

Pre-training General Trajectory Embeddings with Maximum Multi-view Entropy Coding

Yan Lin, Huaiyu Wan, Shengnan Guo, Jilin Hu, Christian S. Jensen, Youfang Lin

Abstract—Spatio-temporal trajectories provide valuable information about movement and travel behavior, enabling various downstream tasks that in turn power real-world applications. Learning trajectory embeddings can improve task performance but may incur high computational costs and face limited training data availability. Pre-training learns generic embeddings by means of specially constructed pretext tasks that enable learning from unlabeled data. Existing pre-training methods face (i) difficulties in learning general embeddings due to biases towards certain downstream tasks incurred by the pretext tasks, (ii) limitations in capturing both travel semantics and spatio-temporal correlations, and (iii) the complexity of long, irregularly sampled trajectories.

To tackle these challenges, we propose Maximum Multi-view Trajectory Entropy Coding (MMTEC) for learning general and comprehensive trajectory embeddings. We introduce a pretext task that reduces biases in pre-trained trajectory embeddings, yielding embeddings that are useful for a wide variety of downstream tasks. We also propose an attention-based discrete encoder and a NeuralCDE-based continuous encoder that extract and represent travel behavior and continuous spatio-temporal correlations from trajectories in embeddings, respectively. Extensive experiments on two real-world datasets and three downstream tasks offer insight into the design properties of our proposal and indicate that it is capable of outperforming existing trajectory embedding methods.

Index Terms—Spatio-temporal data mining, trajectory embedding, pre-training, self-supervised learning, maximum multi-view entropy.

1 INTRODUCTION

A spatio-temporal trajectory is a long sequence of irregularly sampled (location, time) pairs that records the movement and travel behavior of a vehicle or a person. With the widespread availability of GPS-equipped devices and the growing popularity of location-based services, trajectory data from sources such as in-vehicle navigation devices and smartphones have become increasingly available. This wealth of data enables a wide range of spatio-temporal trajectory mining tasks, including predicting future trajectories [1], estimating time of arrival [2], [3], [4], anomaly detection [5], and trajectory clustering [6], [7]. It is attractive to let such downstream tasks operate on embeddings of trajectories— d -dimensional vectors that capture essential aspects of trajectories—rather than on the “native” trajectories themselves. However, when doing so, the accuracy and comprehensiveness of the embeddings are crucial to achieving good performance.

Existing trajectory mining methods that use embeddings commonly adopt an end-to-end embedding approach. This involves training a trajectory encoder alongside other components within the model, using task-specific learning objectives [8], [9], [10]. While intuitive and straightforward to implement, this approach also possesses limitations related

to generalizability, efficiency, and real-world feasibility [11], [12]. Thus, the embeddings learned using this approach are tailored to specific tasks, rendering them less applicable across different tasks. The resulting need to retrain embeddings for each new task reduces computational efficiency. Moreover, the end-to-end embedding approach relies on the availability of extensive task labels to achieve satisfactory outcomes. Obtaining such labels can incur substantial time and financial costs. Instead, we advocate focusing on pre-training general trajectory embedding methods. Such methods entail acquiring foundational embeddings through pretext tasks. This approach offers key advantages. First, it enables utilization of existing unlabeled trajectory data, facilitating the capture of general spatio-temporal information inherent in trajectories. Consequently, the accuracy of predictions and generalization performance in downstream tasks can be enhanced [13], [14], [15]. Second, pre-trained embeddings can be fine-tuned with limited labeled data or can be applied directly to unsupervised tasks. This expedites downstream task training and improves computational efficiency [16], [17], [18].

The design of a self-supervised pretext task for guiding the training of a trajectory encoder is central to the pre-training of general trajectory embeddings. The resulting process facilitates the extraction of universal insights from unlabeled trajectory data, encompassing elements of travel semantics such as purpose, destination, and route preference, alongside spatio-temporal movement patterns like speed and direction. Consequently, the resulting embeddings empower downstream analysis tasks that include prediction or classification. Thus, researchers are increasingly focusing on the pre-training of spatio-temporal trajectory embeddings. However, three challenges remain unsolved.

(1) **Reduced generalizability due to biases introduced into**

Corresponding author: Shengnan Guo

- Yan Lin, Huaiyu Wan, Shengnan Guo and Youfang Lin are with the Beijing Key Laboratory of Traffic Data Analysis and Mining, School of Computer and Information Technology, Beijing Jiaotong University, Beijing 100044, China, and the Key Laboratory of Intelligent Passenger Service of Civil Aviation, CAAC, Beijing, 101318, China. Jilin Hu and Christian S. Jensen are with the Department of Computer Science, Aalborg University, Aalborg, 9220, Denmark.
E-mail: ylincs@bjtu.edu.cn; huywan@bjtu.edu.cn; hujilin@cs.aau.dk; guoshm@bjtu.edu.cn; csj@cs.aau.dk; yflin@bjtu.edu.cn.

Manuscript received xx xx, xxxx; revised xx xx, xxxx.

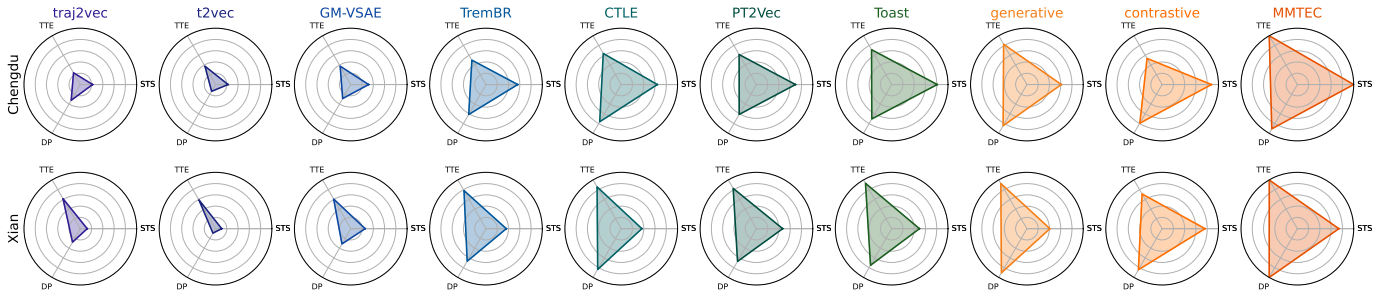


Fig. 1. Comparison of the performance of trajectory embeddings on two datasets and three downstream tasks. The axes represent the relative performance of an embedding on a downstream task. STS, TTE, and DP denote the normalized Acc@1 for similar trajectory search, 100% - MAPE for travel time estimation, and Acc@1 for destination prediction, respectively. The area of a polygon represents the generality of a trajectory embedding method.

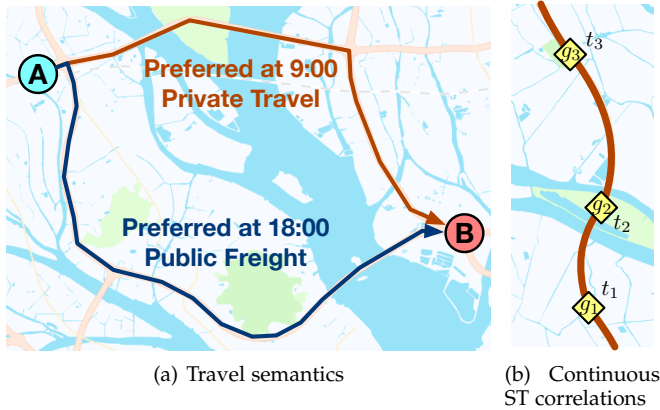


Fig. 2. Two aspects of information embedded in trajectories.

the learned embeddings. To learn general embeddings of spatio-temporal trajectories that can enhance prediction accuracy and accelerate convergence across diverse downstream tasks, self-supervised pretext tasks must avoid introducing biases that favor particular types of downstream tasks. In contrast, current pre-training methods employ pretext tasks that yield improved performance on certain downstream tasks, while causing reduced performance on others, as exemplified in Figure 1.

In particular, the proposals t2vec [19], traj2vec [6], GM-VSAE [7], TremBR [20], and PT2Vec [21] employ a generative auto-regressive-based pretext task, mapping a trajectory into a latent space and subsequently reconstructing its original form from this latent space using an auto-regressive approach. While this pretext task aims at capturing intricate trajectory detail [16], it can inadvertently introduce biases that favor point-level tasks such as travel-time estimation and destination prediction. Also, Toast [22] explores the idea of pre-training trajectory embeddings through sequence discrimination [15], [22]. However, this approach may introduce biases toward sequence-level tasks like similar trajectory search and trajectory-user linking.

(2) **Incomplete capture of informational aspects of trajectories.** To learn trajectory embeddings conducive to enhancing downstream task performance, it is paramount to capture a multitude of informational aspects of trajectories. These aspects can be categorized as either travel semantics and continuous spatio-temporal correlations. Tra-

jectories inherently capture movement in road networks and possess intrinsic semantic attributes that encapsulate travel intentions and route preferences [22], as exemplified in Figure 2(a). Furthermore, trajectories represent continuous movements across time spans, as exemplified in Figure 2(b). Yet, existing trajectory embedding methods predominantly focus on one of these aspects, reducing the ability of the learned embeddings to capture comprehensively the aspects of trajectory information.

More specifically, t2vec [19], traj2vec [6], and GM-VSAE [7] prioritize modeling spatio-temporal correlations in trajectories, while disregarding travel semantics derived from the underlying road networks. Next, TremBR [20], PT2Vec [21], and Toast [22] incorporate travel semantics when embedding road network-constrained trajectories, while their modeling of spatio-temporal correlations is comparatively less developed.

(3) **Difficulties at modeling long and irregularly sampled trajectories.** Long trajectories pose efficiency challenges to sequential models, and temporal irregularity in trajectories render it difficult to capture continuous spatio-temporal correlations.

Existing trajectory embedding methods frequently draw inspiration from Natural Language Processing (NLP) and time series modeling, notably leveraging sequential models. Examples include t2vec [19], traj2vec [6], GM-VSAE [7], TremBR [20], and PT2Vec [21] that all construct their trajectory encoders based on RNNs [23]. Additionally, Toast [22] utilizes a Transformer-based encoder [24]. However, RNN-based encoders suffer from the vanishing gradient issue [25], while Transformer-based encoders have quadratic time and memory complexity, rendering them suboptimal for capturing long-term correlations in trajectories. Conversely, both RNN and Transformer encoders solely update their hidden states upon receiving new input records, limiting their ability to capture the continuous dynamics characterizing trajectories [26].

