

# Predicting Human Mobility via Self-supervised Disentanglement Learning

Qiang Gao, Jinyu Hong, Xovee Xu, Ping Kuang, Fan Zhou, and Goce Trajcevski

**Abstract**—Deep neural networks have recently achieved considerable improvements in learning human behavioral patterns and individual preferences from massive spatial-temporal trajectories data. However, most of the existing research concentrates on fusing different semantics underlying sequential trajectories for mobility pattern learning which, in turn, yields a narrow perspective on comprehending human intrinsic motions. In addition, the inherent sparsity and under-explored heterogeneous collaborative items pertaining to human check-ins hinder the potential exploitation of human diverse periodic regularities as well as common interests. Motivated by recent advances in disentanglement learning, in this study we propose a novel disentangled solution called SSDL for tackling the next POI prediction problem. SSDL primarily seeks to disentangle the potential time-invariant and time-varying factors into different latent spaces from massive trajectories data, providing an interpretable view to understand the intricate semantics underlying human diverse mobility representations. To address the data sparsity issue, we present two realistic trajectory augmentation approaches to enhance the understanding of both the human intrinsic periodicity and constantly-changing intents. In addition, we devise a POI-centric graph structure to explore heterogeneous collaborative signals underlying historical check-ins. Extensive experiments conducted on four real-world datasets demonstrate that our proposed SSDL significantly outperforms the state-of-the-art approaches – for example, it yields up to 8.57% improvements on ACC@1.

**Index Terms**—location-based services, human mobility, graph neural network, disentanglement learning, variational Bayes.



## 1 INTRODUCTION

The proliferation of geo-tagged social media (GTSM) such as Foursquare and WeChat have enabled numerous users to post interesting places, report daily activities, and make like-minded friends, resulting in the accumulation of massive amounts of contextual data (e.g., check-ins). This, in turn, offers unprecedented opportunities to explore human diverse life experiences (e.g., mobility patterns) and facilitate the development of various user-centric downstream applications such as trajectory identification [1], POI recommendation/prediction [2], itinerary prediction [3] – to name a few. As a fundamental task in mining check-in data, predicting *human mobility* (often exemplified as next POI prediction/recommendation) is critical for researchers and practitioners to explore the informative semantics and mutual interactions behind human check-ins [4], [5]. For instance, it enables one to precisely ascertain users' future intentions and draw in more potential customers for new ventures [6], [7].

Spatio-temporal check-in sequences (i.e., trajectories) reflect human daily activities upon a set of POIs, which may include certain (periodic) regularities. The majority of the pioneering works in human mobility prediction aimed at

modeling human sequential behaviors taking into account spatio-temporal preferences. For instance, in order to predict where a certain user will go in the near future, conventional approaches such as Markov Chain [8] and Tensor-based Factorization [9] that rely on data-driven paradigms, attempt to incorporate individual visiting preferences and explore *sequential* patterns. However, these approaches depend heavily on hand-crafted characteristics and face the challenge of comprehending the diverse semantics underlying massive volumes of human trajectories. This, in turn, leads to narrow solutions in disclosing human implicit interactive hints/signals regarding historical check-ins.

More recent deep learning techniques such as recurrent neural networks (RNNs) have brought about encouraging achievements of learning informative check-ins (including POIs) from human trajectories and become a widespread and popular methodology in tackling miscellaneous mobility learning tasks [4], [10], [11]. For example, Wu et al. [10] present a PLSPL model, which leverages a Long-Short Term Memory (LSTM) neural network to model human short-term sequential preferences while learning contextual features of POIs behind human historical check-ins via attention mechanism. To consider the spatial and temporal influences for next POI recommendation, Kong et al. [12] incorporate the spatial and temporal intervals between two successive check-ins into recurrent hidden states to mitigate the data sparsity of human trajectories. Due to the higher model efficiency and the ability to quantify the contribution of each check-in in a given trajectory, several attention mechanisms like, for example, self-attention and vanilla attention emerged for handling long human historical trajectories [2], [6], [13].

Several state-of-the-art methods have employed graph

*Fan Zhou is the corresponding author.*

- Qiang Gao is with the School of Computing and Artificial Intelligence, Southwestern University of Finance and Economics, Chengdu, China. E-mail: qianggao@swufe.edu.cn
- Jinyu Hong, Xovee Xu, Ping Kuang, and Fan Zhou are with the University of Electronic Science and Technology, Chengdu, China. E-mail: {jinyuhong@std., xovee@std., kuangping@, fan.zhou@}uestc.edu.cn
- Goce Trajcevski is with the Iowa State University, Iowa, USA. E-mail: gocet25@iastate.edu

*Manuscript received XX XX, 2022; revised XX XX, 2022.*

structure learning to explicitly uncover spatial correlations or collaborative signals to understand individual human interests. More concretely, they attempt to acquire expressive POI representations by considering the rich contexts of highly correlated POIs. For instance, conventional methods such as word2vec-based [14], [15] and deepwalk-based [16], [17] have successfully uncovered the higher-order correlations between consecutive check-ins and offered contextual POI representations. Other schemes using popular graph neural networks (GNNs), such as graph convolutional networks [11] and graph attention networks [18], primarily seek to incorporate the POI-to-POI correlations (e.g., geographical proximity) behind massive human trajectories.

Despite the recent achievements in deep human mobility learning, we observe that existing solutions still have three significant drawbacks:

(a) *Implicit semantic entanglement.* Although there is a large body of work on human mobility representation learning, the most common scheme is to take a past check-in sequence as input, and either that sequence or the user's next POI is used as the supervision signal. The former can be framed as self-supervised learning, while the latter is standard supervised learning. Nevertheless, both ultimately focus on fusing multiple semantics behind sequential trajectories to predict the user's next POI, which could lead to a myopic perspective and produce a non-diverse recommendation result. We call this phenomenon *semantic entanglement*. In practice, human trajectories, as typical sequential data, contain rich user mobility patterns that reflect diverse periodic regularities or behavioral habits of humans. More importantly, the intrinsic individual patterns/habits of humans are difficult to change over time, but their near-term intentions/behaviors are prone to be influenced/dictated by certain time instants. Thus, we consider that human mobility patterns can be implicitly disentangled into two aspects: *time-independent* and *time-dependent* behaviors. Existing solutions rely on data-driven models to understand limited mobility patterns, which fail to reveal the nature of human visiting intents. As a result, they only provide a narrow scope to become familiar with human future behaviors, which usually carries the risk of prediction bias due to the limited scale of trajectory data available.

