 5.4. Model Generalization Study (*RQ3*)

We evaluate FedTDP's generalization in different numbers of training tasks, as details presented in **Appendix D.3**.

### 5.5. Model Base Study (*RQ4*)

We evaluate the impact of various model bases on FedTDP, as shown in **Appendix D.4**.

### 5.6. Efficiency Study (*RQ5*)

We evaluate the communication costs and running times. The results are shown in **Appendix D.5**.

### 5.7. Parameter Sensitivity Study (*RQ6*)

We evaluate the effect of FedTDP's hyperparameter (i.e., the training layers ratio $m$ of LoRA sparse-tuning in the TKE module). The results are shown in **Appendix D.6**.

## 6. Conclusion

This paper introduces FedTDP, a privacy-preserving, unified framework for trajectory data preparation. It includes a trajectory privacy autoencoder to protect data while maintaining spatio-temporal correlations, a trajectory knowledge enhancer to develop TDP-oriented LLMs, and parallel optimization to boost efficiency. Experiments on 6 datasets and 10 TDP tasks validate its effectiveness. In the future, we will extend FedTDP for more trajectory analysis tasks.

# References

Alistarh, D., Hoefler, T., Johansson, M., Konstantinov, N., Khirirat, S., and Renggli, C. The convergence of sparsified gradient methods. In *NeurIPS*, pp. 5977–5987, 2018.

Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., Hui, B., Ji, L., Li, M., Lin, J., Lin, R., Liu, D., Liu, G., Lu, C., Lu, K., Ma, J., Men, R., Ren, X., Ren, X., Tan, C., Tan, S., Tu, J., Wang, P., Wang, S., Wang, W., Wu, S., Xu, B., Xu, J., Yang, A., Yang, H., Yang, J., Yang, S., Yao, Y., Yu, B., Yuan, H., Yuan, Z., Zhang, J., Zhang, X., Zhang, Y., Zhang, Z., Zhou, C., Zhou, J., Zhou, X., and Zhu, T. Qwen technical report. *CoRR*, abs/2309.16609, 2023.

Beijing. Beijing municipal public data management measures. https://www.beijing.gov.cn/zhengce/zhengcefagui/202312/t20231211_3496032.html, 2023.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *NeurIPS*, 2020.

Chang, Y., Tanin, E., Cong, G., Jensen, C. S., and Qi, J. Trajectory similarity measurement: An efficiency perspective. *Proc. VLDB Endow.*, 17(9):2293–2306, 2024.

Chen, W., Huang, C., Yu, Y., Jiang, Y., and Dong, J. Trajectory-user linking via hierarchical spatio-temporal attention networks. *ACM Trans. Knowl. Discov. Data*, 18(4):85:1–85:22, 2024a.

Chen, Y., Lent, H. C., and Bjerva, J. Text embedding inversion security for multilingual language models. In *ACL*, pp. 7808–7827, 2024b.

Chib, P. S., Nath, A., Kabra, P., Gupta, I., and Singh, P. MS-TIP: imputation aware pedestrian trajectory prediction. In *ICML*, 2024.

China. Personal information protection law of the people's republic of china. https://www.gjbmj.gov.cn/n1/2024/0910/c459362-40317200.html, 2021.

Cho, E., Myers, S. A., and Leskovec, J. Friendship and mobility: user movement in location-based social networks. In *KDD*, pp. 1082–1090, 2011.

Da, L., Gao, M., Mei, H., and Wei, H. Prompt to transfer: Sim-to-real transfer for traffic signal control with prompt learning. In *AAAI*, pp. 82–90, 2024.

Dai, R., Shen, L., He, F., Tian, X., and Tao, D. Dispfl: Towards communication-efficient personalized federated learning via decentralized sparse training. In *ICML*, volume 162, pp. 4587–4604, 2022.

Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *TCC*, pp. 265–284, 2006.

Eltabakh, M. Y., Naeem, Z. A., Ahmad, M. S., Ouzzani, M., and Tang, N. Retclean: Retrieval-based tabular data cleaning using llms and data lakes. *Proc. VLDB Endow.*, 17(12):4421–4424, 2024.

Fan, J., Gu, Z., Zhang, S., Zhang, Y., Chen, Z., Cao, L., Li, G., Madden, S., Du, X., and Tang, N. Combining small language models and large language models for zero-shot NL2SQL. *Proc. VLDB Endow.*, 17(11):2750–2763, 2024.

Fan, Y., Yu, X., Wieser, R., Meakin, D., Shaton, A., Jaubert, J., Flottemesch, R., Howell, M., Braid, J., Bruckman, L. S., French, R. H., and Wu, Y. Spatio-temporal denoising graph autoencoders with data augmentation for photovoltaic data imputation. *Proc. ACM Manag. Data*, 1(1):50:1–50:19, 2023.

Fang, Z., He, C., Chen, L., Hu, D., Sun, Q., Li, L., and Gao, Y. A lightweight framework for fast trajectory simplification. In *ICDE*, pp. 2386–2399, 2023.

Fanì, E., Camoriano, R., Caputo, B., and Ciccone, M. Accelerating heterogeneous federated learning with closed-form classifiers. In *ICML*, 2024.

Fu, F., Wang, X., Jiang, J., Xue, H., and Cui, B. Projpert: Projection-based perturbation for label protection in split learning based vertical federated learning. *IEEE Trans. Knowl. Data Eng.*, 36(7):3417–3428, 2024.

Gao, Q., Wang, X., Liu, C., Trajcevski, G., Huang, L., and Zhou, F. Open anomalous trajectory recognition via probabilistic metric learning. In *IJCAI*, pp. 2095–2103, 2023.

GDPR. General data protection regulation. https://gdpr-info.eu, 2016.

Gu, Z., Fan, J., Tang, N., Cao, L., Jia, B., Madden, S., and Du, X. Few-shot text-to-sql translation using structure and content prompt learning. *Proc. ACM Manag. Data*, 1(2):147:1–147:28, 2023.

Hendrycks, D. and Gimpel, K. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

Hou, C., Shrivastava, A., Zhan, H., Conway, R., Le, T., Sagar, A., Fanti, G., and Lazar, D. Pre-text: Training language models on private federated data in the age of llms. In *ICML*, 2024.

Hu, D., Fang, Z., Fang, H., Li, T., Shen, C., Chen, L., and Gao, Y. Estimator: An effective and scalable framework for transportation mode classification over trajectories. *IEEE Trans. Intell. Transp. Syst.*, pp. 1–12, 2024a.

Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022.

Hu, Y., Du, Y., Zhang, Z., Fang, Z., Chen, L., Zheng, K., and Gao, Y. Real-time trajectory synthesis with local differential privacy. In *ICDE*, pp. 1685–1698, 2024b.

Huang, Y., Tsai, Y., Hsiao, H., Lin, H., and Lin, S. Transferable embedding inversion attack: Uncovering privacy risks in text embeddings without model queries. In *ACL*, pp. 4193–4205, 2024.

Kullback, S. and Leibler, R. A. On information and sufficiency. *The annals of mathematical statistics*, 22(1): 79–86, 1951.

Li, H., Li, S., Hao, F., Zhang, C. J., Song, Y., and Chen, L. Booster: Leveraging large language models for enhancing entity resolution. In *WWW*, pp. 1043–1046, 2024a.

Li, P., He, Y., Yashar, D., Cui, W., Ge, S., Zhang, H., Fainman, D. R., Zhang, D., and Chaudhuri, S. Table-gpt: Table fine-tuned GPT for diverse table tasks. *Proc. ACM Manag. Data*, 2(3):176, 2024b.

Li, Z., Xia, L., Xu, Y., and Huang, C. GPT-ST: generative pre-training of spatio-temporal graph neural networks. In *NeruIPS*, 2023.

Li, Z., Xia, L., Tang, J., Xu, Y., Shi, L., Xia, L., Yin, D., and Huang, C. Urbangpt: Spatio-temporal large language models. In *KDD*, pp. 5351–5362, 2024c.

Li, Z., Xia, L., Xu, Y., and Huang, C. Flashst: A simple and universal prompt-tuning framework for traffic prediction. In *ICML*, 2024d.

Liang, A., Yao, B., Wang, B., Liu, Y., Chen, Z., Xie, J., and Li, F. Sub-trajectory clustering with deep reinforcement learning. *VLDB J.*, 33(3):685–702, 2024.

Lin, Y., Wan, H., Hu, J., Guo, S., Yang, B., Lin, Y., and Jensen, C. S. Origin-destination travel time oracle for map-based services. *Proc. ACM Manag. Data*, 1(3): 217:1–217:27, 2023.

Liu, X., Tien, C., Ding, P., Jiang, S., and Stevens, R. L. Entropy-reinforced planning with large language models for drug discovery. In *ICML*, 2024a.

Liu, Y., Ge, Q., Luo, W., Huang, Q., Zou, L., Wang, H., Li, X., and Liu, C. Graphmm: Graph-based vehicular map matching by leveraging trajectory and road correlations. *IEEE Trans. Knowl. Data Eng.*, 36(1):184–198, 2024b.

Liu, Z., Miao, H., Zhao, Y., Liu, C., Zheng, K., and Li, H. Lighttr: A lightweight framework for federated trajectory recovery. In *ICDE*, pp. 4422–4434, 2024c.

Meng, C., Rambhatla, S., and Liu, Y. Cross-node federated graph neural network for spatio-temporal data modeling. In *KDD*, pp. 1202–1211, 2021.

MS. New data spaces for green mobility. `https://mobispaces.eu/`, 2022.

Musleh, M. and Mokbel, M. F. KAMEL: A scalable bert-based system for trajectory imputation. *Proc. VLDB Endow.*, 17(3):525–538, 2023.

Narayan, A., Chami, I., Orr, L. J., and Ré, C. Can foundation models wrangle your data? *Proc. VLDB Endow.*, 16(4): 738–746, 2022.

OSM. Openstreemap. `https://www.openstreetmap.org`, 2024.

OWM. Openweathermap. `https://www.openweathermap.org`, 2024.

Porto. Ecml/pkdd 15: Taxi trajectory prediction. `https://www.kaggle.com/c/pkdd-15-predict-taxi-service-trajectory-i/data`, 2015.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020.

Rivest, R. L., Shamir, A., and Adleman, L. A method for obtaining digital signatures and public-key cryptosystems. *Communications of the ACM*, 21(2):120–126, 1978.

Rosenblatt, F. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.

Shamir, A. How to share a secret. *Commun. ACM*, 22(11): 612–613, 1979.

Shi, Y., Tong, Y., Zeng, Y., Zhou, Z., Ding, B., and Chen, L. Efficient approximate range aggregation over large-scale spatial data federation. *IEEE Trans. Knowl. Data Eng.*, 35(1):418–430, 2023.

SHL. Sussex-huawei locomotion dataset. http://www.shl-dataset.org/, 2017.

Song, C. and Raghunathan, A. Information leakage in embedding models. In *CCS*, pp. 377–390, 2020.

Sun, J., Xu, Z., Yin, H., Yang, D., Xu, D., Liu, Y., Du, Z., Chen, Y., and Roth, H. R. Fedbpt: Efficient federated black-box prompt tuning for large language models. In *ICML*, 2024.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023.

Uber. Rethinking gps: Engineering next-gen location at uber. https://www.uber.com/blog/rethinking-gps/, 2018.

Wang, S., Sun, Y., Musco, C., and Bao, Z. Public transport planning: When transit network connectivity meets commuting demand. In *SIGMOD*, pp. 1906–1919, 2021.

Wang, Y., Tong, Y., Zhou, Z., Ren, Z., Xu, Y., Wu, G., and Lv, W. Fed-ltd: Towards cross-platform ride hailing via federated learning to dispatch. In *KDD*, pp. 4079–4089, 2022.

Wang, Y., Liang, J., and He, R. Towards eliminating hard label constraints in gradient inversion attacks. In *ICLR*, 2024.

Xiao, G., Lin, J., and Han, S. Offsite-tuning: Transfer learning without full model. *CoRR*, abs/2302.04870, 2023.

Xu, Y., Bazarjani, A., Chi, H., Choi, C., and Fu, Y. Uncovering the missing pattern: Unified framework towards trajectory imputation and prediction. In *CVPR*, pp. 9632–9643, 2023.

Yan, M., Wang, Y., Pang, K., Xie, M., and Li, J. Efficient mixture of experts based on large language models for low-resource data preprocessing. In *KDD*, pp. 3690–3701, 2024.

Yang, L., Chen, W., He, X., Wei, S., Xu, Y., Zhou, Z., and Tong, Y. Fedgtp: Exploiting inter-client spatial dependency in federated graph-based traffic prediction. In *KDD*, pp. 6105–6116, 2024.