To address the above challenges, we propose a novel trajectory embedding method called *Maximum Multi-view Trajectory Entropy Coding* (MMTEC). This method tackles the challenges as follows: (1) It leverages the principle of maximum entropy to learn unbiased, universal trajectory embeddings that benefits a wide range of downstream tasks; (2) It incorporates multiple views of trajectories into the maximum entropy loss, capturing both aspects of the

information embedded within trajectories; (3) It encompasses an attention-based discrete trajectory encoder featuring globally-shared anchors in order to capture travel semantics and a continuous trajectory encoder grounded in NeuralCDE to extract continuous spatio-temporal trajectory dynamics.

In summary, the paper’s main contributions are as follows:

- We propose MMTEC, a method for pre-training spatio-temporal trajectory embeddings that learns general trajectory embeddings beneficial for the performance of various downstream tasks.
- We develop a trajectory multi-view scheme that incorporates both semantic and continuous spatio-temporal information into the learned embeddings.
- We build an attention-based discrete trajectory encoder and a NeuralCDE-based continuous trajectory encoder that can model different aspects of spatio-temporal correlations in trajectories comprehensively and efficiently.
- We report on an experimental study involving two real-world trajectory datasets and three downstream tasks, finding that MMTEC embeddings generalize well and improve downstream tasks’ performance.

The remainder of the paper is organized as follows: Section 2 reviews related work. Section 3 covers the problem statement, definitions, and basic maximum entropy coding framework. Section 4 details the implementation of our proposed method. Section 5 reports on the experimental study, and Section 6 concludes the paper.

2 RELATED WORK

2.1 Learning Embeddings for Spatio-Temporal Data Mining

Embedding learning is an inherent aspect of neural network data mining, as most neural network models are feature-based and require the input targets to be represented by latent embedding tensors. In spatio-temporal data mining, embedding learning is also crucial. Spatio-temporal objects, such as locations, road networks, and trajectories, often contain complex spatio-temporal information while being geographically and dynamically correlated with other objects. Capturing and fusing such information into embeddings is challenging.

Locations are fundamental units in spatio-temporal datasets, and location embedding techniques are used widely in spatio-temporal data mining models. A straightforward idea is to use an index-fetching embedding layer [1], [8], as locations such as POIs and road segments are often represented by discrete indices. However, this kind of embedding is trained through task-specific supervision, cannot properly capture spatio-temporal correlations between locations comprehensively, and often requires a large labeled dataset to work well.

To tackle the above problems, several methods pre-train embeddings of POIs utilizing large, unlabeled datasets. For example, DeepMove [27] proposes word2vec [28] that works on unlabeled movement data, while POI2Vec [29]

and TALE [30] incorporate spatial and temporal correlations between locations into the embedding vectors. Next, CTLE [31] applies the idea of contextual embedding from language models to location embedding. These pre-training methods support downstream tasks such as location recommendation [32], [33] and location classification [34], [35].

Road segments and intersections in road networks can also be seen as locations. They contain complex information inherited from the topological structure of the road network. Using only an index-fetching embedding layer for road network representation will eliminate the important topological information. Therefore, methods such as IRN2Vec [36], Toast [22], and HRNR [37] incorporate topological correlations in road segment embedding by coupling random walks on road networks with word2vec, by utilizing real-world trajectories to fuse network topology with traffic patterns, and by proposing a hierarchical graph neural network and an auto-encoding pre-training objective to learn more comprehensive embeddings of road segments.

In this paper, we focus on the pre-training of trajectory embeddings, thus emphasizing the travel semantics and spatio-temporal movements in a trajectory sequence as a whole.

2.2 Pre-training Trajectory Embeddings

The quality of trajectory embeddings is crucial for achieving good performance of trajectory analysis tasks. The most widely applied trajectory embedding methodology is to construct an end-to-end sequence encoder for generating latent representations [1], [8]. However, the resulting trajectory embeddings often do not generalize and are difficult to migrate to other models or tasks. This methodology also assumed that substantial labeled training data is available, which may not be the case for some downstream tasks.

In part to address these issues, there is a growing interest in the pre-training of trajectory embeddings with self-supervised learning objectives. The notion of pre-training embeddings via self-supervised objectives stems from Language Modeling (LM) [13], [38] and Computer Vision (CV) [15], [16], domains that embrace self-supervised pre-training techniques. For instance, t2vec [19] builds an RNN-based trajectory encoder and applies it to GPS sequences to infer latent representations of trajectories based on sequential correlations of coordinates in the trajectories. GM-VSAE [7] further maps trajectories into a Gaussian latent space, so that the distributions of trajectories can be modeled. Traj2vec [6] considers relative spatio-temporal properties, using the differences in a set of attributes between consecutive trajectory points as input features. TremBR [20] models travel semantics by mapping trajectories onto road networks and additionally incorporates temporal information by concatenating the travel time of each point with its latent embedding vector. PT2Vec [21] partitions road networks into sub-networks to model map-matched trajectories more effectively.

The above methods build pretext tasks based on generative objectives, including auto-encoding [39] and auto-regression [40], where a decoder is trained to recover the trajectories given their embeddings. Since the generative objectives focus on modeling the detailed sequential features,

the learned embeddings are usually biased towards point-level tasks, such as travel time estimation.

Recently, Toast [22] proposed a trajectory embedding pre-training method that combines a masked language model (MLM) [13] pretext task with a sequence discrimination contrastive pretext objective [13], [15]. Trajectory embeddings learned by contrastive objectives tend to be biased towards sequence-level tasks, such as trajectory-user linking. It is evident that existing pretext tasks are inevitably biased towards certain types of downstream tasks, which conflicts with the goal of learning general trajectory embeddings that performs optimally across diverse tasks.

We note that several existing methods were initially tailored to address specific tasks within their original contexts. However, these methods are united by their underlying adherence to self-supervised pre-training paradigms, inherently enabling them to be used in other downstream tasks. This paper's coverage is not confined to the initial application scopes of these related studies. Instead, we focus on the technical frameworks that form the core of these methods. A comparison of three attributes of the methods is provided in Table 2.

3 PRELIMINARIES

3.1 Definitions

Definition 1 (Trajectory). A trajectory is a sequence of timestamped GPS points $\mathcal{T} = \langle (g_1, t_1), (g_2, t_2), \dots, (g_{|\mathcal{T}|}, t_{|\mathcal{T}|}) \rangle$, where $g_i = (\text{lng}_i, \text{lat}_i)$ denotes the i -th GPS point, and $|\mathcal{T}|$ is the total number of GPS points in the trajectory.

Definition 2 (Map-matched Trajectory). A map-matched trajectory is a sequence of timestamped road segments $\mathcal{R} = \langle (s_1, t_1), (s_2, t_2), \dots, (s_{|\mathcal{R}|}, t_{|\mathcal{R}|}) \rangle$, where $s_i \in \mathbb{S}$ denotes the i -th road segment, \mathbb{S} is the set of all road segments in the area of interest, and $|\mathcal{R}|$ is the total number of road segments in the trajectory. We can obtain the corresponding \mathcal{R} of a given \mathcal{T} using map-matching methods, such as Markov-based methods [41].

3.2 Problem Statement

Pre-training Trajectory Embeddings. Given a set of spatio-temporal trajectories, we aim to pre-train a trajectory encoder f to generate an embedding vector $e_{\mathcal{T}} \in \mathbb{R}^d$ when given a trajectory \mathcal{T} , i.e., $e_{\mathcal{T}} = f(\mathcal{T})$, where d is the embedding dimensionality. Our focus is on self-supervised pre-training of the encoder. The goal is to achieve embeddings that possess adaptability across a variety of downstream tasks. Specifically, the learned embeddings aim to enhance the convergence speed and accuracy in diverse trajectory mining tasks, rendering them effective tools in an array of applications.

3.3 Maximum Entropy Coding

Self-supervised pretext tasks are crucial for pre-training embedding methods. To ensure that embeddings generalize well across diverse downstream tasks, it is essential to minimize biases introduced during pre-training. The information theoretical principle of maximum entropy can be used for

measuring the generality of embeddings. According to this principle, the probability distribution that best represents the current state of knowledge about a system is the one with the largest entropy, in the context of precisely stated prior data [42]. Leveraging this principle, maximum entropy coding (MEC) [43] aims to maximize the entropy of embeddings under a specific prior, with the end goal of producing embeddings that enhance performance across diverse tasks.

Suppose we have a set of embeddings $\mathbf{E} \in \mathbb{R}^{d \times N}$, where d is the embedding dimensionality and N is the number of vectors. To provide a tractable formulation of the entropy of the embeddings, MEC proposes the following coding length function:

$$L = \frac{N + d}{2} \log \det(\mathbf{I}_N + \frac{d}{N\epsilon^2} \mathbf{E}^\top \mathbf{E}) \quad (1)$$

Here, \mathbf{I}_N is an identity matrix with dimension N , and ϵ is the upper bound of the decoding error. This function calculates the coding length of a lossy data coding, which can be seen as an estimate of the entropy of the embeddings \mathbf{E} . The training objective under the MEC framework is to maximize Equation 1.

4 METHODOLOGY

4.1 Overview

Figure 3 shows the proposed framework. Given a spatio-temporal trajectory, our goal is to develop a self-supervised pretext task using maximum entropy coding, with the aim of maximizing the information entropy of the trajectory's embeddings. To incorporate multiple aspects of information embedded within the trajectory, we propose a multi-view scheme for modeling and fusing two aspects of information.

Specifically, we map-match the trajectory and use an attention-based discrete trajectory encoder to efficiently extract correlations from the map-matched trajectory, mapping it into the semantic embedding. We also recover the underlying continuous movements of the trajectory and use a NeuralCDE-based continuous trajectory encoder to model the recovered movements and map these into the continuous spatio-temporal embedding. We incorporate these two embeddings into our pretext task by defining a view consistency prior constraint, which we use to pre-train the two embeddings with our final pretext task. After pre-training, we combine the two embeddings for use in downstream tasks. The remainder of the section provides a detailed explanation of the proposed method.