(b) *Sparsity in representation learning.* Only when a user decides or is willing to check in via location-based applications can a POI be recorded, which inevitably leads to the sparsity problem when gathering historical human footprints. As a result, the sparsity problem hinders the model from learning a good representation of human mobility. Existing methods either use the next POI as the sole supervisor or complement the representation learning in a semi-supervised manner with unlabeled trajectories. These paradigms mostly follow the merit of text representation in the field of natural language processing (NLP) which, however, easily fails in capturing the innate rules underlying human trajectories such as individual periodic regularity.

(c) *Heterogeneous collaborative signals.* Most existing efforts concentrate on learning POI-to-POI relationships (a.k.a. the connectivity of POIs) from a large number of trajectories, such as consecutive correlation and geographic proximity [17], [18], [19]. Despite the successful collaboration of individual human interests via these homogeneous graph

structure learning, a notable limitation is that heterogeneous semantics affiliated with the POIs are not investigated well, yielding a limit in the exploration of affluent common preferences behind human diverse trajectories. For example, people may have similar visit time preferences for certain POIs, such as going to a café after lunch. In addition, each POI is associated with a textual description (e.g., POI category), reflecting the underlying human activity interest. We conjecture that incorporating the heterogeneous correlations between POI and their category can provide us with a coarse-grained view of the higher-order connective between POIs. For example, people often go to several fashion stores to buy clothes at a time.

To address the aforementioned limitations, we present a novel solution called **SSDL**, a self-supervised disentanglement learning framework for understanding human mobility. Rather than previous data-driven representation learning, SSDL performing self-supervision in the latent space aims at seeking a clean separation of the time-independent and time-dependent vectors for diverse human trajectories, which is inspired by the recent advances of variational inference and contrastive learning. Specifically, SSDL operates the sequential variational autoencoder (VAE) with a mutual information regularization to guide the training of evidence lower bound (ELBO), aiming at promoting the disentanglement of human mobility-related representations. In particular, we provide two realistic trajectory augmentation strategies to alleviate the sparsity issue in representation learning, which can further help us enhance the understanding of human intrinsic periodicity and constantly-changing intents. In addition, we also present a POI-centric graph structure to explore human common interests underlying diverse check-ins, which primarily seeks human consecutive, geospatial, temporal-aspect, and activity-aspect interests. In sum, our contributions can be summarized as follows:

- We introduce a novel disentangled representation learning framework to understand human time-independent and time-dependent behaviors of their individual mobility patterns. To the best of our knowledge, this study is the first work to disentangle human mobility and investigate how it can be used for the prediction of the next POI.
- We propose two practical trajectory augmentation methods, guided by the inherent characteristics of individual human mobility patterns, to promote disentanglement learning.
- To capture heterogeneous collaborative signals behind historical check-ins, We devise a flexible POI-centric network structure to explore rich human interests in trajectories, which enhances the performance of downstream next POI prediction task.
- We conduct extensive experiments on four real-world datasets to evaluate the performance of our proposed SSDL. The results demonstrate that our approach outperforms state-of-the-art methods.

## 2 RELATED WORK

### 2.1 Next POI Prediction in Deep Learning

Recent deep learning solutions have stimulated many researchers and practitioners to learn human periodic regu-

larities from massive historical check-ins. Especially, deep (recurrent) neural networks such as LSTM [20] and GRU [21] have received widespread interest in the next POI prediction task as they are able to capture the sequential dependencies for mobility pattern understanding. For instance, [22] extends the vanilla RNN model and integrates the spatial-temporal impacts into each RNN cell, yielding promising results on the next location prediction. Zhao et al. propose a novel ST-LSTM that implements time gates and distance gates into standard LSTM, aiming at capturing the spatio-temporal relation between consecutive check-ins [4]. To learn more contextual information, Wu et al. [10] propose a personalized long- and short-term preference learning scheme to learn the specific user context, where the different influences of locations and categories of POIs are considered. While most of endeavors focus on pruning or modifying the RNN-based modules [23], [24], [25], researches also tried to adopt other popular deep neural networks for next POI prediction, e.g., attention-based neural networks [2], [26] and convolutional neural networks [6], [27]. Xue et al. [2] build the Transformer architecture as the mobility feature extractor in which it regards the historical trajectory and semantic contexts as the input to handle multiple factors such as temporal and geographic contexts.

## 2.2 Mobility Representation Learning

POI embedding and trajectory embedding, as two core components in mobility representation learning, have been investigated in recent studies.

For POI embeddings, the earlier studies such as [22] and [28] set a fixed or learnable matrix as the initial representations of POIs, primarily seeking to alleviate the “Curse of Dimensionality” concern. However, any semantic information between POIs is under-explored. As word embeddings, especially word2vec-based [29], have achieved great performance in NLP, recent studies also proposed various word2vec-based solutions aimed at capturing the proximity semantics of POIs from human check-in sequences (or real-world trajectories). For instance, [30] and [31] regard each POI as a “word” while each human trajectory as a “sequence”, and use word2vec to obtain a low-dimensional vector for each POI. POI2Vec is a latent representation model that incorporates geographic influence when using word2vec method for POI embedding [14]. However, training sparse trajectories to obtain POI representations often confronts the problem of poor capability of POI semantics. More recently, the extraordinary success of graph neural networks (GNNs) has inspired tremendous researchers to turn to devise graph-based models to facilitate the learning of human trajectories [11], [18]. For instance, [18] proposes a graph-based model to explore the spatial, temporal, and preference factors behind the POIs. However, it only considers homogeneous interactions among the POIs and ignores heterogeneous interactions with other key entities such as activity and check-in time.

Regarding trajectory representation learning, the majority of existing research concentrates on taking the historical trajectory as input and using the next POI as the sole supervision signal [4], [7], [22], [32]. To address the narrow scale of trajectory data, some efforts attempt to employ the

unlabeled trajectories as supplements and train them with the labeled trajectories jointly in an unsupervised or self-supervised manner to acquire a good representation for each trajectory [6], [33], [34]. Especially, to operate the trajectories in a latent space, recent studies employ generative models such as variational inference or adversarial models to learn the intrinsic distribution underlying massive trajectory data and then turn to fine-tune the model for the next POI prediction tasks. For instance, VANext extended the variational autoencoder (VAE) to consider the uncertainty of user preferences for regularized representation of historical trajectories [6]. A meta-learning technique called METAODE also employed variational Bayes to encode past human movement patterns into latent space [35]. In essence, these approaches principally rely on integrating numerous semantics including sequential information into a unified space while omitting the possibility of disentangling it to expose the characteristics of human mobility patterns.