Yuan, H., Li, G., Bao, Z., and Feng, L. Effective travel time estimation: When historical trajectories over road networks matter. In *SIGMOD, June 14-19, 2020*, pp. 2135–2149, 2020.

Yuan, H., Cong, G., and Li, G. Nuhuo: An effective estimation model for traffic speed histogram imputation on A road network. *Proc. VLDB Endow.*, 17(7):1605–1617, 2024a.

Yuan, J., Zheng, Y., Zhang, C., Xie, W., Xie, X., Sun, G., and Huang, Y. T-drive: driving directions based on taxi trajectories. In *SIGSPATIAL*, pp. 99–108, 2010.

Yuan, Y., Ding, J., Feng, J., Jin, D., and Li, Y. Unist: A prompt-empowered universal model for urban spatio-temporal prediction. In *KDD*, pp. 4095–4106, 2024b.

Zeng, Z., Fang, Z., Chen, L., Gao, Y., Zheng, K., and Chen, G. Fedctq: A federated-based framework for accurate and efficient contact tracing query. In *ICDE*, pp. 4628–4642, 2024.

Zhang, J., Vahidian, S., Kuo, M., Li, C., Zhang, R., Yu, T., Wang, G., and Chen, Y. Towards building the federatedgpt: Federated instruction tuning. In *ICASSP*, pp. 6915–6919, 2024a.

Zhang, W., Zhou, Z., Wang, Y., and Tong, Y. DM-PFL: hitchhiking generic federated learning for efficient shift-robust personalization. In *KDD*, pp. 3396–3408, 2023a.

Zhang, X., Wang, Q., Xu, C., Peng, Y., and Xu, J. Fedknn: Secure federated k-nearest neighbor search. *Proc. ACM Manag. Data*, 2(1):V2mod011:1–V2mod011:26, 2024b.

Zhang, Z., Zhao, X., Liu, Q., Zhang, C., Ma, Q., Wang, W., Zhao, H., Wang, Y., and Liu, Z. Promptst: Prompt-enhanced spatio-temporal multi-attribute prediction. In *CIKM*, pp. 3195–3205, 2023b.

Zheng, L., Cao, Y., Jiang, R., Taura, K., Shen, Y., Li, S., and Yoshikawa, M. Enhancing privacy of spatiotemporal federated learning against gradient inversion attacks. In *DASFAA*, volume 14850, pp. 457–473, 2024.

Zheng, Y. Trajectory data mining: An overview. *ACM Trans. Intell. Syst. Technol.*, 6(3):29:1–29:41, 2015.

Zheng, Y., Xie, X., and Ma, W. Geolife: A collaborative social networking service among user, location and trajectory. *IEEE Data Eng. Bull.*, 33(2):32–39, 2010.

Zhu, J., Yao, R., and Blaschko, M. B. Surrogate model extension (SME): A fast and accurate weight update attack on federated learning. In *ICML*, volume 202, pp. 43228–43257, 2023.

# Appendix

# Appendix

In the subsequent sections, we present supplementary materials to provide more details of this paper, offering deeper insights and additional technical details for readers seeking further clarification. The appendix is organized as follows.

In **Section A**, we present a systematic review of related work to help readers understand the key development in areas relevant to this paper, including (i) the latest trajectory data preparation methods, (ii) applications of large language models in other data preparation, and (iii) recent advancements in large language models for spatio-temporal data analysis.

In **Section B**, we summary the preliminary of notations and trajectory data preparation tasks for better understanding our work, including (i) the frequently used notations and (ii) detailed description of trajectory data preparation tasks.

In **Section C**, we provide the additional methodology details to support the analysis shown in the main body of this paper, including (i) complete theoretical proofs of Theorems 4.1 and 4.2, (ii) practical examples of trajectory data preparation tasks using the trajectory prompt engineering of the proposed trajectory knowledge enhancer, and (iii) the training process of FedTDP.

In **Section D**, we describe the extensive experimental details to provide more information about experimental settings and further demonstrate the superiority performance of the proposed FedTDP framework, including (i) datasets description, (ii) compared baselines introduction, and (iii) the details experimental results of model generalization, model base, efficiency, and parameter sensitivity studies.

## A. Related Work

**Trajectory Data Preparation.** Numerous Trajectory Data Preparation (TDP) methods have been proposed to improve the quality of trajectory data for trajectory data preparation tasks. For the anomalous detection task, ATROM (Gao et al., 2023) addresses the critical challenge of anomaly recognition in open-world scenarios through the development of a probabilistic metric learning model, which significantly improves the accuracy of anomaly detection in complex environments. For the trajectory imputation task, Kamel (Musleh & Mokbel, 2023) proposes a scalable architecture that incorporates additional real trajectory points to predict the missing trajectory data, improving the accuracy of trajectory imputation. For the map matching task, GraphMM (Liu et al., 2024b) leverages the graphical structure using a graph neural network to effectively model the topology of road network and trajectory features, improving the accuracy of map matching. For the trajectory-user link task, AttnTUL (Chen et al., 2024a) introduces a hierarchical spatio-temporal attention neural network, which co-encodes local trajectory transition patterns and global spatial dependencies to establish links between trajectories and users more accurately. For the travel mode identification task, Estimator (Hu et al., 2024a) proposes an effective and scalable framework that partitions the traffic space into disjoint spatial regions based on traffic conditions, improving the accuracy of travel mode identification. For the trajectory simplification task, S3 (Fang et al., 2023) presents a lightweight framework consisting of two chained sequence-to-sequence modules, which is integrated within a graph neural architecture, improving the accuracy and efficiency of trajectory simplification. For the trajectory recovery task, LightTR (Liu et al., 2024c) presents an efficient framework using a local trajectory embedding module, robust feature extraction capabilities while significantly reducing computational overhead. However, none of these studies have considered data privacy constraints. They typically assume that the trajectory data collection is centralized, which introduces a significant risk of privacy leakage, especially in federated learning environments. In addition, all of them are single-task methods. When handling multiple TDP tasks, different models need to be trained for each specific task. It not only demands substantial time and computational resources but also results in poor model generalization ability. ***In contrast, we aim to propose a privacy-preserving and unified framework to support various trajectory data preparation tasks while safeguarding trajectory data privacy.***