4.2 Pre-training General Embeddings by Maximizing Entropy

Current pre-training methods for trajectory embedding utilize generative or contrastive losses in their pretext tasks, which can introduce biases and restrict the generality of the learned embeddings to apply across different types of downstream tasks. To address this limitation, we draw inspiration from maximum entropy coding (MEC) as discussed in Section 3.3 and develop a pretext task that minimizes the bias introduced into pre-trained embeddings.

Suppose we have a set of trajectories $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_N\}$ and the corresponding embeddings $\mathbf{E}_{\mathcal{T}} = \{e_{\mathcal{T}_1}, e_{\mathcal{T}_2}, \dots, e_{\mathcal{T}_N}\}$, where $e_{\mathcal{T}_i} \in \mathbb{R}^d$. To minimize

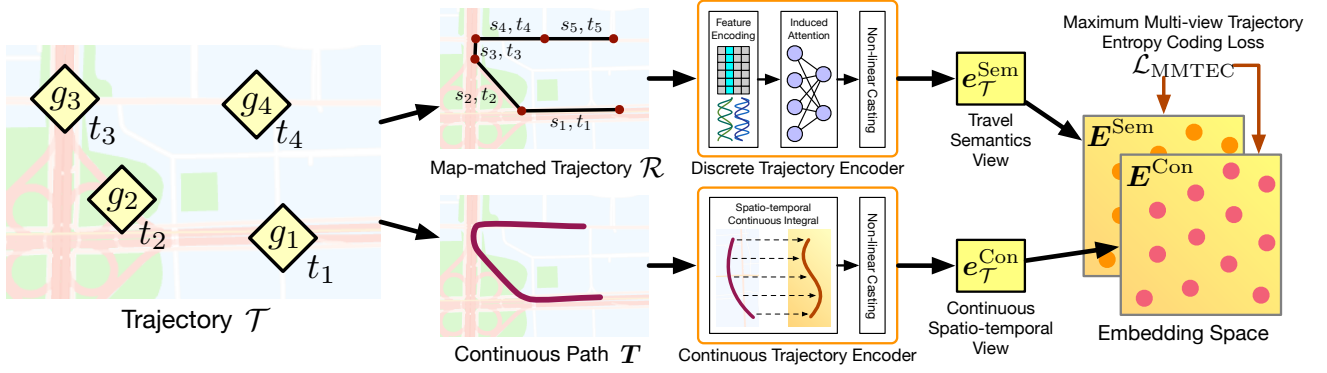


Fig. 3. The overall framework of our proposed model.

bias in the learned embeddings, we utilize MEC to maximize the entropy of the embeddings. Therefore, the negative of Equation 1 is used as our pre-training loss function. We simplify Equation 1 using Taylor expansion following [43] to obtain a more tractable loss function:

$$\mathcal{L} = -\text{trace}\left(\frac{N+d}{2} \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \left(\frac{d}{N\epsilon^2} \mathbf{E}_{\mathcal{T}}^{\top} \mathbf{E}_{\mathcal{T}}\right)^k\right) \quad (2)$$

According to the principle of maximum entropy, the resulting embeddings $\mathbf{E}_{\mathcal{T}}$ will have the maximum entropy among all possible embeddings. This characteristic facilitates optimal accuracy and convergence performance across diverse downstream tasks.

Recall that the principle of maximum entropy requires a precisely stated prior. Simply pre-training trajectory embeddings with Equation 2 without a well-defined prior will lead to all trajectories being represented uniformly [42]. Additionally, training only one form of trajectory embedding limits the scope of their applicability. For example, embedding only the map-matched trajectories captures only the semantic information of a road network, disregarding the continuous spatio-temporal correlations of trajectories. A narrow focus can limit the performance of the embeddings on downstream tasks.

Based on these considerations, we propose a multi-view trajectory scheme that models different aspects of the information in trajectories. This allows us to introduce a testable prior based on the similarity between the different views, leading to more comprehensive and effective trajectory embeddings.

4.3 Multi-view Trajectory Scheme

We propose a multi-view trajectory scheme that aims to incorporate multiple aspects of trajectory information into embeddings. Specifically, we seek to incorporate both travel semantics and continuous spatio-temporal correlations. Instead of augmenting trajectories into pairs of samples, we model the two aspects separately using specifically designed trajectory encoders and construct intermediate trajectory embeddings.

4.3.1 Embedding Travel Semantics

Obtaining semantically rich representations using map-matching. To extract travel semantics from trajectories, we

first map the trajectories onto the road network to utilize the semantic information in road segments. We use the Fast Map Matching (FMM) algorithm [41] to obtain the corresponding map-matched trajectory \mathcal{R} for a given trajectory \mathcal{T} , shown as:

$$\mathcal{R} = \text{FMM}(\mathcal{T}) \quad (3)$$

The resulting \mathcal{R} is a sequence of timestamped road segments, which is a semantically rich representation of the original \mathcal{T} . However, \mathcal{R} is not appropriate as an embedding of \mathcal{T} because it contains both discrete features (i.e., road segment indices) and numerical features (i.e., timestamps and road segment coordinates), and because different \mathcal{R} have different lengths.

To address this, we propose an attention-based discrete encoder with a feature encoding layer and induced attention layers. This encoder efficiently and effectively models the correlations between road segments in \mathcal{R} and maps \mathcal{R} into a fixed-length embedding, as shown in Figure 4. This encoder is realized as follows.

Encoding features in timestamped road segments. Each timestamped road segment in \mathcal{R} contains several features: the road segment index s_i , the timestamp t_i when the road segment is visited, and the spatial coordinates $(\text{lng}_{s_i}, \text{lat}_{s_i})$ of the road segment.

To embed the discrete road segment index s_i , we provide an index fetching embedding module so that the inherent semantic information in road segments can be captured and learned. Specifically, we initialize an embedding matrix $\mathbf{E}^{\text{Road}} \in \mathbb{R}^{d \times |\mathcal{S}|}$, where each column $\mathbf{E}_{s_i}^{\text{Road}}$ corresponds to the embedding vector of road segment s_i .

We then introduce a trigonometric function $\Psi : \mathbb{R} \rightarrow \mathbb{R}^d$ to encode the numerical features $t_i, \text{lng}_{s_i}, \text{lat}_{s_i}$, thereby emphasizing their periodic characteristics, for example, the similarity in travel behavior during the same period of the day. The inspiration for this approach stems from both the positional encoding used in the Transformer [24] and the temporal encoding proposed in recent studies [31], [44]. The encoding function $\Psi(v)$ is defined as follows:

$$\Psi(v) = (\cos(\omega_1 v), \sin(\omega_1 v), \dots, \cos(\omega_{d/2} v), \sin(\omega_{d/2} v)), \quad (4)$$

where $v \in \{t_i\} \cup \{\text{lng}_{s_i}\} \cup \{\text{lat}_{s_i}\}$, $i \in \{1, 2, \dots, |\mathcal{R}|\}$ denotes the input numerical feature and $\omega_1, \omega_2, \dots, \omega_{d/2}$ are learnable parameters. Importantly, distinct sets of parameters are

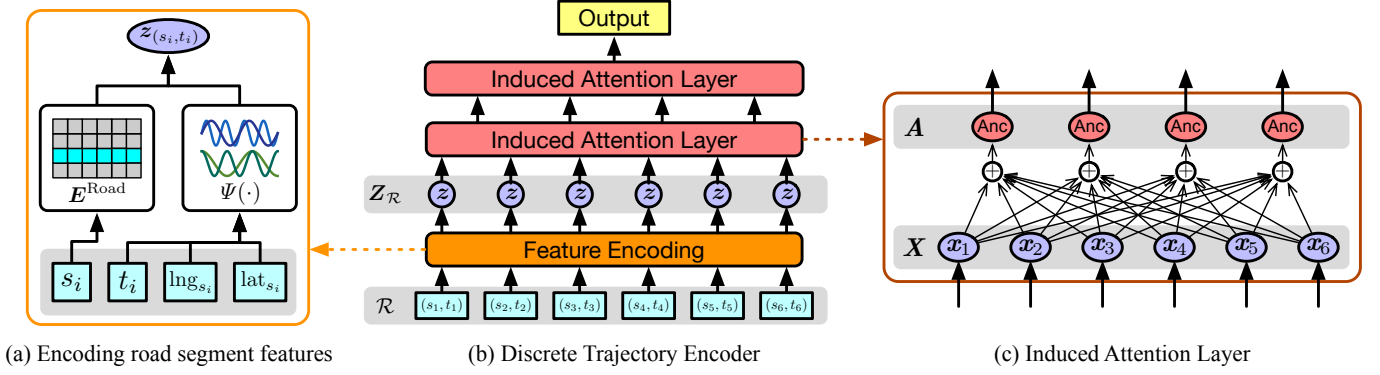


Fig. 4. The architecture of our proposed discrete trajectory encoder and its components.

employed for different types of input features. By using Ψ , the encoder adeptly captures the spatial and temporal periodic characteristics of the numerical features via the sin and cos functions. Additionally, this encoding enhances the model's capture of relative correlations between features of different segments in \mathcal{R} . This is predominantly due to the central role played by dot-product computations in encoder architectures, including multi-head attention [24]. Specifically, the dot-product of two encoding vectors $\Psi(v)$ and $\Psi(v + \delta)$ is calculated as:

$$\Psi(v) \cdot \Psi(v + \delta) = \cos(\omega_1 \delta) + \cos(\omega_2 \delta) + \dots + \cos(\omega_{d_v/2} \delta) \quad (5)$$

This calculation illustrates that the distance (measured by the dot-product) between $\Psi(v)$ and $\Psi(v + \delta)$ relies exclusively on their numerical difference δ , obviating the influence of offset v [45], [46]. Furthermore, it shows that the distance information is learnable, guided by the set of parameters $\{\omega_1, \dots, \omega_{d/2}\}$. Our encoder then effectively captures this distance information, enabling the modeling of correlations between timestamped road segments within \mathcal{R} , such as the time gap between the traversal of two road segments.