## 2.3 Disentanglement Learning

The privilege of disentanglement learning is that it enables an interpretable perspective to understand the multiple inherent motions/factors behind the intricate data representations in addition to notable expressiveness. To disentangle the learned representations, most recent studies developed VAEs such as  $\beta$ -VAE to optimize the mutual interaction between different latent factors [36], [37], [38]. For example,  $\beta$ -VAE [36] is a simple but effective variant of the ordinary VAE that severely penalizes the Kullback–Leibler (KL) divergence term for disentanglement learning. Li et al. presented a Disentangled Sequential Autoencoder (DSVAE) approach for sequential data (e.g., video), aiming at factorizing the latent variables into static and dynamic parts [39]. To make the latent variables interpretable and controllable, a latent variable guidance-based generative model called Guided-VAE makes an effort to utilize VAE to learn a transparent representation [40]. Bai et al. presents a sequential VAE to learn disentangled representations in a self-supervised manner [41]. Bai et al. also extend the sequential VAE with a self-supervised learning approach to facilitate the factorization of video representations [38]. In addition, The newly developed self-supervised learning offers a new avenue to drive the acquisition of semantic representations [42]. For example, Ma et al. employ the ideas of latent self-supervision and intention disentanglement to boost the convergence of representation learning and utilize it in sequential recommendation tasks [43]. In sum, the success of these approaches suggests that, in addition to facilitating the understanding of rich semantics underlying data, disentangling the representation into distinct parts can make the representation more transparent and interpretable.

## 3 PRELIMINARIES

### 3.1 Problem Definition

**Definition 1 (POI).** Let  $l \in \mathcal{L}$  denotes a POI tagged by the location-based systems, and each POI corresponds to a geographic coordinate (e.g., longitude  $l_o$  and latitude  $l_a$ ) and a category  $ca$  (e.g., restaurant, museum, or park).

**Definition 2 (Check-in Sequence).** A check-in sequence (or trajectory)  $T_u = \{l_1^u, l_2^u, \dots, l_n^u\}$  left by user  $u$  is a

sequence of  $n$  POIs ordered by visiting time, where  $l_\tau^u$  means a user  $u$  visit POI  $l$  at time  $t_\tau$  ( $\tau \in \{1, 2, \dots, n\}$ ). Let  $\mathcal{T}_u = \{T_u^1, T_u^2, \dots, T_u^m\}$  denote  $m$  historical trajectories of user  $u$ , where each trajectory  $T_u^i$  contains a sequence of POIs ordered by visiting time, e.g.,  $T_u^i = \{l_1^{i,u}, l_2^{i,u}, \dots, l_n^{i,u}\}$ .

Formally, given a user  $u$  with his/her recently visited check-in sequence  $T_u^m = \{l_1^{m,u}, l_2^{m,u}, \dots, l_n^{m,u}\}$  and entire historical trajectory  $\mathcal{T}_u$ , our goal is to predict a POI  $l_{n+1}^{m,u}$  for user  $u$  to visit next. Notably, we mainly target disentangled representation learning for users' recently visited POI sequences. For simplicity, we will omit user identity (i.e.  $u$ ) and trajectory index (i.e.  $m$ ) in the following sections.

### 3.2 Variational Bayes

Variational Autoencoder (VAE) [44] containing an encoder and a decoder operates the input data  $\mathbf{x}$  into a latent space, where the latent variables are denoted by  $\mathbf{z}$ . Thus, the marginal likelihood  $\log p(\mathbf{x})$  can be obtained by maximizing the Evidence Lower BOund (ELBO), which is defined as:

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - KL[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]. \quad (1)$$

Herein,  $q_\phi(\mathbf{z}|\mathbf{x})$  is an approximate posterior distribution, parameterized by  $\phi$ ,  $p_\theta(\mathbf{x}|\mathbf{z})$  with parameters  $\theta$  is a likelihood function, and  $p(\mathbf{z})$  is a prior (e.g., Gaussian prior) over the latent variables.

### 3.3 Contrastive Estimation

In recent self-supervised learning paradigms [33], [45], mutual information (MI) is a common measure of the mutual dependence or compatibility between two variables. Specifically, they usually employ the noise contrastive estimation (NCE) [46], [47] to maximize the lower bound on the mutual information, which can be denoted as follows:

$$\mathcal{L}_{NCE} = \mathbb{E} \left[ -\log \left( \frac{\exp^{g(x)^\top g(x^+)}}{\exp^{g(x)^\top g(x^+)} + \sum_{j=1}^J \exp^{g(x)^\top g(x^-)}} \right) \right], \quad (2)$$

where  $x$ ,  $x^+$ , and  $x^-$  respectively denote the *anchor*, *positive*, and *negative* instances. Besides,  $\exp^{g(\cdot)^\top g(\cdot)}$  is a similarity measure (e.g., cosine similarity) between two instances.

## 4 ARCHITECTURE DESIGN

We make an overview of our proposed framework SSDL in Fig. 1, which mainly comprises three components. First, we build a POI-centric Graph (PGraph) to explore the common interests from the entire user trajectories and make interest aggregation to obtain both homogeneous and heterogeneous semantics underlying each POI. Then, our Self-supervised Disentanglement Learning component attempts to produce the time-invariant and time-varying variables for each trajectory. At last, SSDL uses the disentangled representations as well as the user's long-term preference modeled by an attentive network to predict the next POI.

### 4.1 Common Interest Distillation

To distill multiple correlations behind the POIs and their affiliated context, we build a POI-centric graph (PGraph).

#### 4.1.1 Graph Structure and Building Process

Incorporating prior correlations and multiple common interests are critical to obtain a good POI representation and understand human diverse mobility patterns. As several elements are recorded by LBSN, e.g., POI identity, geographical coordinate, visiting time, and POI category, we concentrate on exploring four contextual semantics to build our PGraph, including *consecutive*, *geospatial*, *time-aspect*, and *activity-aspect* interests.

Let  $\mathcal{G} = (\mathcal{V}, E)$  denotes our PGraph that models the human common interests, where  $\mathcal{V} = (\mathcal{V}_l \cup \mathcal{V}_t \cup \mathcal{V}_a)$  is the set of nodes, and  $E = (E_c \cup E_g \cup E_t \cup E_a)$  is the set of edges. Here  $\mathcal{V}_l = \mathcal{L}$  represents a collection of different POIs,  $\mathcal{V}_t$  is the set of time bins,  $\mathcal{V}_a$  denotes the set of POI categories, and  $E_c, E_g, E_t, E_a$  indicate the above four contextual semantics, respectively. That is to say,  $\mathcal{G}$  contains four sub-graphs, each of which represents an important user interest. Four sub-graphs are described in the following four paragraphs.