**Large Language Models for Other Data Preparation.** A few works on table data preparation using Large Language Models (LLMs) have been proposed recently. For instance, MELD (Yan et al., 2024) introduces a general solver for low-resource table data preparation, leveraging a mixture-of-experts architecture to support merging and augmentation of domain-specific experts trained on limited annotated examples. Similarly, TableGPT (Li et al., 2024b) presents a table-tuning paradigm, where LLMs are fine-tuned using various table tasks synthesized from real tables to enhance the model's ability to understand and process table-related tasks. Additionally, (Narayan et al., 2022) explores the performance of LLMs for table data preparation tasks, which evaluates their performance on five data cleaning and integration tasks through prompt-based methods. However, these works are specifically tailored for table data preparation and are not directly applicable to trajectory data preparation tasks. They lack the necessary understanding of trajectory data and do not account for the spatio-temporal characteristics and complexity of trajectory data preparation tasks, making them unsuitable for such applications. ***In contrast,***

*we aim to develop a TDP-oriented LLM to effectively support various trajectory data preparation tasks.*

**Large Language Models for Spatio-Temporal Data Analysis.** There are a few spatio-temporal LLMs proposed (Li et al., 2024d; Zhang et al., 2023b; Li et al., 2023), which have achieved superior performance in spatio-temporal downstream applications. Specifically, UrbanGPT (Li et al., 2024c) integrates a spatio-temporal dependency encoder with a command adjustment paradigm to enhance the LLMs' understanding of complex temporal and spatial interdependencies. Besides, UniST (Yuan et al., 2024b) develops a general-purpose model for urban spatio-temporal prediction through diverse data utilization, effective pre-training, and knowledge-guided prompts. In addition, PromptGAT (Da et al., 2024) employs prompt-based grounded action transformations to analyze system dynamics by leveraging reasoning capabilities of large language models to understand environmental impacts on traffic patterns. However, none of them have considered trajectory data quality. If the quality of trajectory data is extremely poor, the performance of spatio-temporal large language models in downstream tasks will not be satisfactory either. *In contrast, we aim to explore the powerful capabilities of large language models for trajectory data preparation to enhance the quality of trajectory data.*

## B. Notation and Trajectory Data Preparation Task

**Notation and Description.** We first present the frequently used notations and descriptions in this paper, as listed in Table 2.

*Table 2.* Notation and Description

| Notation | Description |
|---|---|
| $p$ | A spatio-temporal point consisting of location and time $\langle l, t \rangle$ |
| $T$ | A trajectory consisting of multiple spatio-temporal points $\{p_1, p_2, \ldots\}$ |
| $ST$ | A sub-trajectory of the trajectory $T$ |
| $S$ | A data silo that represents a region |
| $C$ | A client that represents a region |
| $D$ | The trajectory database in a client $C$ |
| $\mathcal{C}$ | A set of clients $\{C_1, C_2, \ldots\}$ |
| $\mathcal{D}$ | A set of trajectory datasets $\{D_1, D_2, \ldots\}$ |
| $\theta_{LLM}, \theta_{SLM}$ | The large language model and small language model |

**Trajectory Data Preparation Task.** Besides, we summary the supported trajectory data preparation tasks of the proposed FedTDP in this paper, as shown in Table 3, with the rough processing of each task shown in Fig. 9.

*Table 3.* Trajectory Data Preparation Task

| Task Cetegory | Task | Description |
|---|---|---|
| Data Cleaning | Anomaly Detection (AD) | Detect anomalous trajectory |
| | Trajectory Imputation (TI) | Predict missing points |
| | Noise Filtering (NF) | Filter point noise |
| | Stay Point Detection (SPD) | Identify stationary points |
| Data Matching | Map Matching (MM) | Align a trajectory to road network |
| | Trajectory-User Linking (TUL) | Associate trajectories with users |
| Data Annotation | Travel Mode Identification (TMI) | Identify transportation mode |
| Data Reduction | Trajectory Simplification (TSim) | Remove number of points |
| | Trajectory Segmentation (TSeg) | Divide a trajectory to segments |
| Data Augmentation | Trajectory Recovery (TR) | Recovery complete trajectory |

## C. Additional Methodology Details

### C.1. Proof of Theorem 4.1

We provide the complete theoretical proof of Theorem 4.2 proposed in this paper that provides the correctness of the trajectory privacy autoencoder model aggregation, as detailed below.

*Proof.* According to Eq. 2, we can get the result of aggregating masked parameter blocks $\{\tilde{P}_1^{(k)}, \tilde{P}_2^{(k)}, \ldots, \tilde{P}_{|\mathcal{C}|}^{(k)}\}$ for all

(a) Anomaly Detection  (b) Trajectory Imputation  (c) Noise Filtering  (d) Stay Point Detection  (e) Map Matching

(f) Trajectory-User Link  (g) Travel Mode Identification  (h) Trajectory Simplification  (i) Trajectory Segmentation  (j) Trajectory Recovery
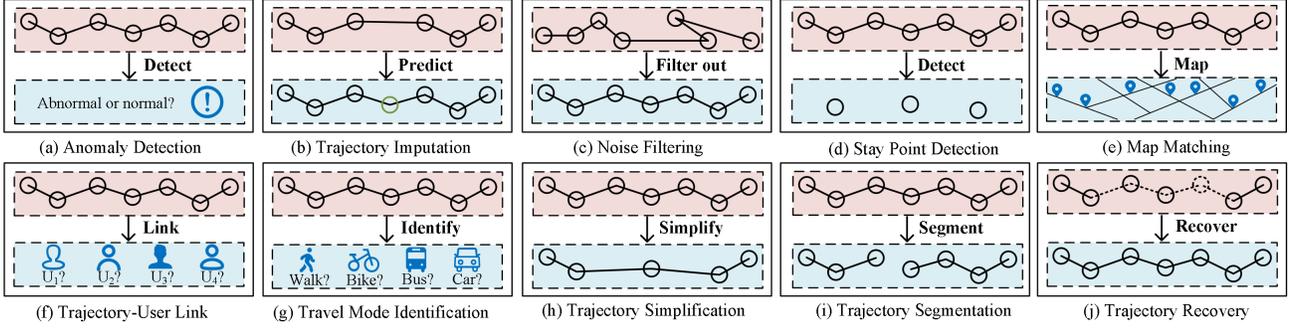
*Figure 9.* The Supported Trajectory Data Preparation Tasks

clients, as formally shown below:

$$\sum_{i=1}^{|\mathcal{C}|} \tilde{P}_i^{(k)} = \sum_{i=1}^{|\mathcal{C}|} P_i^{(k)} + \sum_{i=1}^{|\mathcal{C}|} \sum_{j=1 \& j \neq i}^{|\mathcal{C}|} a_{i,j} * sk_{i,j} \quad (11)$$