Finally, we obtain the latent embedding of each timestamped road segment (s_i, t_i) , denoted as:

$$\mathbf{z}_{(s_i, t_i)} = \mathbf{E}_{s_i}^{\text{Road}} + \Psi(t_i) + \Psi(\text{lng}_{s_i}) + \Psi(\text{lat}_{s_i}) \quad (6)$$

We can then transform a map-matched trajectory \mathcal{R} into a sequence of latent embeddings:

$$\mathbf{Z}_{\mathcal{R}} = \langle \mathbf{z}_{(s_1, t_1)}, \mathbf{z}_{(s_2, t_2)}, \dots, \mathbf{z}_{(s_{|\mathcal{R}|}, t_{|\mathcal{R}|})} \rangle \in \mathbb{R}^{d \times |\mathcal{R}|} \quad (7)$$

Efficient sequence encoder based on induced attention.

To capture the correlation between road segments using a latent embedding sequence $\mathbf{Z}_{\mathcal{R}}$, an intuitive approach is to use the Transformer encoder, which has been adopted widely. However, the self-attention mechanism in the Transformer incurs quadratic time and memory complexity, meaning that it does not scale well to long sequences. Instead, we draw inspiration from recent developments in natural language processing for sequential models [47]. Specifically, we use an attention mechanism that incorporates long-term correlations and has low computational complexity. Our attention scheme, which we refer to as induced attention, utilizes globally-shared anchor sequences for attention weight calculation.

For the induced attention layer, we initialize an anchor sequence $\mathbf{A}_i \in \mathbb{R}^{d \times L_{A_i}}$ as the inherent parameter of this layer, where i denotes the i -th layer and L_{A_i} is the length of the anchor sequence. Then, the induced attention layer IA_i is realized using attention and feed-forward networks:

$$\begin{aligned} \text{IA}_i(\mathbf{X}_i) &= \text{Norm}(\mathbf{H}_i + \text{FFN}_i(\mathbf{H}_i)) \\ \mathbf{H}_i &= \text{Norm}(\mathbf{A}_i + \text{Att}_i(\mathbf{A}_i, \mathbf{X}_i, \mathbf{X}_i)), \end{aligned} \quad (8)$$

where Norm , FFN_i , and Att_i represent the layer normalization [48], the feed-forward network, and the multi-head dot-product attention of i -th layer, respectively. We use 8 attention heads. Next, $\mathbf{X}_i \in \mathbb{R}^{d \times L_{X_i}}$ is the input to i -th layer. During attention calculation, \mathbf{A}_i is regarded as the query, while \mathbf{X}_i is regarded as the key and value. The parameters in FFN_i , Att_i and the anchor sequence \mathbf{A}_i are shared for every specific IA layer between mini-batches, but are not shared between different layers.

We build our discrete encoder DisEnc by stacking two IA layers. Given a latent embedding sequence $\mathbf{Z}_{\mathcal{R}}$, the implementation of DisEnc is defined as:

$$\text{DisEnc}(\mathcal{R}) = \text{IA}_2(\text{IA}_1(\mathbf{Z}_{\mathcal{R}})), \quad (9)$$

where the anchor sequence length of the second IA layer L_{A_2} is set to 1 so that the output dimensionality of DisEnc is d .

By sharing the anchor sequences across mini-batches for attention score calculation rather than adopting self-attention, we build a more efficient discrete encoder that is suitable for modeling long-term correlations in map-matched trajectories. The computational complexity of Equation 9 is $O(|\mathcal{R}|)$. In contrast, the complexity of a vanilla Transformer encoder is $O(|\mathcal{R}|^2)$, which does not scale well to large $|R|$. An empirical study of the efficiency improvement achieved by the proposed model is covered in Section 5.5.2.

Final semantic embedding. We obtain the final semantic embedding by appending a non-linear casting to the output of DisEnc :

$$\mathbf{e}_{\mathcal{T}}^{\text{Sem}} = \mathbf{W}_2(\sigma(\mathbf{W}_1 \text{DisEnc}(\mathcal{R}))), \quad (10)$$

where $\mathbf{e}_{\mathcal{T}}^{\text{Sem}} \in \mathbb{R}^d$ denotes the final semantic embedding of trajectory \mathcal{T} , representing the travel semantics view. $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times d}$ are mapping matrices, and σ denotes a non-linear activation function—we use ReLU.

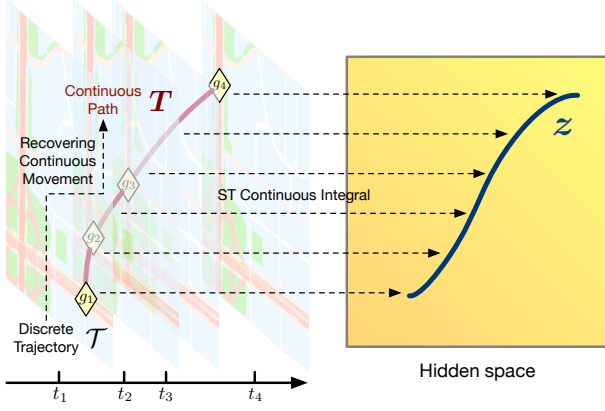


Fig. 5. Illustration of the continuous trajectory encoder.

4.3.2 Embedding Continuous Spatio-Temporal Correlations.

Recovering continuous movements of trajectories. To properly model the spatio-temporal correlations in trajectories, it is essential to account for their continuous spatial and temporal dynamics. However, the irregular and discrete nature of GPS points in trajectories causes challenges for conventional sequential models such as RNNs [23] and Transformers [24]. Instead, we adopt an approach that can inherently consider the underlying continuous paths of trajectories, as shown in Figure 5.

We first recover the continuous movements of trajectories based on their GPS points. Given a trajectory $\mathcal{T} = \langle (g_1, t_1), (g_2, t_2), \dots, (g_{|\mathcal{T}|}, t_{|\mathcal{T}|}) \rangle$, we approximate the underlying continuous movements of the trajectory by calculating its cubic spline:

$$\mathbf{T} = \text{Spline}(\mathcal{T}), \quad (11)$$

where Spline is the cubic Hermite spline that interpolates the underlying continuous process of the trajectory. In particular, \mathbf{T} is the natural cubic spline with knots at $t_1, t_2, \dots, t_{|\mathcal{T}|}$ that satisfies:

$$\begin{aligned} \mathbf{T} : t &\rightarrow \mathbb{R}^2, t \in [t_1, t_{|\mathcal{T}|}], \\ \mathbf{T}_{t_i} &= g_i, i \in \{1, 2, \dots, |\mathcal{T}|\} \end{aligned} \quad (12)$$

Since \mathbf{T} estimates the underlying continuous movement of a trajectory \mathcal{T} , it can serve as an intermediate step for mapping \mathcal{T} to the continuous spatio-temporal embedding. However, \mathbf{T} is unsuitable for most trajectory analysis tasks due to its continuous definition. To address this, we propose a continuous encoder ConEnc that can capture the continuous dynamics represented by \mathbf{T} .

Spatio-temporal continuous integral. To construct the continuous encoder ConEnc, we leverage the family of ordinary differential equation methods, specifically the neural controlled differential equation (NeuralCDE) [49]. This method captures the underlying sequential processes in continuous time and models both the input features and hidden states in

a continuous manner. We adopt an existing formulation [49] and define our continuous encoder as follows:

$$\begin{aligned} \text{ConEnc}(\mathbf{T}) &= \mathbf{z}_{t_1} + \int_{t_1}^{t_{|\mathcal{T}|}} \text{CDENet}(\mathbf{z}_t) d\mathbf{T}_t \\ &= \mathbf{z}_{t_1} + \int_{t_1}^{t_{|\mathcal{T}|}} \text{CDENet}(\mathbf{z}_t) \frac{d\mathbf{T}}{dt}(t) dt \end{aligned} \quad (13)$$

Here, $\frac{d\mathbf{T}}{dt}$ denotes the derivative of the cubic spline with regard to time, which we implement by employing the torchdiffeq package [50]. Next, $\mathbf{z}_{t_1} = \text{InitNet}(g_1, t_1)$ is a parameterized initial state, $\text{InitNet} : \mathbb{R}^3 \rightarrow \mathbb{R}^d$ denotes the neural network used for calculating the initial state, and $\text{CDENet} : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times 3}$ is a neural network that serves as the integral kernel of CDE. InitNet and CDENet are realized as follows.

Parameterized CDE initial state and integral kernel. The initial state \mathbf{z}_{t_1} determines the starting point of the CDE integral and depends on the first timestamped GPS point (g_1, t_1) of a trajectory. Since all features in a GPS point are numerical, we calculate the initial state using the trigonometric function from Equation 4. Thus, the neural network InitNet is defined as:

$$\text{InitNet}(g_1, t_1) = \Psi(\text{lng}_1) + \Psi(\text{lat}_1) + \Psi(t_1) \quad (14)$$

The neural network CDENet serves as the CDE integral kernel. Here, we use a two-layer non-linear mapping:

$$\text{CDENet}(\mathbf{z}_t) = \mathbf{W}_6(\sigma(\mathbf{W}_5 \mathbf{z}_t + \mathbf{b}_5)) + \mathbf{b}_6, \quad (15)$$

where $\mathbf{W}_5 \in \mathbb{R}^{3d \times d}$ and $\mathbf{W}_6 \in \mathbb{R}^{3d \times 3d}$ are mapping matrices and $\mathbf{b}_5, \mathbf{b}_6 \in \mathbb{R}^{3d}$ are biases.

Final continuous spatio-temporal embedding. Akin to Equation 10, we pair the continuous encoder with a non-linear casting to obtain the continuous spatio-temporal embedding vector:

$$\mathbf{e}_{\mathcal{T}}^{\text{Con}} = \mathbf{W}_4(\sigma(\mathbf{W}_3 \text{ConEnc}(\mathbf{T}))), \quad (16)$$

where $\mathbf{e}_{\mathcal{T}}^{\text{Con}}$ denotes the final continuous spatio-temporal embedding of trajectory \mathcal{T} , representing the continuous spatio-temporal view and $\mathbf{W}_3, \mathbf{W}_4 \in \mathbb{R}^{d \times d}$ are mapping matrices.

4.3.3 Combining the Intermediate Embeddings

The final embedding of a trajectory \mathcal{T} is obtained by combining the semantic and continuous spatio-temporal embeddings as follows:

$$\mathbf{e}_{\mathcal{T}} = \mathbf{e}_{\mathcal{T}}^{\text{Sem}} \oplus \mathbf{e}_{\mathcal{T}}^{\text{Con}}, \quad (17)$$

where \oplus denotes the concatenation operation along feature dimension. The resulting embeddings incorporate both informational aspects captured in trajectories. This enables a comprehensive modeling of the informational aspects of trajectories, yielding embeddings that are more capable of enhancing the prediction performance of downstream tasks.