**Consecutive Interest.** According to [31], among millions of POIs in location-based systems, (1) people typically visit only a small subset of POIs that appeal to them; and (2) some POIs are visited more frequently than others. This phenomenon demonstrates that human mobility contains some common transitional regularities behind their past check-ins. Therefore, we consider that it is necessary to capture the consecutive correlations between distinct POIs to reveal human motion-based interests. Correspondingly, we formulate a weighted sub-graph  $\mathcal{G}_c = (\mathcal{V}_l, E_c, \mathbf{A}_c)$  to describe such diverse correlations, where  $\mathcal{V}_l$  is the set of distinct POIs,  $E_c$  is the edge set, and  $\mathbf{A}_c \in \mathbb{R}^{|\mathcal{L}| \times |\mathcal{L}|}$  refers to the adjacency matrix. Given two POIs (e.g., POI  $l_i$  and POI  $l_j$ ) that are successively visited, we create an edge between them and then calculate the edge weight (i.e., entry  $A_c^{ij} \in \mathbf{A}_c$ ) using the corresponding transitional probability. Formally, such an edge weight can be defined as:

$$A_c^{ij} = f_c^{ij} / f_c^i, \quad (3)$$

where  $f_c^{ij}$  refers to the frequency of edge  $l_i \rightarrow l_j$  appeared in the check-in data, and  $f_c^i$  denotes the frequency of POI  $l_i$  appeared in the check-in data. As such, we are able to acquire the matrix  $\mathbf{A}_c$  to preserve the consecutive interests underlying the trajectories.

**Geographical Interest.** People are more likely to visit nearby POIs than distant ones [17]. Motivated by this, we formulate an undirected sub-graph  $\mathcal{G}_g = (\mathcal{V}_l, E_g, \mathbf{A}_g)$  to describe such interactions, where  $E_g$  is the set of edges and  $\mathbf{A}_g \in \mathbb{R}^{|\mathcal{L}| \times |\mathcal{L}|}$  denotes the adjacency matrix regarding geographical interests. Given POI  $l_i$  and  $l_j$ , the edge weight  $A_g^{ij} \in \mathbf{A}_g$  can be calculated as:

$$A_g^{ij} = \begin{cases} 0, & g(l_i, l_j) > \Delta g; \\ 1, & \text{otherwise.} \end{cases} \quad (4)$$

Herein,  $g(l_i, l_j)$  is the great-circle distance function,  $\Delta g$  is a predefined threshold to restrict the impact of geographical noise. In this paper, we set  $\Delta g = 3 \text{ km}$ .

**Time-aspect Interest.** For each check-in, it is associated with a visiting timestamp, reflecting the human temporal semantics. As it is a key factor for understanding human periodic regularity, we propose to investigate the mutual interactions between POI and visiting time to obtain the

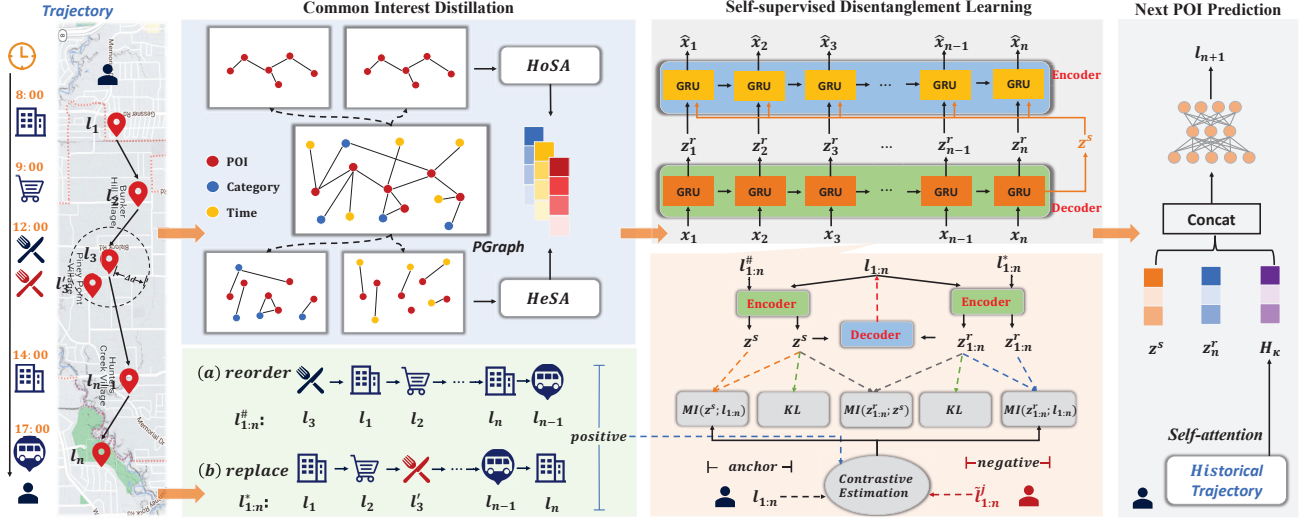


Fig. 1. The pipeline of proposed SSDL.

time-aspect interest. However, each visiting timestamp is actually a continuous value, we follow previous studies and aggregate all of the visiting timestamps into the hour-level time bins [48], [49]. Meanwhile, people may respectively show different preferences on weekday and weekend, we thus assign 48 time bins to replace the original visiting timestamps, where the weekday and weekend are specified. We thus formulate a weighted sub-graph  $\mathcal{G}_t(\mathcal{V}_t \cup \mathcal{V}_t, E_t, \mathbf{A}_t)$ , where  $\mathbf{A}_t$  maintains human time-aspect interest. Similar to the above graph  $\mathcal{G}_c$ , we can also calculate the time-aspect interest between the POI  $l_i$  and time bin  $t_\tau$  by:

$$A_t^{i\tau} = f_t^{i\tau} / f_t^i, \quad (5)$$

where  $f_t^{i\tau}$  denotes the frequency of visiting POI  $l_i$  at time  $t_\tau$ , and  $f_t^i$  is the total number that POI  $l_i$  has been visited.

**Activity-aspect Interest.** A user who wants to post a check-in to LBSNs indicates that he/she is engaged in a specific type of activity that appeals to him/her. In practice, each POI has a contextual description (i.e., POI category) that reflects a real-world activity, we consider that taking into account such contextual interactions is an essential addition to understanding human preferences. Notably, the number of POI categories is much smaller than the number of POIs. As a result, linking a POI to its category can offer a coarse-grained perspective on the higher-order interactions between various POIs. To this end, we build an undirected graph  $\mathcal{G}_a(\mathcal{V}_l \cup \mathcal{V}_a, E_a, \mathbf{A}_a)$  to describe the activity-aspect interest. To be more precise, we explicitly build an edge between a POI and the contextual category it belongs to, and then we treat each category as a regular node in  $\mathcal{G}_a$ .