Here, $sk_{i,j} = sk_{j,i}$ and $a_{i,j} = -a_{j,i}$, thus we can get the result formally shown below:

$$\sum_{i=1}^{|\mathcal{C}|} \sum_{j=1 \& j \neq i}^{|\mathcal{C}|} a_{i,j} * sk_{i,j} = \sum_{i=j+1}^{|\mathcal{C}|} \sum_{j=1}^{|\mathcal{C}|} (a_{i,j} * sk_{i,j} + a_{j,i} * sk_{j,i}) = 0 \quad (12)$$

The aggregation $\tilde{P}^{(k)}$ of the masked parameter block is formally shown below:

$$\sum_{i=1}^{|\mathcal{C}|} \tilde{P}_i^{(k)} = \sum_{i=1}^{|\mathcal{C}|} P_i^{(k)} + \sum_{i=1}^{|\mathcal{C}|} \sum_{j=1 \& j \neq i}^{|\mathcal{C}|} a_{i,j} * sk_{i,j} = \sum_{i=1}^{|\mathcal{C}|} P_i^{(k)} \quad (13)$$

$\square$

## C.2. Examples of Trajectory Prompt Engineering



| **T-3: Noise Filtering** | **T-7: Travel Mode Identification** |
| --- | --- |
| ***Task:*** it is a noise filtering task. It targets identifying and removing irrelevant spatio-temporal points that deviate from a trajectory. These noisy points can result from GPS inaccuracies, signal interference, or sensor malfunctions. | ***Task:*** it is a travel mode identification task. It aims to identify the travel mode of a trajectory based on the moving pattern of data. The travel mode is usually the walk, bike, bus, or car. |
| ***Data:*** the trajectory data consisting of spatio-temporal points is {*(lon₁, lat₁, t₁), (lon₂, lat₂, t₂), ...*}. | ***Data:*** the trajectory data consisting of spatio-temporal points is {*(lon₁, lat₁, t₁), (lon₂, lat₂, t₂), ...*}. |
| ***Information:*** the weather is ***sunny*** with an average temperature of ***15***. | ***Information:*** It is ***sunny*** with an average temperature of ***15***. The road network of the map is {*(id₁, start-lon₁, start-lat₁, stop-lon₁, stop-lat₁), (id₂, start-lon₂, start-lat₂, stop-lon₂, stop-lat₂), ...*}. |
| ***Format:*** the output should be the trajectory data. | ***Format:*** the output should be the travel mode. |
| **Model response:** the result of noise filtering is {(lon₁, lat₁, time₁), (lon₂, lat₂, time₂), ...}. | **Model response:** the travel mode is car. |

*Figure 10.* The Example of Trajectory Prompt Engineering in Clients

For better understanding the trajectory prompt engineering of the proposed trajectory knowledge enhancer, we show practical examples of noise filtering and travel mode identification tasks using the trajectory prompt engineering in the small language model of the client.

As shown in Fig. 10, *Task* shows the task name with its description listed in Section 3. Besides, *Data* uses the raw trajectory data in clients, consisting of spatio-temporal points i.e., $T = \{p_1, p_2, \ldots\}$. Additionally, *Information* includes optional road network and weather data. Specifically, *sunny* and 15 represent the weather tokens during the trajectory, *id* denotes the segment ID, *start-lon* and *start-lat* indicate the segment's longitude and latitude at the starting point, and *stop-lon* and *stop-lat* denote the segment's longitude and latitude at the stopping point. Specifically, we utilize the weather data for these tasks and the road network for the travel mode identification task. Moreover, *Format* refers to the expected output data of the task. We anticipate the outputs of noise filtering and travel mode identification to be the trajectory data and travel mode, respectively. Finally, the model's response should return the expected results in *Format*. Note that, in the server *Data* uses embeddings uploaded by clients, consisting of encoded spatio-temporal points i.e., $E = \{e_1, e_2, \ldots\}$.

## C.3. Proof of Theorem 4.2

We provide the complete theoretical proof of Theorem 4.2 proposed in this paper that determines the probability of parameters that need to be trained in the LoRA sparse-tuning of the proposed trajectory knowledge enhancer, as detailed below.

*Proof.* To derive the probability that each layer is selected when $N_m$ layers are chosen for the next round of training, we first calculate the probability that each layer is selected in the first time, which is equal to the ratio of the layer, as formally shown below:

$$Pr_1^{(r+1)}(L_i, N_m) = R^{(r)}(L_i) \tag{14}$$

Then, we calculate the probability that the layer is not selected in the first time and is selected in the second time, as formally shown below:

$$Pr_2^{(r+1)}(L_i, N_m) = \sum_{j_1=1}^{N} \frac{R^{(r)}(L_i) * R^{(r)}(L_{j_1})}{1 - R^{(r)}(L_{j_1})}, j_1 \neq i \tag{15}$$

Then, we calculate the probability that the layer is not selected in the first two times and is selected in the third time, as formally shown below:

$$Pr_3^{(r+1)}(L_i, N_m) = \sum_{j_1=1}^{N} \sum_{j_2=1}^{N} \frac{R^{(r)}(L_i) * R^{(r)}(L_{j_1}) * R^{(r)}(L_{j_2})}{1 - R^{(r)}(L_{j_1}) - R^{(r)}(L_{j_2})}, \tag{16}$$

where $j_1 \neq j_2 \neq i$. Based on the above inductive steps, we can calculate the probability that the layer is not selected in the first $N_m - 1$ times and is selected in the $N_m$ time by inductive reasoning, as formally shown below:

$$Pr_{N_m}^{(r+1)}(L_i, N_m) = \sum_{j_1=1}^{N} \cdots \sum_{j_{N_m}=1}^{N} \frac{R^{(r)}(L_i) * R^{(r)}(L_{j_1}) * \ldots * R^{(r)}(L_{j_{N_m}})}{1 - R^{(r)}(L_{j_1}) - \ldots - R^{(r)}(L_{j_{N_m}})}, \tag{17}$$