4.4 Maximum Multi-view Trajectory Entropy Coding

To ensure that we satisfy the premise of the principle of maximum entropy and learn meaningful trajectory embeddings, we introduce testable prior information to Equation 2 based on the two embeddings obtained in Equations 10 and 16.

Given a mini-batch of trajectories $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_N\}$ used during batch gradient descent training, where N is the batch size, we obtain their corresponding semantic embeddings $\mathbf{E}^{\text{Sem}} = \{e_{\mathcal{T}_1}^{\text{Sem}}, e_{\mathcal{T}_2}^{\text{Sem}}, \dots, e_{\mathcal{T}_N}^{\text{Sem}}\} \in \mathbb{R}^{d \times N}$ and their corresponding continuous spatio-temporal embeddings $\mathbf{E}^{\text{Con}} = \{e_{\mathcal{T}_1}^{\text{Con}}, e_{\mathcal{T}_2}^{\text{Con}}, \dots, e_{\mathcal{T}_N}^{\text{Con}}\} \in \mathbb{R}^{d \times N}$. As the two types of embeddings correspond to the same set of trajectories, they can be regarded as different views on that set. To incorporate the view consistency prior between \mathbf{E}^{Sem} and \mathbf{E}^{Con} into Equation 2, we obtain:

$$\mathcal{L}' = -\text{trace}\left(\frac{N+d}{2} \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \left(\frac{d}{N\epsilon^2} \mathbf{E}^{\text{Sem}\top} \mathbf{E}^{\text{Con}}\right)^k\right) \quad (18)$$

To make the calculation feasible, we approximate Equation 18 by truncating the infinite Taylor expansion, obtaining the following MMTEC loss function:

$$\mathcal{L}_{\text{MMTEC}} = -\text{trace}\left(\frac{N+d}{2} \sum_{k=1}^K \frac{(-1)^{k+1}}{k} \left(\frac{d}{N\epsilon^2} \mathbf{E}^{\text{Sem}\top} \mathbf{E}^{\text{Con}}\right)^k\right) \quad (19)$$

Here, K is the order of the Taylor expansion, which is a hyper-parameter. This loss function maximizes the entropy of embeddings while ensuring view consistency. Complemented by the incorporation of both aspects of trajectory information, as captured in Equation 17, the learned embeddings are able to enhance the convergence speed and accuracy of a range of downstream tasks.

5 EXPERIMENTS

To evaluate proposed method for learning trajectory embeddings, we conduct experiments on two real-world datasets and apply the learned embeddings in three downstream tasks. We compare our results with those obtained using other state-of-the-art methods.

5.1 Datasets

We conduct the experiments on two real-world datasets, Chengdu and Xian, that contain taxi trip trajectories collected from the central areas of the respective cities. The datasets are provided by DiDi¹. To ensure a reasonable size of the dataset, we sort the drivers descendingly according to how many trajectories they are associated with. We then select the trajectories produced by the top 20,000 drivers. We also download the road networks of Chengdu and Xian from OpenStreetMap (OSM) [51] for use in the map matching presented in Equation 3. We retain only the road segments covered by at least one trajectory. Table 1 presents statistics of the datasets.

TABLE 1
Statistics of datasets.

Dataset	#Trajectories	#Segments	#Records
Chengdu	298,995	3,791	7,025,468
Xian	376,407	3,558	10,198,837

5.2 Comparison Methods

To assess the benefits of the proposed model over the state-of-the-art, we include comparisons with six state-of-the-art pre-training trajectory embedding methods.

- **traj2vec** [6]: constructs feature sequences by calculating the differences in the spatio-temporal attributes between consecutive trajectory points and applies an auto-regressive pretext task to learn embeddings.
- **t2vec** [19]: pre-trains an RNN trajectory encoder based on a de-noising auto-encoder, aiming to recover fine-grained trajectories given sparse ones.
- **GM-VSAE** [7]: casts trajectories into a multi-dimension Gaussian space using an RNN variational trajectory encoder and pre-trains the model following variational auto-encoders.
- **TremBR** [20]: constructs an RNN auto-encoder to learn trajectory embeddings with both road segment and timestamp recovery tasks.
- **CTLE** [31]: pre-trains a bi-directional Transformer with two MLM-based tasks and then calculates the embeddings of locations based on their contexts. We apply a mean pooling on the embeddings of all points in a trajectory to get its embedding.
- **PT2Vec** [21]: combines a road network partitioning and an RNN auto-encoder to effectively learn trajectory embeddings on large-scale road networks.
- **Toast** [22]: applies the node2vec model to the road network to pre-train a set of embeddings for the road segments and then pre-trains a Transformer encoder with an MLM-based task and a sequence discrimination task to generate trajectory embeddings.

Table 2 offers a comparison of three fundamental attributes of the pre-training trajectory embedding methods considered in the experiments. In the table, the column labeled *Initial task* refers to the downstream tasks the methods targeted in within their respective publications, while the column titled *Pre-training framework* denotes the pre-training approach embraced by each method.

TABLE 2
Comparison of pre-training trajectory embedding methods.

Methods	Encoder structure	Initial task	Pre-training framework
traj2vec	RNN	Clustering	Auto-encoding
t2vec	RNN	Similarity Measurement	Auto-encoding
GM-VSAE	RNN	Anomaly Detection	Variational AE
TremBR	RNN	Embedding Learning	Auto-encoding
CTLE	Transformer	Destination Prediction	MLM
PT2Vec	RNN	Similarity Measurement	Auto-encoding
Toast	Transformer	Embedding Learning	MLM+Contrastive
MMTEC	Multi-view	Embedding Learning	Max Entropy

1. The dataset was published at <https://gaia.didichuxing.com/>.

5.3 Downstream Tasks

To assess the generality of the trajectory embeddings learned by our method, we evaluate their performance on three downstream tasks: similar trajectory search, travel time estimation, and destination prediction.

Similar trajectory search. Our goal is to rank similar trajectories based on their distances in the learned embedding space. Since ground truth labels are not available, we adopt the evaluation strategy from Toast [22] and randomly select 10,000 trajectories from the full dataset as the target trajectories, denoted as \mathbb{T} . For each $\mathcal{T}_i \in \mathbb{T}$, we randomly select a 20% portion, and replace it with a shortest path on the road network from the origin to the destination of the portion, ensuring it is distinct from the original portion, this way obtaining a positive sample \mathcal{T}'_i . We then calculate the embedding $e_{\mathcal{T}'_i}$ of \mathcal{T}'_i and determine its similarity with all target trajectories $\mathcal{T}_j \in \mathbb{T}$ in the embedding space using cosine similarity. The target trajectory most similar \mathcal{T}'_i , i.e., the ground truth, should be \mathcal{T}_i .

Travel time estimation. Our aim is to estimate the travel time of trajectories using their embeddings. We feed the embedding methods with trajectories excluding timestamps and then pass the resulting embeddings to a fully-connected network for travel time regression.

Destination prediction. Our aim is to predict the destination road segments for trajectories using their embeddings. We feed the embedding methods a map-matched trajectory \mathcal{R} excluding its last 5 segments. The output embedding vector is then passed to a fully-connected network to classify the last road segment index of \mathcal{R} .

TABLE 3
Hyper-parameter ranges and optimal values.

Parameter	Range
d	32, <u>64</u> , 128, 192, 256
N	128, 256, 384, <u>512</u> , 640, 768, 896, 1024
ϵ^2	128, 256, 384, <u>512</u> , 640, 768, 896, 1024
L_{A_1}	0, 4, 8, 12, 16, 20, 24, 28, 32
K	1, 2, 3, 4, <u>5</u> , 6, 7, 8, 9, 10

Underline denotes the optimal value.

5.4 Settings

For all datasets, we sort trajectories by their start time and split the whole trajectory set into training, evaluation, and testing sets by 8:1:1. Both pre-training embedding methods and downstream predictors are trained with the training set. The embedding methods are pre-trained for 30 epochs, while the downstream predictors are stopped early on the evaluation set. The final metrics are calculated on the testing set. We use Top- N accuracy (i.e., $\text{Acc}@N$, $N = 1, 5$) for the similar trajectory search task; $\text{Acc}@N$, $N = 1, 5$, and macro-F1 for the destination prediction task; and mean absolute error (MAE), root mean square error (RMSE), and mean absolute percentage error (MAPE) for the travel time estimation task. We run each set of experiments 10 times and report the mean and deviation of the metrics.

All models are implemented using PyTorch [52].² The proposed method has five key hyper-parameters, whose ranges and optimal values are listed in Table 3. We choose parameters based on the $\text{Acc}@1$ results of the similar trajectory search task on the validation set of the Chengdu dataset. We report the effectiveness of these parameters in Section 5.5.5.

5.5 Experimental Results

5.5.1 Overall Performance

Tables 4 and 5 present a comprehensive comparison of the performance of different trajectory embedding methods across the three tasks. We also provide a visual illustration of the results in Figure 1. Our proposed method consistently outperforms the other methods, offering evidence that it is the best at generating general and comprehensive trajectory embeddings.

Generality of the pretext task. traj2vec, t2vec, GM-VSAE, TremBR, and PT2Vec all adopt the auto-encoding or auto-regressive frameworks. Among these, t2vec re-samples the original trajectories and adds noise to the encoder input to make the learned embeddings more robust to sparse and noisy trajectories. GM-VSAE uses a multi-dimensional Gaussian space as its latent embedding space and implements a VAE-like encoder-decoder architecture [53] to better model the trajectory distribution. However, auto-regressive pre-training methods focus mainly on recovering the detailed raw features of trajectories, biasing the learned embeddings towards point-level downstream tasks, such as travel time estimation and destination prediction. As the results show, auto-regressive pre-trained embedding methods perform poorly on sequence-level tasks, like similar trajectory search.

CTLE and Toast use the masked language model (MLM) pretext task from BERT [13]. Compared to the auto-regressive task, MLM utilizes the bi-directionality of Transformers and helps the model better understand the correlations between points in a trajectory. However, as a generative pretext task, it also biases the learned embeddings towards certain types of downstream tasks. Interestingly, Toast additionally implements a sequence-level pretext task that helps Toast perform better on the similar trajectory search task.