#### 4.1.2 Interest Aggregation

To extract the semantic contexts underlying POIs from the PGraph, we propose to adopt graph neural networks (GNNs) which have been widely applied in numerous graph-based tasks and obtained remarkable success.

**Homogeneous Semantic Aggregation (HoSA).** According to the structure of the built PGraph, we can find that consecutive interest and geospatial interest that belong to the homogeneous semantics as they only contain the nodes

of POI identities. Thus, HoSA attempts to aggregate the underlying interest from the nodes of the same type, i.e., POI identity. First, the consecutive correlation matrix  $\mathbf{A}_c$  reflects human real-world transitional preferences, we can naturally regard each POI's transitional distribution as its prior feature to describe the relationship between a specific POI and its neighbors. To this end, we set each  $\mathbf{A}_c^i$  as the initial feature of the POI node  $\mathcal{V}_p^i$ . Besides, the geospatial correlation matrix  $\mathbf{A}_g$  preserves the geographical closeness between different POIs, providing the weak signal of human potential transitional tendencies. Hence,  $\mathbf{A}_g$  can be regarded as an augmentation of the consecutive correlation matrix  $\mathbf{A}_c$ . Therefore, we merge these two matrices into a unified matrix  $\mathbf{A}_h$  to reveal observed and unobserved preferences of transitional dependencies. Specifically, given two distinct POI nodes  $\mathcal{V}_i^i$  and  $\mathcal{V}_j^j$ , its correlation score  $A_h^{ij}$  is defined as:

$$A_h^{ij} = \begin{cases} A_c^{ij}, & \text{if } A_c^{ij} \neq 0; \\ A_g^{ij}, & \text{others.} \end{cases} \quad (6)$$

For any POI node  $\mathcal{V}_l^l$ , we embed each POI node to a unified representation:

$$\mathbf{s}_l^l = \mathbf{A}_h^l \mathbf{W}_l + \mathbf{b}_l, \quad (7)$$

where  $\mathbf{W}_l \in \mathbb{R}^{|\mathcal{L}| \times d}$  and  $\mathbf{b}_l \in \mathbb{R}^d$  are trainable matrices. The dimension of  $\mathbf{s}_l^l$  is  $d$ . Afterwards, each POI has its unique initial representation. To bridge the correlation between POI  $\mathcal{V}_l^l$  and each of its neighbor  $\mathcal{V}_j^j \in \Omega(\mathcal{V}_l^l)$ , we devise a scoring function to evaluate the different contributions of neighboring nodes. For instance, given POI node  $\mathcal{V}_l^l$  and its neighbor  $\mathcal{V}_j^j$ , we define contribution measure as:

$$a(\mathbf{s}_i^l, \mathbf{s}_j^l) = \mathbf{b}_a^T [\mathbf{s}_i^l \oplus \mathbf{s}_j^l], \quad (8)$$

where  $\oplus$  is the concatenation operation and  $\mathbf{b}_a \in \mathbb{R}^{2d}$  is a learnable vector. Then, we follow the standard GAT [50] and use softmax function to normalize the attention scores across all neighbors of POI  $\mathcal{V}_l^l$ , where each attention score regarding its neighbour  $\mathcal{V}_k^k$  can be formulated as:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(a(\mathbf{s}_i^l, \mathbf{s}_j^l)))}{\sum_{k \in \Omega(\mathcal{V}_l^l)} \exp(\text{LeakyReLU}(a(\mathbf{s}_i^l, \mathbf{s}_k^l)))}. \quad (9)$$

In the end, we obtain the aggregated representation  $e_i^l$  of POI node  $\mathcal{V}_i^l$  by a sum operation:

$$e_i^l = \sigma\left(\sum_{j \in \Omega(\mathcal{V}_i^l)} \alpha_{ij} s_j^l \mathbf{W}_e\right), \quad (10)$$

where  $\sigma$  is the sigmoid activation function and  $\mathbf{W}_e \in \mathbb{R}^{d \times d}$  is a set of trainable parameters.

**Heterogeneous Semantic Aggregation (HeSA).** HeSA is to aggregate the associated information of POIs from the neighboring nodes with different types. In our PGraph, there are two correlations that describe the heterogeneous semantics between different types of nodes, i.e., the time-aspect and activity-aspect interests. In contrast to HoSA, we do not involve the attention mechanism to quantify the different contributions of POI's heterogeneous neighbors. The reason is that the number of them are extremely smaller than that of the POIs, we thus attempt to capture all of the possible heterogeneous neighbors of a given POI directly to enhance the semantic information.

(1) For time-aspect interest, each POI  $p_i$  is associated with a probability distribution  $\mathbf{A}_t^i$  ( $\in \mathbf{A}_t$ ) that describes the preference strengths between POI and time bins. We leverage the message-passing neural network inspired by [11] to incorporate the time-aspect preference of each POI, which can be formulated as:

$$e_i^t = \tanh\left(\mathbf{A}_t^i \mathbf{W}_t\right), \quad (11)$$

where  $\tanh$  is the activation function and  $\mathbf{W}_t$  is a trainable matrix. Finally, we can obtain each POI's temporal context.

(2) For activity-aspect interest, we obtain each POI's activity-aware semantic by:

$$e_i^a = \tanh\left(\mathbf{A}_a^i \mathbf{W}_a\right), \quad (12)$$

where  $\mathbf{W}_a$  is a trainable matrix. Finally, the homogeneous and heterogeneous semantics behind each POI are acquired by HoSA and HeSA, respectively. In the following mobility encoding procedures, we will use these contextual representations as the embeddings of POIs in user trajectories. And these embeddings can be jointly optimized during self-supervised learning and task learning.

## 4.2 Context-aware Mobility Encoding

Existing studies usually choose the recurrent neural networks such as Long-short Term Memory (LSTM) or Gated Recurrent Unit (GRU) to capture human transitional regularities. Since the complex stacked gate operations in LSTM typically struggle with the gradient vanishing problem, we select GRU as the kernel of our mobility encoder. For each  $l_\tau$  in a given trajectory  $T = \{l_1, l_2, \dots, l_n\}$ , we have collected the homogeneous and heterogeneous semantics behind it. In this way, they can be viewed as reflections of different interest in different domains. Therefore, we extend the GRU cell to capture the sequential information as well as the contextual information behind each POI. Correspondingly, the recursive process with GRU can be formulated as follows:

$$\mathbf{c}_\tau = [e_\tau^l \oplus e_\tau^t \oplus e_\tau^a] \mathbf{W}_f + \mathbf{b}_f, \quad (13)$$

$$\mathbf{h}_\tau = \text{GRU}(\mathbf{c}_\tau, \mathbf{h}_{\tau-1}), \quad (14)$$

where  $\mathbf{h}_\tau$  and  $\mathbf{h}_{\tau-1}$  are the hidden states of the current POI  $l_\tau$  and the last visited POI  $l_{\tau-1}$ , respectively. Herein,  $\mathbf{c}_\tau$  is the contextual embedding of POI  $l_\tau$ , which is a unified representation that integrates the homogeneous and heterogeneous semantics of POI  $l_\tau$  (they include  $e_\tau^l$ ,  $e_\tau^t$ , and  $e_\tau^a$ ). In addition,  $\mathbf{W}_f$  and  $\mathbf{b}_f$  are trainable parameters.