where $j_1 \neq \ldots \neq j_{N_m} \neq i$. Thus, the probability that layer $L_i$ is selected to train in the next round $r + 1$ can be calculated as follows:

$$
\begin{aligned}
Pr^{(r+1)}(L_i, N_m) = \sum_{j=1}^{N_m} Pr_j^{(r+1)}(L_i, N_m) = R^{(r)}(L_i) + \sum_{j_1=1}^{N} \frac{R^{(r)}(L_i) * R^{(r)}(L_{j_1})}{1 - R^{(r)}(L_{j_1})} + \ldots + \\
\sum_{j_1=1}^{N} \cdots \sum_{j_{N_m}=1}^{N} \frac{R^{(r)}(L_i) * R^{(r)}(L_{j_1}) * \ldots * R^{(r)}(L_{j_{N_m}})}{1 - R^{(r)}(L_{j_1}) - \ldots - R^{(r)}(L_{j_{N_m}})},
\end{aligned}
\tag{18}
$$

where $j_1 \neq \ldots \neq j_{N_m} \neq i$. $\qquad\square$

## C.4. Multi-Task Training

Due to the diverse range of trajectory data preparation tasks that need to be addressed, we propose a multi-task training strategy to enhance the model's learning and generalization capabilities. Specifically, we prepare a trajectory dataset applicable to most trajectory data preparation tasks and construct labels for each task. During the training phase, we execute

multiple trajectory data preparation tasks on the same trajectory data input, calculate the loss for each task, and jointly optimize the model formulaically shown below:

$$\mathcal{L} = \mathcal{L}_{T\text{-}1} + \mathcal{L}_{T\text{-}2} + \ldots + \mathcal{L}_{T\text{-}10} \,, \tag{19}$$

where $\mathcal{L}_{T\text{-}i}$ is the loss of trajectory data preparation task T-$i$ listed in the Section 3. Note that the proposed FedTDP framework can be easily extended to support other trajectory data preparation tasks benefiting from its modular architecture, decoupled data processing pipeline, and variable model base.

**Training Algorithm.** For convenient method reproduction, we provide a detailed training process of the entire FedTDP framework, which can be divided into the server and client, as shown in Algorithms 1 and 2.

---

**Algorithm 1** The training on the server

1: **Input:** the number of training rounds *TR*
2: **for** round $r = 0, \ldots, TR - 1$ **do**
3:     $f \leftarrow IsFrozen(r)$ // Get the frozen status of this round.
4:     **if** $f$ *is False* **then**
5:         $E \leftarrow Get(\mathcal{C})$ // Get the embeddings data from clients.
6:         $E \leftarrow Connect(E)$ // Connect into a complete embeddings.
7:     **else**
8:         $E \leftarrow GetFrozenData(r - 1)$ // Get the frozen data.
9:     **end if**
10:    $prompt \leftarrow TKE(E)$ // Construct the prompt of the embeddings.
11:    $o \leftarrow \theta_{LLM}(prompt)$ // Input the prompt data to the LLM.
12:    $o \leftarrow Split(o)$ // Split the output of the LLM.
13:    **if** $f$ *is False* **then**
14:       **for** client number $i = 0, \ldots, |\mathcal{C}| - 1$ **do**
15:         $Send(o_i, C_i)$ // Send split results to respective clients.
16:       **end for**
17:    **end if**
18: **end for**

---

**Algorithm 2** The training on the client

1: **Input:** the number of training rounds *TR* and server $s$
2: **for** round $r = 0, \ldots, TR - 1$ **do**
3:     $D \leftarrow GetData()$ // Get the local trajectory data.
4:     $f \leftarrow IsFrozen(r)$ // Get the frozen status of this round.
5:     **if** $f$ *is False* **then**
6:         $E \leftarrow Enc(D)$ // Encode the trajectory into embeddings.
7:         $Send(E, s)$ // Send the embeddings data to the server.
8:     **end if**
9:     $prompt \leftarrow TKE(D)$ // Construct the prompt of the data.
10:    $o^{'} \leftarrow \theta_{SLM}(prompt)$ // Input the prompt data to the SLM.
11:    **if** $f$ *is False* **then**
12:       $o \leftarrow Get(s)$ // Get the result from the server.
13:       $o \leftarrow Dec(o)$ // Decode the server's result.
14:    **else**
15:       $o \leftarrow GetFrozenData(r - 1)$ // Get the frozen data.
16:    **end if**
17:    $result \leftarrow TKE(o^{'}, o)$ // Compute the distillation result.
18: **end for**

---

In the server (i.e., Algorithm 1), the input is the number of training rounds (line 1). For each training round $r$, it begins to get the frozen state $f$ (lines 2–3). If $f$ is not frozen, the server gets the trajectory embeddings $E$ from clients $\mathcal{C}$ and connects

them, or it gets local $E$ frozen in the last training round $r - 1$ (lines 4–9). Then, the server uses TKE to construct the TDP prompt for the LLM and get the output $o$ (lines 10–11). Finally, the server splits it into several parts and sends them to respective clients if $f$ is not frozen (lines 12–18).

In the client (i.e., Algorithm 2), the input are the number of training rounds and the server (line 1). For each training round $r$, it begins to get the trajectory data $D$ and frozen state $f$ (lines 2–4). If $f$ is not frozen, clients encode $D$ into embeddings and send them to the server (lines 5–8). Then, clients use TKE to construct the TDP prompt for the SLM and get the output $o'$ (lines 9–10). If $f$ is not frozen, clients get the LLM's output $o$ from the server and decode it, or it gets local $o$ frozen in the last round $r - 1$ (lines 11–16). Finally, clients use TKE to compute the final result between $o'$ and $o$ (lines 17–18).

**Complexity Analyses.** We also give complexity analyses for Algorithms 1 and 2. Specifically, given the number of trajectory embeddings data $|E|$ from all clients, the complexity of Algorithm 1 is $O(|E| * TR * MC)$, where $MC$ is the model complexity of the LLM. Given the number of trajectories $|D|$ in the client, the complexity of Algorithm 2 is $O(|D| * TR * MC')$, where $MC'$ is the model complexity of the SLM.

# D. Experimental Details

## D.1. Dataset

We evaluate the proposed FedTDP framework using 6 datasets, including GeoLife (Zheng et al., 2010), Porto (Porto, 2015), T-Drive (Yuan et al., 2010), Tencent (Liu et al., 2024b), Gowalla (Cho et al., 2011), and SHL (SHL, 2017), with their statistics shown in Table 4, as detailed below.