The pretext task of our proposed method is built based on the principle of maximum entropy, aiming to learn embeddings that generalize well across different types of downstream tasks. This makes our method better suited for pre-training general trajectory embeddings. We further investigate the effectiveness of pretext tasks by pre-training our embedding models generatively and contrastively—the detailed analysis is presented in Section 5.5.4.

Comprehensiveness of the incorporated information. Methods traj2vec, t2vec, and GM-VSAE are designed to model the spatio-temporal features of trajectories based on GPS coordinates and timestamps, but they do not consider the semantic information in the underlying road networks. This limitation reduces their performance on tasks such as

² The code and sample datasets are available in the supplementary material, and will be published on GitHub after the review.

TABLE 4
Performance comparison of trajectory embedding methods on the Chengdu dataset.

Task	Similar Trajectory Search		Travel Time Estimation			Destination Prediction		
Metric	Acc@1(%)	Acc@5(%)	MAE(min)	RMSE(min)	MAPE(%)	Acc@1(%)	Acc@5(%)	F1(%)
traj2vec	61.38±0.78	71.19±0.59	2.570±0.022	3.740±0.020	35.11±0.24	45.54±0.27	68.23±0.20	17.83±0.09
t2vec	61.52±1.14	71.70±1.56	2.518±0.024	3.695±0.031	33.49±0.30	39.52±0.13	63.36±0.40	13.51±0.12
GM-VSAE	64.28±0.87	76.17±0.49	2.515±0.033	3.732±0.041	33.49±0.32	44.37±0.63	67.30±0.86	17.64±0.11
TremBR	71.51±2.51	79.49±1.65	2.445±0.043	3.633±0.063	32.01±0.45	55.04±0.55	77.74±0.26	26.89±0.13
CTLE	73.75±1.41	82.69±1.62	2.341±0.035	3.558±0.071	30.37±0.51	59.88±0.26	79.92±0.13	27.66±0.13
PT2Vec	75.27±1.18	85.03±1.48	2.391±0.027	3.618±0.039	30.64±0.36	54.98±0.91	75.27±0.82	24.19±0.42
Toast	78.67±0.71	86.40±0.60	2.320±0.028	3.477±0.069	29.44±0.44	57.99±0.64	77.14±0.85	26.17±0.99
generative	72.89±0.84	83.20±1.21	2.239±0.042	3.341±0.033	28.10±0.93	62.52±0.12	89.26±0.14	30.04±0.14
contrastive	80.62±1.99	93.13±1.41	2.419±0.124	3.590±0.184	31.59±1.79	60.84±0.31	87.66±0.13	27.35±0.28
MMTEC	84.27±3.02	95.15±1.63	2.079±0.109	3.104±0.194	26.02±1.00	64.61±0.13	91.71±0.06	33.35±0.11

Bold denotes the best result, underline denotes the second best result.

TABLE 5
Performance comparison of trajectory embedding methods on the Xian dataset.

Task	Similar Trajectory Search		Travel Time Estimation			Destination Prediction		
Metric	Acc@1(%)	Acc@5(%)	MAE(min)	RMSE(min)	MAPE(%)	Acc@1(%)	Acc@5(%)	F1(%)
traj2vec	58.57±0.72	77.03±0.62	3.142±0.050	5.159±0.048	30.64±0.79	44.04±0.43	67.20±0.21	13.56±0.49
t2vec	58.19±1.27	73.78±1.14	3.178±0.016	5.248±0.011	30.93±0.30	37.99±0.16	62.47±0.41	8.73±0.34
GM-VSAE	62.47±0.30	78.24±1.03	3.143±0.032	5.123±0.039	30.73±0.43	45.21±0.89	68.82±0.88	13.73±0.79
TremBR	65.55±1.34	79.55±1.63	3.037±0.028	5.073±0.022	28.65±0.31	56.78±0.61	79.33±0.98	22.14±0.90
CTLE	65.74±1.95	79.26±1.25	2.949±0.042	4.912±0.108	27.74±0.58	62.25±0.53	85.21±0.53	27.17±0.81
PT2Vec	68.64±1.40	80.34±1.79	3.000±0.027	4.972±0.062	28.11±0.53	56.88±1.28	79.14±1.40	22.61±1.27
Toast	69.43±2.12	83.41±1.67	2.941±0.066	4.894±0.173	<u>26.86±0.80</u>	59.47±0.65	81.77±0.71	24.71±0.93
generative	66.95±1.93	82.28±1.24	2.900±0.090	4.724±0.138	26.91±0.92	64.52±0.18	92.09±0.11	30.14±0.74
contrastive	77.27±1.65	89.92±1.22	3.158±0.049	5.086±0.111	29.47±1.38	62.50±0.14	87.46±0.07	26.75±0.41
MMTEC	<u>76.57±3.05</u>	91.38±2.04	2.870±0.023	4.653±0.073	26.06±0.44	67.62±0.21	93.68±0.08	36.38±0.79

Bold denotes the best result, underline denotes the second best result.

destination prediction, which requires predicting destination road segments. Additionally, their trajectory encoders are based on RNNs, which are unable to effectively capture continuous spatio-temporal correlations.

In comparison, TremBR, CTLE, PT2Vec, and Toast incorporate semantic information by mapping trajectories onto road networks and using trajectory encoders that encode map-matched trajectories into embeddings. TremBR also considers temporal information by including the visiting time of road segments as an input feature. CTLE incorporates relative and absolute temporal information through a temporal encoding layer and a masked time pre-training task. PT2Vec applies a graph partition algorithm to road networks to build a more effective model on large-scale networks. Toast improves its incorporation of semantic information by pre-training embeddings for the road network using node2vec. However, their trajectory encoders still rely on RNNs or Transformers, which update their hidden states discretely, making it challenging to accurately model the continuous spatio-temporal correlations in trajectories.

In contrast, our method incorporates both the semantic information of travel and the continuous spatio-temporal correlations present in trajectories. This enables comprehensive information to be incorporated into the learned embeddings, thereby improving the performance of downstream tasks.

Performance comparison summary. The proposed method

leverages the benefits of maximum entropy coding and the trajectory multi-view approach to form the MMTEC pretext task, which enables the learning of trajectory embeddings that are highly generalizable across different downstream tasks. The proposed method captures both the travel semantic information and continuous spatio-temporal correlations in trajectories, providing comprehensive information that improves the performance of downstream tasks. Furthermore, the method’s discrete and continuous trajectory encoders effectively and efficiently embed trajectories. These advantages yield superior performance across multiple downstream tasks.

5.5.2 Efficiency Comparison

When investigating the efficiency of the trajectory embedding methods, we consider the model size, training speed, and embedding time of all embedding methods on the Chengdu dataset. The training speed is the time it takes to complete one epoch of training, and the embedding time is the time in seconds it takes to calculate embeddings for all trajectories in the test set. The results are presented in Table 6.

Among the methods that incorporate road networks, our method has the smallest model size. This is because we use maximum entropy coding and multi-view consistency to reduce redundancy in the learned trajectory embeddings.

The RNN-based methods, such as traj2vec, t2vec, GM-VSAE, TremBR, and PT2Vec, are relatively slower at train-

TABLE 6
Efficiency comparison between methods on the Chengdu dataset.

Property	Model size (MBytes)	Train speed (min/epoch)	Embed time (sec)
traj2vec	1.106	0.924	1.039
t2vec	1.106	0.851	<u>1.028</u>
GM-VSAE	3.029	0.931	1.482
TremBR	6.711	0.973	1.428
CTLE	3.218	0.795	1.132
PT2Vec	2.462	0.874	1.231
Toast	3.218	2.733	1.129
Transformer	1.909	<u>0.713</u>	1.033
MMTEC	1.761	0.371	0.510

Bold denotes the best result, underline denotes the second best result.

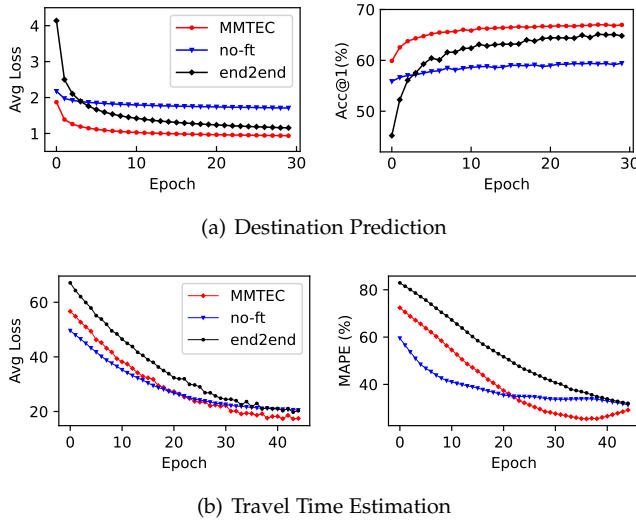


Fig. 6. Average loss value and validation metrics through the training processes of different variances.

ing and embedding, as RNNs do not parallelize on long trajectory sequences. In comparison, the Transformer-based embedding method CTLE is faster at training thanks to its parallelized attention mechanism. Toast is the slowest at training as the pre-training of road segment embeddings takes up additional time. Our method is the fastest at training and embedding due to its sharing of anchor sequences across mini-batches. To assess the efficiency of the proposed discrete trajectory encoder, we replace the induced attention layers within the encoder by vanilla Transformer encoder layers. The efficiency comparison of the modified encoder is shown in the *Transformer* row in Table 6. A comparison of the efficiency metrics with those of MMTEC indicates that a robust acceleration is achieved by the incorporation of the induced attention layers.

Overall, our method is efficient while outperforming the other trajectory embedding methods on downstream tasks.

5.5.3 Effectiveness of Pre-training

We evaluate the effectiveness of the pre-training by comparing the complete MMTEC method with two variants of MMTEC.

- 1) *no-ft*: not implementing fine-tuning. After pre-training, the parameters of the trajectory encoders

are fixed and do not participate in the backward propagation. Only the parameters of the downstream predictor are trained in downstream tasks.

- 2) *end2end*: train the model in end-to-end manner. The parameters of the trajectory encoders are initialized randomly and are trained directly with supervision from downstream tasks.