## 4.3 Self-supervised Disentanglement Learning

Now we describe in detail the self-supervised disentanglement learning in SSDL.

### 4.3.1 Disentanglement via Variational Inference

Given any recent trajectory  $T = \{l_1, l_2, \dots, l_n\}$ , we attempt to learn a set of time-varying variables  $z_{1:n}^r = \{z_1^r, z_2^r, \dots, z_n^r\}$  and a time-invariant variable  $z^s$ , where  $z_{1:n}^r$  aims at exploring the dynamics of human time-dependent interests while  $z^s$  undertakes the role of learning human inherent time-independent periodicity (habits). Formally, let  $z_\tau$  be the entangled latent code of check-in  $l_\tau$ , and we have  $z_\tau = (z_\tau^r, z^s)$ . For consistency, let  $l_{1:n}$  denote the check-in sequence  $\{l_1, l_2, \dots, l_n\}$ . As people's future movements are affected by their previous check-in behaviors, we assume that each  $z_\tau$  depends on its previous states  $z_{<\tau} = \{z_1, z_2, \dots, z_{\tau-1}\}$ . In addition, as user's long-standing interests will not be changed dramatically by recent activities, we assume that  $z_{1:n}^r$  and  $z^s$  are independent from each other, i.e.,  $p(z_{1:n}) = p(z_{1:n}^r)p(z^s)$ . Hence, we formulate our probabilistic generative model as follows:

$$\begin{aligned} \text{Prior: } p(l_{1:n}, z_{1:n}) &= p(z_{1:n})p(l_{1:n}|z_{1:n}) \\ &= [p(z^s) \prod_{\tau=1}^n p(z_\tau^r|z_{<\tau}^r)] \cdot \prod_{\tau=1}^n p(l_\tau|z_\tau^r, z^s), \end{aligned} \quad (15)$$

where  $p(z_{1:n})$  is a prior. Herein, we choose the Gaussian distribution as the prior  $p(z^s)$ , i.e.,  $p(z^s) \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ . We follow the rule of standard variational Bayes and set  $\mathcal{N}(\mu(z_{<\tau}), \sigma^2(z_{<\tau}))$  as  $p(z_\tau^r|z_{<\tau}^r)$ , where  $\mu(\cdot)$  and  $\sigma^2(\cdot)$  can be modeled by popular recursive networks. In practice, we also use GRU to obtain  $z_\tau^r$  as follows:

$$\mathbf{h}_\tau^r = \text{GRU}(\mathbf{h}_\tau, z_{\tau-1}^r), z_\tau^r = \mu(\mathbf{h}_\tau^r) + \sigma(\mathbf{h}_\tau^r) \odot \epsilon, \quad (16)$$

where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$  and  $\odot$  is element-wise multiplication.

Subsequently, we expect to produce a posterior distribution  $q(z_{1:n}|l_{1:n})$  to cater to the learning manner of variational inference. Thus, we define the posterior as follows:

$$\begin{aligned} \text{Posterior: } q(z_{1:n}|l_{1:n}) &= q(z^s, z_{1:n}^r | l_{1:n}) \\ &= q(z_{1:n}^r | l_{1:n})q(z^s | l_{1:n}) \\ &= q(z^s | l_{1:n}) \prod_{\tau=1}^n q(z_\tau^r | z_{<\tau}^r, l_{\leq \tau}). \end{aligned} \quad (17)$$

We note that the above process is also operated in an auto-regressive manner. We use another GRU cell that has the same architecture as the Mobility Encoding network to generate posterior distributions. At last, we obtain the Evidence Lower Bound (ELBO) as follows:

$$\begin{aligned} \text{ELBO: } \max_{p, q} \mathbb{E}_{l_{1:n} \sim p_D} \mathbb{E}_{q(z_{1:n}|l_{1:n})} [\log p(l_{1:n} | z_{1:n})] \\ - KL[q(z_{1:n} | l_{1:n}) \| p(z_{1:n})], \end{aligned} \quad (18)$$

where  $p_D$  is the empirical trajectory distribution. As  $z_{1:n}$  is comprised of mutually independent  $z^s$  and  $z_{1:n}^r$ , the second term of KL-divergence can be disentangled as:

$$\begin{aligned} KL[q(z_{1:n} | l_{1:n}) \| p(z_{1:n})] &= \\ KL[q(z^s | l_{1:n}) \| p(z^s)] + KL[q(z_{1:n}^r | l_{1:n}) \| p(z_{1:n}^r)]. \end{aligned} \quad (19)$$

Following the principle of VAE [38], [44], we present a theoretical proof to illustrate above modeling processes.

**Proof 1.**

$$\begin{aligned} &\log p(l_{1:n}) \\ &\geq -KL[q(z_{1:n} | l_{1:n}) \| p(z_{1:n} | l_{1:n})] + \log p(l_{1:n}) \\ &= -KL[q(z^s, z_{1:n}^r | l_{1:n}) \| p(z^s, z_{1:n}^r | l_{1:n})] + \log p(l_{1:n}) \\ &= \mathbb{E}_{q(z^s, z_{1:n}^r | l_{1:n})} [\log p(l_{1:n} | z^s, z_{1:n}^r) \\ &\quad - \log q(z^s, z_{1:n}^r | l_{1:n}) + \log p(l_{1:n})] \\ &= \mathbb{E}_{q(z^s, z_{1:n}^r | l_{1:n})} [\log p(l_{1:n} | z^s, z_{1:n}^r) \\ &\quad - \log q(z^s, z_{1:n}^r | l_{1:n}) + \log p(z^s, z_{1:n}^r)] \\ &= \mathbb{E}_{q(z^s, z_{1:n}^r | l_{1:n})} [\log p(l_{1:n} | z^s, z_{1:n}^r) - \log q(z^s | l_{1:n}) \\ &\quad - \log p(z_{1:n}^r | l_{1:n}) + \log p(z^s) + \log p(z_{1:n}^r)] \\ &= \mathbb{E}_{q(z^s, z_{1:n}^r | l_{1:n})} [\log p(l_{1:n} | z^s, z_{1:n}^r) \\ &\quad - KL[q(z^s | l_{1:n}) \| p(z^s)] \\ &\quad - KL[q(z_{1:n}^r | l_{1:n}) \| p(z_{1:n}^r)]. \end{aligned} \quad (20)$$