*Table 4.* The Statistics of Dataset

| Dataset | # trajectories | # points | Quality Issue | Task |
|---------|----------------|----------|---------------|------|
| **GeoLife** | 182 | 24,876,978 | Positional inaccuracies, data noise, and lower precision | AD, TI, MM, TUL, TMI, TSim, and TR |
| **Porto** | 442 | 83,409,386 | Anomalies and missing data | AD and TI |
| **T-Drive** | 10,336 | 17,662,984 | Noisy and incomplete points | NF, TR, SPD, and TSim |
| **Tencent** | 40,966 | 1,610,216 | Inaccurate points | MM |
| **Gowall** | 107,092 | 6,442,890 | Sparse and non-continuous data | TUL |
| **SHL** | 3 | 109,390 | Duplicate records | TSeg and TMI |

- **GeoLife (Zheng et al., 2010).** It collected 182 trajectories with 24,876,978 spatio-temporal points, used for the training tasks, including Anomaly Detection (AD), Trajectory Imputation (TI), Map Matching (MM), Trajectory-User Link (TUL), Travel Mode Identification (TMI), Trajectory Simplification (TSim), and Trajectory Recovery (TR) tasks. It contains various quality issues such as positional inaccuracies, data noise, and lower precision due to irregular sampling intervals and sensor limitations, which make it suitable for various trajectory data preparation tasks.

- **Porto (Porto, 2015).** It collected 442 trajectories with 83,409,386 spatio-temporal points, which contains quality issues such as anomalies and missing data, used for AD and TI testing tasks.

- **T-Drive (Yuan et al., 2010).** It collected 10,336 trajectories with 17,662,984 spatio-temporal points, which contains quality issues such as noisy and incomplete points, used for NF, TR, SPD, and TSim testing tasks.

- **Tencent (Liu et al., 2024b).** It collected 40,966 trajectories with 1,610,216 spatio-temporal points, which contains quality issues such as inaccurate points due to the low sampling rate, used for the MM testing task.

- **Gowalla (Cho et al., 2011).** It collected 107,092 trajectories with 6,442,890 spatio-temporal points, which contains quality issues such as sparse and non-continuous data, used for the TUL testing task.

- **SHL (SHL, 2017).** It collected 3 trajectories with 109,390 spatio-temporal points, which contains quality issues such as duplicate records, used for TSeg and TMI testing tasks.

*Table 5.* The Compared Baselines

| Category | Method | Task | Year |
|---|---|---|---|
| **S-TDP** | ATROM | Anomaly Detection | 2023 |
| | Kamel | Trajectory Imputation | 2023 |
| | GraphMM | Map Matching | 2024 |
| | AttnTUL | Trajectory-User Linking | 2024 |
| | Estimator | Travel Mode Identification | 2024 |
| | S3 | Trajectory Simplification | 2023 |
| | LightTR | Trajectory Recovery | 2024 |
| **Large Language Models for Table Data Preparation** | FM4DP | All tasks for evaluation | 2022 |
| | MELD | | 2024 |
| | TableGPT | | 2024 |
| **Large Language Models for Spatio-Temporal Data Analysis** | PromptGAT | | 2024 |
| | UniST | | 2024 |
| | UrbanGPT | | 2024 |

## D.2. Baseline

We compare the proposed FedTDP framework with state-of-the-art baselines, as shown in Table 5.

First, we compare FedTDP with various Trajectory Data Preparation (TDP) methods in a single TDP task (referred to S-TDP), including ATROM (Gao et al., 2023) for the AD task, Kamel (Musleh & Mokbel, 2023) for the TI task, GraphMM (Liu et al., 2024b) for the MM task, AttnTUL (Chen et al., 2024a) for the TUL task, Estimator (Hu et al., 2024a) for the TMI task, S3 (Fang et al., 2023) for the TSeg task, and LightTR (Liu et al., 2024c) for the TR task, as detailed below.

- **ATROM (Gao et al., 2023).** It solves the anomaly detection task in open-world scenarios and introduces a new probabilistic metric learning model.

- **Kamel (Musleh & Mokbel, 2023).** It proposes a scalable system that inserts additional real trajectory points to improve the accuracy of the trajectory imputation task.

- **GraphMM (Liu et al., 2024b).** It develops the graphical nature of the map matching task to exploit the road network and trajectory graphical topology.

- **AttnTUL (Chen et al., 2024a).** It proposes a hierarchical trajectory attention neural network for co-encoding local trajectory transition patterns and global spatial dependencies to solve the trajectory-user link task.

- **Estimator (Hu et al., 2024a).** It partitions the entire traffic space into disjoint spatial regions based on the traffic conditions for the travel mode identification task.

- **S3 (Fang et al., 2023).** It presents a lightweight trajectory segmentation task framework to augment the trajectory representation paradigm with geo-semantics.

- **LightTR (Liu et al., 2024c).** It develops a local trajectory embedding module that provides higher computational efficiency for the trajectory recovery task.

Besides, we compare FedTDP with three methods using Large Language Models (LLMs) for table data preparation, including FM4DP (Narayan et al., 2022), MELD (Yan et al., 2024), and TableGPT (Li et al., 2024b), as detailed below.

- **FM4DP (Narayan et al., 2022).** It helps LLMs understand table DP tasks, which uses 5 data cleaning and integration table DP tasks as prompt tasks and evaluates the performance of LLMs on these tasks.

- **MELD (Yan et al., 2024).** It employs the mixture-of-experts architecture to support the merging and augmentation of experts trained on the domain-specific experts trained on limited annotated examples.

- **TableGPT (Li et al., 2024b).** It proposes a table-tuning paradigm using various table tasks synthesized from real tables as the training data to help LLMs understand the table data and perform table tasks.

Moreover, we compare FedTDP with three LLM-based for spatio-temporal data analysis, including PromptGAT (Da et al., 2024), UniST (Yuan et al., 2024b), and UrbanGPT (Li et al., 2024c), as detailed below.

- **PromptGAT (Da et al., 2024).** It uses the LLM to analyze system dynamics, leveraging the context and spatio-temporal data to understand how weather, traffic, and road conditions affect traffic dynamics.

- **UniST (Yuan et al., 2024b).** It proposes a general-purpose model, which is designed for urban spatio-temporal prediction in various urban scenarios to capture the complex spatio-temporal relationship.