The average loss values and validation metrics for destination prediction and travel time estimation on the Xian dataset are presented in Figure 6. They show that the pre-training and fine-tuning approach facilitates accelerated convergence, signifying improved generalization and swifter progress in downstream tasks following the pre-training phase.

Pre-training without fine-tuning does not reach the performance of the fine-tuned and the end-to-end trained methods. Nevertheless, it still achieves satisfactory results, indicating that the pre-training embeddings capture general features of trajectories that are independent of task-specific supervision.

5.5.4 Effectiveness of the Pretext Task

To evaluate the effectiveness of the pretext task in our proposed method, we compare the trajectory embeddings learned by the complete MMTEC method with the embeddings learned when MMTEC uses the following pretext tasks for training:

- 1) *generative*: recover the original trajectory sequence given the embeddings. Here, we utilize two Transformer decoders to pair with the discrete and continuous trajectory encoders for trajectory recovery in an auto-regressive manner. The two encoders are trained separately.
- 2) *contrastive*: discriminating the positive sample from a set of negative samples given the target. For each trajectory in a mini-batch, we regard its continuous spatio-temporal embedding as the target, its semantic embedding as the positive sample, and the semantic embeddings of other trajectories in the batch as negative samples. InfoNCE [16] is chosen as the loss function.

Their performance on various downstream tasks is shown in Tables 4 and 5 and in Figure 1. We observe that when trained using the generative pretext task, the embeddings are biased towards point-level tasks such as travel time estimation and destination prediction. Conversely, when trained using the contrastive pretext task, the embeddings are biased towards sequence-level tasks such as similar trajectory search. Our full method, trained using the proposed MMTEC pretext task, possess consistent performance across different downstream tasks.

5.5.5 Impact of Hyper-parameters

We study the effects of the hyper-parameters listed in Table 3 on the test set of the Chengdu dataset. We use the Acc@1 and Acc@5 metrics of the similar trajectory search task, since the task does not involve fine-tuning and can indicate the quality of the learned trajectory embeddings directly. The

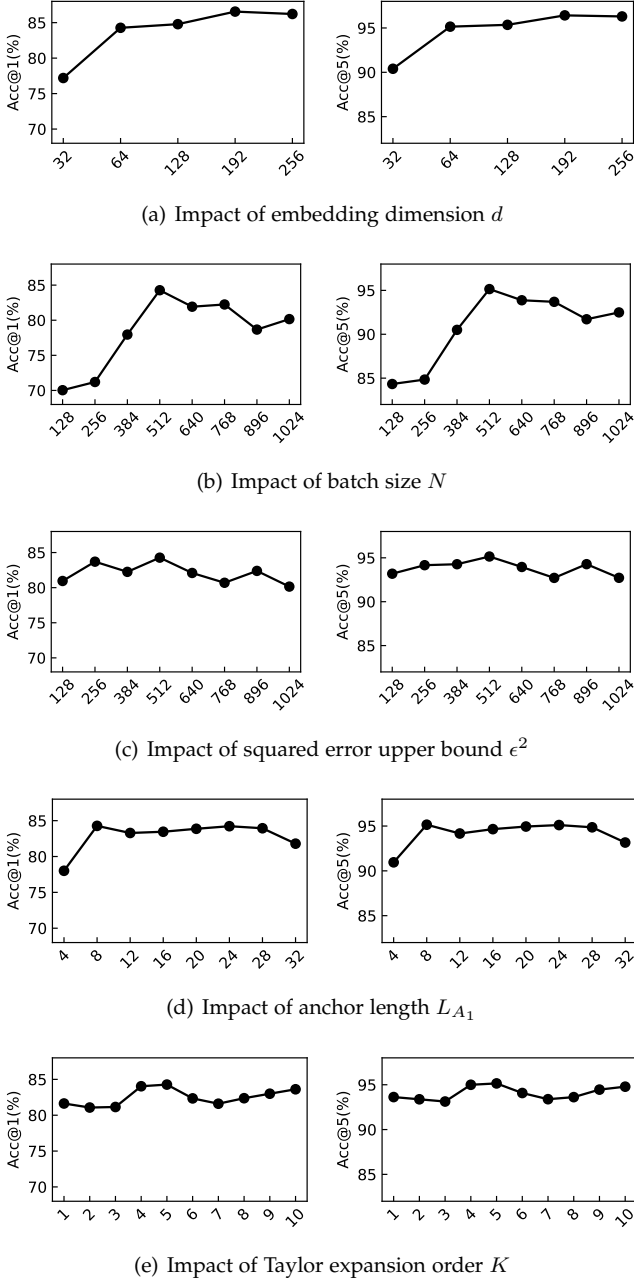


Fig. 7. Impact of hyper-parameters for the similar trajectory search task on the Chengdu dataset.

results are presented in Figure 7, and we make the following observations:

- 1) As illustrated in Figure 7(a), increasing the dimensionality of embeddings generally improves their performance. However, beyond $d = 64$, the accuracy improvement is negligible, while the computation and memory requirements increase significantly. Therefore, we set $d = 64$ to balance performance and efficiency.
- 2) Both the batch size N and the squared error upper bound ϵ^2 affect the pre-training loss function’s coefficients (Equation 19). As shown in Figure 7(b) and Figure 7(c), their optimal value is 512 in both cases.
- 3) The anchor length L_{A_1} of the first IA layer in the

TABLE 7
Performance comparison of different variances on the similar trajectory search task, Chengdu dataset.

Metric	Acc@1(%)	Acc@5(%)
only-sem	64.84±11.53	77.93±13.62
only-con	59.36±4.67	71.23±5.32
no-time	67.38±5.28	78.14±7.50
no-spatial	53.10±0.51	64.23±1.04
dis-no-time	70.72±0.49	82.53±0.91
dis-no-spatial	72.71±2.14	83.93±2.85
con-no-time	69.69±0.57	80.24±0.83
con-no-spatial	62.79±0.59	73.56±0.91
Dis-Trans	82.86±2.57	94.36±1.39
Con-Trans	74.97±3.55	88.75±1.91
MMTEC	84.27±3.02	95.15±1.63

Bold denotes the best result, underline denotes the second best result.

discrete trajectory encoder has an optimal value of 8, as seen in Figure 7(d). A smaller length decreases model complexity, making it challenging to extract correlations from trajectories, while a longer length increases the model capacity, leading to overfitting.

- 4) The Taylor expansion order K determines the accuracy of the estimation of Equation 18. We set $K = 5$ since it achieves good results.

5.5.6 Ablation Study

To evaluate the effectiveness of different components of the proposed method, we compare the performance of learned embeddings obtained by the complete version and the following variants of the method:

- 1) *only-sem*: only use the semantic embeddings of trajectories in downstream tasks.
- 2) *only-con*: only use the continuous spatio-temporal embeddings of trajectories in downstream tasks.
- 3) *no-time*: remove the timestamp feature from both trajectory encoders.
- 4) *no-spatial*: remove the spatial coordinate features from both trajectory encoders.
- 5) *dis-no-time*: remove the timestamp feature from the discrete trajectory encoder.
- 6) *dis-no-spatial*: remove the spatial coordinate features from the discrete trajectory encoder.
- 7) *con-no-time*: remove the timestamp feature from the continuous trajectory encoder.
- 8) *con-no-spatial*: remove the spatial coordinate features from the continuous trajectory encoder.
- 9) *Dis-Trans*: replace the discrete trajectory encoder with a Transformer encoder.
- 10) *Con-Trans*: replace the continuous trajectory encoder with a Transformer encoder that takes the trajectory sequence \mathcal{T} as input.

For the reason stated in Section 5.5.5, we compare the performance of these variants at the similar trajectory search task on the Chengdu dataset. The results, shown in Table 7, lead to the following observations:

- 1) Both the semantic and the continuous spatio-temporal embeddings are useful in downstream

tasks. Using only one of them at similar trajectory search fails to match the accuracy of the full method. This indicates that both types of embeddings capture important and complementary information about the trajectories.

- 2) The spatial and temporal features are essential for both trajectory encoders. Removing either of these significantly degrades the accuracy. Since the two encoders are trained together during pre-training, even removing features from one encoder impacts the performance of the learned embeddings substantially.
- 3) Replacing the proposed discrete trajectory encoder with a vanilla Transformer encoder degrades the accuracy somewhat. This suggests that the proposed encoder is capable of effectively and efficiently handling sequences of map-matched trajectories better than the Transformer. Additionally, we observe a decrease in performance when replacing the continuous trajectory encoder with a Transformer encoder, as Transformers do not inherently model spatio-temporal correlations continuously.

6 CONCLUSION

We propose MMTEC, a novel pre-training trajectory embedding method that learns general trajectory embeddings that generalize to a wide range of downstream trajectory mining tasks, improving their performance. We propose a discrete trajectory encoder and a continuous trajectory encoder to extract semantic travel information and continuous spatio-temporal information from trajectories, respectively, improving the comprehensiveness of the learned trajectory embeddings. The experimental study offers evidence of the improved effectiveness and performance over existing methods of the proposed method at pre-training high-quality trajectory embeddings.

ACKNOWLEDGMENTS

This work was supported by the Fundamental Research Funds for the Central Universities (Grant No. 2021YJS030).