Recall that the results of the above proof are similar to the results in standard VAE, which usually confront the agnostic prior distribution that causes posterior collapse problem and leaves the learned latent space still entangled [51]. This phenomenon has been revealed in recent studies, e.g.,  $\beta$ -VAE [36] and  $\beta$ -TCVAE [37]. Additionally, [52] provides us with a clearer perspective that reveals the challenges with disentangled representation in variational inference. Thus, we conjecture that the last two terms regularized by KL-divergence in Eq.(20) are hard to close to their corresponding prior, which would make each posterior become non-informative. For the purpose of receiving clean disentanglement of  $z^s$  and  $z_{1:n}^r$ , we are inspired by recent self-supervised learning and enforce disentanglement mobility learning from the perspective of mutual information.

#### 4.3.2 Mutual Information Regularization

We now turn to detail on how to combine contrastive learning with disentangled mobility learning. We first introduce variational mobility learning from the perspective of Mutual Information (MI). The goal of MI is a measure of the mutual dependence between two variables. Since both  $z^s$  and  $z_{1:n}^r$  are derived from the original trajectory, we thus add three additional MI terms to regularize the latent space of them, which can be defined as follows:

$$\begin{aligned} \mathcal{J}_{self} &= \max_{p,q} \mathbb{E}_{l_{1:n} \sim p_D} \mathbb{E}_{q(z_{1:n} | l_{1:n})} [\log p(l_{1:n} | z_{1:n}) \\ &\quad - \alpha(KL[q(z^s | l_{1:n}) \| p(z^s)] + KL[q(z_{1:n}^r | l_{1:n}) \| p(z_{1:n}^r)]) \\ &\quad + \beta(MI(z^s; l_{1:n}) + MI(z_{1:n}^r; l_{1:n})) - \gamma MI(z_{1:n}^r; z^s), \end{aligned} \quad (21)$$

where  $\alpha$ ,  $\beta$  and  $\gamma$  are weight coefficients.  $MI(\cdot, \cdot)$  refers to MI term. For instance,  $MI(z^s; l_{1:n})$  is defined as:

$$\mathbb{E}_{q(z^s, l_{1:n})} \left[ \log \frac{q(z^s | l_{1:n})}{q(z^s)} \right]. \quad (22)$$

We note that other MI terms have the similar formulation. Now our goal become enforcing the posteriors matching with their corresponding priors while ensuring that  $z^s$  and  $z_{1:n}^r$  are disentangled from each other. Note that the complete proof of Eq.(21) is provided in Appendix part. To estimate the MI terms, we follow most of recent studies [33], [45], [46] and employ the NCE loss to make contrastive estimation. For instance, a contrastive estimation of  $MI(z^s; l_{1:n})$  can be defined as follows,

$$\mathcal{C}_{z^s} \approx \mathbb{E}_{p_D} \log \frac{\psi(z^s, l_{1:n}^+)}{\psi(z^s, l_{1:n}^+) + \sum_{j=1}^m \psi(z^s, \tilde{l}_{1:n}^j)}, \quad (23)$$

where  $\psi(\cdot, \cdot) = \exp(\text{sim}(\cdot, \cdot)/\eta)$ ,  $\text{sim}(\cdot, \cdot)$  denotes the cosine similarity function,  $m$  is the number of negative trajectories, and  $\eta = 0.5$  is a temperature parameter. Notably, we treat  $l_{1:n}$  as the *positive* trajectory sequence regarding  $z^s$  and specify it using  $l_{1:n}^+$ . Besides,  $\tilde{l}_{1:n}^j$  refers to a *negative* sample (trajectory), which is generated from other users.

**Augmentation for time-invariant factor.** However, due to the limited scale of *positive* samples, we try to generate more realistic trajectories to augment the original samples. As  $z^s$  reveals human intrinsic periodicity and should not be affected by recent moving behaviors, i.e., time-independent, we can thus randomly change the order of a given trajectory  $l_{1:n}$  and formulate several augmentation versions w.r.t  $l_{1:n}$ . We claim that it is a simple but efficient strategy to obtain rich augmented samples. Correspondingly, the contrastive estimation regarding these augmented samples can be denoted as follows:

$$\mathcal{C}_{z^s} \approx \mathbb{E}_{p_D} \log \frac{\psi(z^s, l_{1:n}^\#)}{\psi(z^s, l_{1:n}^\#) + \sum_{j=1}^m \psi(z^s, l_{1:n}^j)}, \quad (24)$$

where  $l_{1:n}^\#$  indicates it is an augmentation version of  $l_{1:n}$ . For time-invariant factor  $z^s$ , we can use the collected samples including augmented samples to make a final estimates as follows:

$$MI(z^s; l_{1:n}) \approx \frac{1}{2}(\mathcal{C}_{z^s} + \mathcal{C}_{z^s}^\#). \quad (25)$$

**Augmentation for time-varying factor.** As for  $z_{1:n}^r$  is a set of latent variables regarding  $l_{1:n}$ , showing human human time-dependent interests. Similar to Eq. (23), we can obtain the contrastive estimation of  $MI(z_{1:n}^r; l_{1:n})$  as  $\mathcal{C}_{z_{1:n}^r}$ . Furthermore, we provide another data augmentation method to enhance the optimization of  $MI(z_{1:n}^r; l_{1:n})$ . The intuition is that  $z_{1:n}^r$  is a set of time-dependent variables. In practice, real-world check-in data could be subject to noise and uncertainty due to the presence of collective POIs [53]. Hence, a user usually posts a fuzzy POI to replace her accurate position, which could weaken human mobility pattern learning and even result in inaccurate predictions. Motivated by [53], [54], it is encouraging that we can use any member of related collective POI to replace the original POI in a trajectory to obtain an augmentation trajectory for time-varying factor training, which will not change any temporal semantics. In addition, another potential benefit of such a practice is to alleviate the uncertainty issue behind diverse human check-in behaviors. In our implementation, we use neighbors of the same category within 300m of a given POI as members

of its collective POI and replace about 30% of the POIs in a given trajectory with their related collective POIs, which does not heavily affect the full semantics behind the original trajectory. As a result, we can obtain a large number of synthetic trajectories that provide multiple views of a given trajectory. Similar to Eq.(25), we can get the final estimation regarding  $z_{1:n}^r$  as follows:

$$MI(z_{1:n}^r; l_{1:n}) \approx \frac{1}{2}(C_{z_{1:n}^r} + C_{z_{1:n}^r}^*), \quad (26)$$

where  $C_{z_{1:n}^r}^*$  is the contrastive estimation of the augmented trajectories regarding time-varying factors. As for the final term  $MI(z_{1:n}^r; z^s)$ , the variables in it are all in the latent space, we thus can directly choose the standard mini-batch weighted sampling (MWS) [37] for comparative estimation.