- **UrbanGPT (Li et al., 2024c).** It integrates a spatio-temporal dependency encoder with a command adjustment paradigm, which enables LLMs to understand complex spatio-temporal interdependencies.

## D.3. Model Generalization Study

To evaluate the proposed FedTDP framework generalization in different numbers of training tasks, we systematically remove the training task from back to front based on their order in Table 1. As illustrated in Fig.11, the results indicate that the performance of FedTDP across various tasks declines as the number of training tasks decreases. This decline is primarily attributed to the reduced acquisition of TDP knowledge, which adversely affects generalization. Notably, when the number of tasks is reduced to one (i.e., training FedTDP solely on the anomaly detection task using GeoLife), the performance also falls below that of S-TDP in Fig. 5
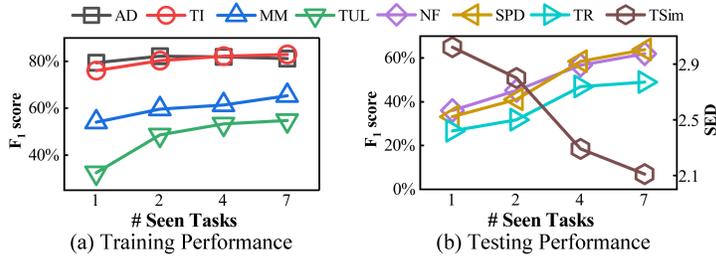


*Figure 11.* Scalability Study

## D.4. Model Base Study

To evaluate the impact of various model bases on the proposed FedTDP framework, we choose widely used model bases for the LLM (Llama-8B (Touvron et al., 2023), GPT3-7B (Brown et al., 2020), and Qwen-7B (Bai et al., 2023)) and SLM (GPT3-Small-125M (Brown et al., 2020), GPT2-Small-117M (Radford et al., 2019), and T5-Small-60M (Raffel et al., 2020)). Besides, to evaluate the impact of different LLMs on FedTDP, we use GPT3-Small as the client's SLM while we use Llama as the server's LLM to evaluate the impact of different SLMs on FedTDP. The results are shown in Fig. 12. As observed, Llama achieves optimal performance in most TPD tasks for the server's LLM, followed by GPT3 and then Qwen. In contrast, GPT3-Small demonstrates the best performance for the client's SLM, succeeded by T5-Small and then GPT2-Small. Consequently, we adopt Llama and GPT3-Small as the default server's LLM and client's SLM in other experiments, respectively.
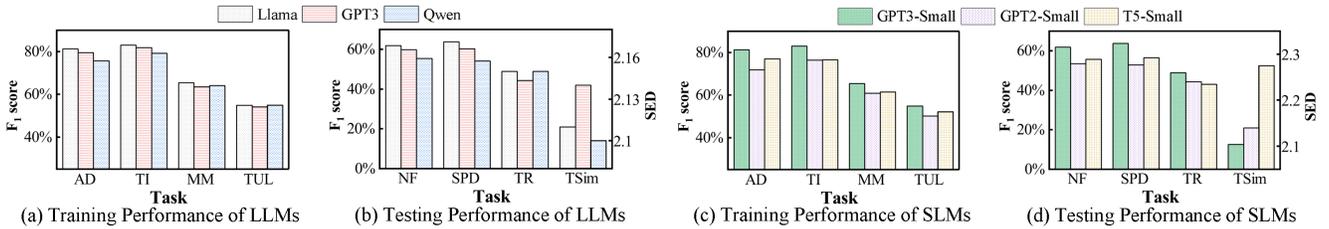


*Figure 12.* Large Language Model and Small Language Model Base Study

## D.5. Efficiency Study

Fig. 13 shows the communication costs (in GB) and running times (in hours) of various methods across all TDP tasks. During the training process, the proposed framework incurs the largest communication size because it must transfer embeddings

and model parameters, whereas LLM-based methods (i.e., LLMs for table table data preparation and spatio-temporal LLMs) transmit all perturbed data, generated through differential privacy, to the server in the first round of training and do not require data transmission in the following rounds of training. Furthermore, the runtime of FedTDP is reduced by a factor of **11.3 to 14.2** compared to other LLM-based methods, demonstrating its efficiency in the F-TDP context. In the testing phase, the communication size of FedTDP is nearly identical to that of S-TDP and **1.4 to 1.8** times less than that of other LLM-based methods, which require transferring all data to the server, whereas FedTDP only transmits the data necessary for cross-client TDP. Additionally, the runtime of FedTDP is **2.6 to 4.8** times lower than that of other LLM-based methods, further underscoring its efficiency.

## D.6. Parameter Sensitivity Study

We evaluate the effects of hyperparameters of the proposed FedTDP framework (i.e., the training layers ratio $m$ of LoRA sparse-tuning in the TKE module), as shown in Fig. 14, where we change the $m$ from 25% to 100%. We can observe that as $m$ increases, the performance of FedTDP improves slightly. However, this improvement comes at the cost of increased training time and communication size, as the number of parameters that need to be trained and transmitted also rises when $m$ is increased. Therefore, the suggestion value of $m$ is 25% or less, as long as the model performance with the value of $m$ is acceptable.
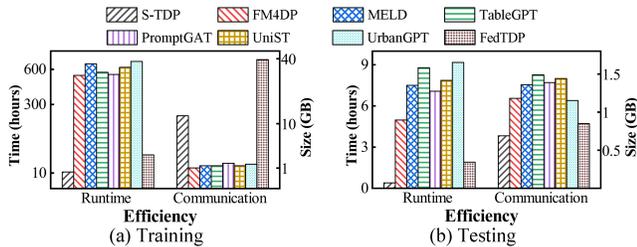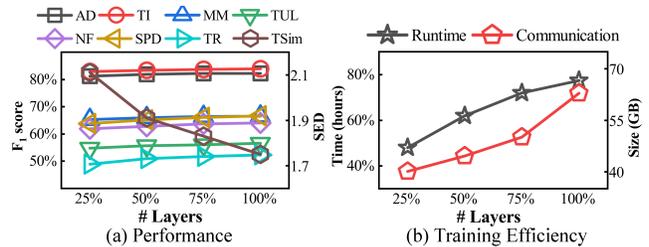


*Figure 13.* Efficiency Study of Different Methods



*Figure 14.* Parameter Sensitive Study