REFERENCES

- [1] D. Kong and F. Wu, "HST-LSTM: A hierarchical spatial-temporal long-short term memory network for location prediction." in *International Joint Conference on Artificial Intelligence*, 2018, pp. 2341–2347.
- [2] X. Li, G. Cong, A. Sun, and Y. Cheng, "Learning travel time distributions with deep generative model," in *World Wide Web Conference*, 2019, pp. 1017–1027.
- [3] H. Hong, Y. Lin, X. Yang, Z. Li, K. Fu, Z. Wang, X. Qie, and J. Ye, "HetETA: heterogeneous information network embedding for estimating time of arrival," in *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 2444–2454.
- [4] G. Jin, H. Yan, F. Li, Y. Li, and J. Huang, "Hierarchical neural architecture search for travel time estimation," in *International Conference on Advances in Geographic Information Systems*, 2021, pp. 91–94.
- [5] X. Han, R. Cheng, C. Ma, and T. Grubenmann, "DeepTEA: effective and efficient online time-dependent trajectory outlier detection," *Proceedings of the VLDB Endowment*, vol. 15, no. 7, pp. 1493–1505, 2022.
- [6] D. Yao, C. Zhang, Z. Zhu, J. Huang, and J. Bi, "Trajectory clustering via deep representation learning," in *International Joint Conference on Neural Networks*, 2017, pp. 3880–3887.
- [7] Y. Liu, K. Zhao, G. Cong, and Z. Bao, "Online anomalous trajectory detection with deep generative sequence modeling," in *International Conference on Data Engineering*, 2020, pp. 949–960.
- [8] P. Zhao, H. Zhu, Y. Liu, J. Xu, Z. Li, F. Zhuang, V. S. Sheng, and X. Zhou, "Where to go next: A spatio-temporal gated network for next POI recommendation," in *AAAI Conference on Artificial Intelligence*, 2019, pp. 5877–5884.
- [9] F. Sun, R. Gao, W. Xing, Y. Zhang, W. Lu, J. Fang, and S. Liu, "Deep fusion for travel time estimation based on road network topology," *IEEE Intell. Syst.*, vol. 37, no. 3, pp. 98–107, 2022.
- [10] W. Dang, H. Wang, S. Pan, P. Zhang, C. Zhou, X. Chen, and J. Wang, "Predicting human mobility via graph convolutional dual-attentive networks," in *ACM International Conference on Web Search and Data Mining*, 2022, pp. 192–200.
- [11] A. Ashukha, A. Lyzhov, D. Molchanov, and D. P. Vetrov, "Pitfalls of in-domain uncertainty estimation and ensembling in deep learning," in *International Conference on Learning Representations*, 2020.
- [12] T. Ishida, I. Yamane, T. Sakai, G. Niu, and M. Sugiyama, "Do we need zero training loss after achieving zero training error?" in *International Conference on Machine Learning*, 2020, pp. 4604–4614.
- [13] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 4171–4186.
- [14] Z. Yang, Z. Dai, Y. Yang, J. G. Carbonell, R. Salakhutdinov, and Q. V. Le, "XLNet: Generalized autoregressive pretraining for language understanding," in *Advances in Neural Information Processing Systems*, 2019, pp. 5754–5764.
- [15] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," in *European Conference on Computer Vision*, 2020, pp. 776–794.
- [16] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [17] P. Bachman, R. D. Hjelm, and W. Buchwalter, "Learning representations by maximizing mutual information across views," *Advances in neural information processing systems*, 2019.
- [18] X. Qiu, T. Sun, Y. Xu, Y. Shao, N. Dai, and X. Huang, "Pre-trained models for natural language processing: A survey," *CoRR*, vol. abs/2003.08271, 2020.
- [19] X. Li, K. Zhao, G. Cong, C. S. Jensen, and W. Wei, "Deep representation learning for trajectory similarity computation," in *International Conference on Data Engineering*, 2018, pp. 617–628.
- [20] T.-Y. Fu and W.-C. Lee, "TremBR: Exploring road networks for trajectory representation learning," *ACM Transactions on Intelligent Systems and Technology*, vol. 11, no. 1, pp. 1–25, 2020.
- [21] J. Li, M. Wang, L. Li, K. Xin, W. Hua, and X. Zhou, "Trajectory representation learning based on road network partition for similarity computation," in *International Conference on Database Systems for Advanced Applications*, 2023, pp. 396–413.
- [22] Y. Chen, X. Li, G. Cong, Z. Bao, C. Long, Y. Liu, A. K. Chandran, and R. Ellison, "Robust road network representation learning: When traffic patterns meet traveling semantics," in *ACM International Conference on Information & Knowledge Management*, 2021, pp. 211–220.
- [23] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [25] T. Linzen, E. Dupoux, and Y. Goldberg, "Assessing the ability of lstms to learn syntax-sensitive dependencies," *Transactions of the Association for Computational Linguistics*, vol. 4, pp. 521–535, 2016.
- [26] Y. Liang, K. Ouyang, H. Yan, Y. Wang, Z. Tong, and R. Zimmermann, "Modeling trajectories with neural ordinary differential equations," in *International Joint Conference on Artificial Intelligence*, 2021, pp. 1498–1504.
- [27] Y. Zhou and Y. Huang, "Deepmove: Learning place representations through large scale movement data," in *IEEE International Conference on Big Data*, 2018, pp. 2403–2412.

- [28] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *International Conference on Learning Representation*, 2013.
- [29] S. Feng, G. Cong, B. An, and Y. M. Chee, "POI2Vec: Geographical latent representation for predicting future visitors." in *AAAI Conference on Artificial Intelligence*, 2017, pp. 102–108.
- [30] H. Wan, Y. Lin, S. Guo, and Y. Lin, "Pre-training time-aware location embeddings from spatial-temporal trajectories," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 11, pp. 5510–5523, 2021.
- [31] Y. Lin, H. Wan, S. Guo, and Y. Lin, "Pre-training context and time aware location embeddings from spatial-temporal trajectories for user next location prediction," in *AAAI Conference on Artificial Intelligence*, 2021, pp. 4241–4248.
- [32] S. Feng, X. Li, Y. Zeng, G. Cong, Y. M. Chee, and Q. Yuan, "Personalized ranking metric embedding for next new POI recommendation," in *International Joint Conference on Artificial Intelligence*, 2015, pp. 2069–2075.
- [33] C. Zhou, J. Bai, J. Song, X. Liu, Z. Zhao, X. Chen, and J. Gao, "ATRank: An attention-based user behavior modeling framework for recommendation," in *AAAI Conference on Artificial Intelligence*, 2018, pp. 4564–4571.
- [34] Z. Yao, Y. Fu, B. Liu, W. Hu, and H. Xiong, "Representing urban functions through zone embedding with human mobility patterns." in *International Joint Conference on Artificial Intelligence*, 2018, pp. 3919–3925.
- [35] T. Shimizu, T. Yabe, and K. Tsubouchi, "Learning fine grained place embeddings with spatial hierarchy from human mobility trajectories," *arXiv preprint arXiv:2002.02058*, 2020.
- [36] M.-x. Wang, W.-C. Lee, T.-y. Fu, and G. Yu, "Learning embeddings of intersections on road networks," in *ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 2019, pp. 309–318.
- [37] N. Wu, X. W. Zhao, J. Wang, and D. Pan, "Learning effective road network representation with hierarchical graph neural networks," in *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 6–14.
- [38] Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, and J. Tang, "GLM: General language model pretraining with autoregressive blank infilling," in *Annual Meeting of the Association for Computational Linguistics*, 2022, pp. 320–335.
- [39] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [40] A. Pauls and D. Klein, "Faster and smaller n-gram language models," in *Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 2011, pp. 258–267.
- [41] C. Yang and G. Gidofalvi, "Fast map matching, an algorithm integrating hidden markov model with precomputation," *International Journal of Geographical Information Science*, vol. 32, no. 3, pp. 547–570, 2018.
- [42] Principle of maximum entropy. [Online]. Available: https://en.wikipedia.org/wiki/Principle_of_maximum_entropy
- [43] X. Liu, Z. Wang, Y.-L. Li, and S. Wang, "Self-supervised learning via maximum entropy coding," in *Advances in Neural Information Processing Systems*, 2022.
- [44] D. Xu, C. Ruan, E. Körpeoglu, S. Kumar, and K. Achan, "Inductive representation learning on temporal graphs," in *International Conference on Learning Representations*, 2020.
- [45] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng, "Fourier features let networks learn high frequency functions in low dimensional domains," in *Advances in Neural Information Processing Systems*, 2020.
- [46] Y. Li, S. Si, G. Li, C. Hsieh, and S. Bengio, "Learnable fourier features for multi-dimensional spatial positional encoding," in *Advances in Neural Information Processing Systems*, 2021, pp. 15 816–15 829.
- [47] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh, "Set transformer: A framework for attention-based permutation-invariant neural networks," in *International Conference on Machine Learning*, 2019, pp. 3744–3753.
- [48] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [49] P. Kidger, J. Morrill, J. Foster, and T. Lyons, "Neural controlled differential equations for irregular time series," *Advances in Neural Information Processing Systems*, pp. 6696–6707, 2020.
- [50] R. T. Q. Chen. (2018) torchdiffeq: PyTorch implementation of differentiable ODE solvers. [Online]. Available: <https://github.com/rtqichen/torchdiffeq>
- [51] Openstreetmap. [Online]. Available: <https://www.openstreetmap.org/>
- [52] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems*, 2019, pp. 8024–8035.
- [53] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.



Yan Lin received the B.S. degree in computer science from Beijing Jiaotong University, Beijing, China, in 2019.

He is currently working toward the Ph.D. degree in the School of Computer and Information Technology, Beijing Jiaotong University. His research interests include spatio-temporal data mining and graph neural networks.



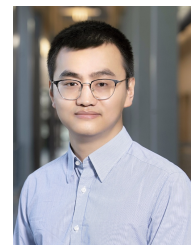
Huaiyu Wan received the Ph.D. degree in computer science and technology from Beijing Jiaotong University, Beijing, China, in 2012.

He is a Professor with the School of Computer and Information Technology, Beijing Jiaotong University. His current research interests focus on spatio-temporal data mining, social network mining, information extraction, and knowledge graph.



Shengnan Guo received the B.S. degree in computer science from Beijing Jiaotong University, Beijing, China, in 2015.

She is currently working toward the Ph.D. degree in the School of Computer and Information Technology, Beijing Jiaotong University. Her research interests focus on the area of deep learning and spatio-temporal data mining.



Jilin Hu received the Ph.D. degree in computer science from Aalborg University, Aalborg, Denmark in 2019.

He is an Associate Professor at the Department of Computer Science, Aalborg University. His research interests include spatio-temporal data analytics and transportation data mining.



Christian S. Jensen received the Ph.D. degree from Aalborg University in 1991 after 2 1/2 years of study at University of Maryland, and he received the Dr.Techn. degree from Aalborg University in 2000.

He is a Professor at the Department of Computer Science, Aalborg University. His research concerns primarily temporal and spatio-temporal data management and analytics, including indexing and query processing, data mining, and machine learning.



Youfang Lin received the Ph.D. degree in signal and information processing from Beijing Jiaotong University, Beijing, China, in 2003.

He is a Professor with the School of Computer and Information Technology, Beijing Jiaotong University. His main fields of expertise and current research interests include big data technology, intelligent systems, complex networks, and traffic data mining.