#### 4.4 Task Learning

So far, we have obtained a set of time-varying variables and a time-invariant variable for each trajectory. We turn to use our task learning network to predict the next POI. For each user, we actually own her entire historical trajectory. Inspired by [6], [55] modeling such a long trajectory would boost the capture of human long-term transitional preferences. It is natural to adopt the RNN to encode the transitional regularity underlying human historical trajectory. But in practice, there are massive time-ordered POIs in their historical trajectories, which usually result in a serious time cost problem. Therefore, we employ a self-attention layer with position encoding to capture the taste of the transitional behavior of a user as well as the long-distance dependencies. Given a user’s entire historical trajectory  $\mathcal{T}_{1:\mathcal{K}}$  containing  $K$  ordered POIs. We first reuse the linear layer (cf. Eq.13) to obtain the dense representation of each POI in  $\mathcal{T}_{1:\mathcal{K}}$ . Correspondingly, we use  $\mathcal{T}_{1:\mathcal{K}}$  to denote the trajectory with embedded POIs. To determine the order of POIs in  $\mathcal{T}_{1:\mathcal{K}}$ , we follow [56] and use the sine/cosine function-based position embedding to formulate the final representation of each POI, which can be denoted as:

$$\mathcal{T}'_i = \mathcal{T}_i + \Phi(\mathcal{T}_i), \quad (27)$$

where  $\Phi(\mathcal{T}_i)$  is the position embedding of POI  $l_i$  in  $\mathcal{T}_{1:\mathcal{K}}$ . Then, we employ one-layer self-attention network to receive a set of hidden states regarding  $\mathcal{T}$ , as follows:

$$\mathbf{H}_{1:\mathcal{K}} = \text{self-att}(\mathcal{T}'_{1:\mathcal{K}}). \quad (28)$$

In our study, we use the last state  $\mathbf{H}_{\mathcal{K}}$  to represent  $\mathcal{T}_{1:\mathcal{K}}$  and regard it as one of the inputs for task prediction.

Now we take  $z_{1:n}^r, z^s$ , and  $\mathbf{H}_{1:\mathcal{K}}$  as the input and employ a one-layer fully-connected network with softmax function to obtain the predict POI. the process can be expressed as:

$$\tilde{l}_{n+1} = \arg \max(\text{softmax}([z_n^r \oplus z^s \oplus \mathbf{H}_{\mathcal{K}}] \mathbf{W}_t + \mathbf{b}_t)) \quad (29)$$

Correspondingly, the loss function for trajectory  $T$  can be expressed as:

$$\mathcal{L}_T = -l_{n+1} \log \tilde{l}_{n+1}. \quad (30)$$

To minimize the above cross-entropy loss, we employ Adam algorithm to optimize the parameters. We outline the complete pipeline of training SSDL in Algorithm 1.

---

#### Algorithm 1: The pipeline of training SSDL.

---

```

Input: POI set  $\mathcal{L}$ ; Historical trajectories  $\mathcal{T}$  and current trajectory  $T$  of users.
/* Common Interest Distillation */
1 Build the PGraph from entire trajectory data;
2 Generate the homogeneous and heterogeneous semantics via HoSA and HeSA for each POI;
/* Disentanglement learning */
3 repeat
4   foreach  $T$  do
5     Compute each POI embedding via Eq. (13);
6     Obtain each hidden state  $h_\tau$  via Eq.(14);
7     Compute each  $z_\tau^r$  via Eq.(16);
8     Compute  $z^s$  based on the last hidden state;
9     Make trajectory augmentation regarding  $z^s$ ;
10    Compute  $MI(z^s; l_{1:n})$  through Eq. (25);
11    Make trajectory augmentation regarding  $z_{1:n}^r$ ;
12    Compute  $MI(z_{1:n}^r; l_{1:n})$  through Eq. (26);
13    Update the parameters by maximizing Eq.(21);
14  end
15 until convergence;
/* Task learning */
16 repeat
17   foreach  $T$  do
18     Compute each POI embedding via Eq. (13);
19     Obtain each hidden state  $h_\tau$  via Eq.(14);
20     Compute each  $z_\tau^r$  via Eq.(16);
21     Compute  $z^s$  based on the last hidden state;
22     Model user historical trajectory via Eq.(28);
23     Obtain the predicted POI through Eq.(29);
24     Update the parameters according to Eq.(30);
25   end
26 until convergence;
Output: Trained model.

```

---

## 5 EXPERIMENTS

We now conduct experiments to evaluate the performance of our proposed SSDL on four real-world datasets.

### 5.1 Experimental Settings

#### 5.1.1 Datasets

To facilitate reproducible results, we conduct all experiments on two publicly available LBS applications: Foursquare [57] and Gowalla [58]. Foursquare contains check-ins in New York and Tokyo collected from 12 April 2012 to 16 February 2013. Each check-in has a timestamp, GPS coordinates and semantics about it. In Gowalla, we select the data from two cities, i.e., Los Angeles and Houston. Following previous studies [6], [33], we filter out the POIs visited by fewer than eight times. For each user, we concatenate his/her all chronological check-ins and divide each trajectory into subsequence with the time interval of 24 hours. To specify whether check-ins are collected on weekdays or weekends, we further assign 48 time slots to each check-in time. We take each user’s first 80% trajectories as the training set, the remaining 20% as the test set. The statistics of four datasets are summarized in Table 1.

TABLE 1  
Statistics of the datasets.

City	Users	POIs	Check-ins	Trajectories
Tokyo	2102	6789	240056	60365
New York	990	4211	79006	23252
Los Angeles	2346	8676	195231	61542
Houston	1351	6994	121502	37514



### 5.1.2 Baselines

We compare our SSDL with several representative approaches for next POI prediction task.

- GRU [21] is a common approach for sequential data learning as its superiority in incorporating the semantics of long-term dependencies.
- ST-RNN [22] is an RNN-based method that incorporates spatio-temporal contexts when predicting the next POI.
- HST-LSTM [12] employs sequence-to-sequence learning scheme to include spatial-temporal influence in LSTM and makes use of contextual information to enhance model performance for sparse data prediction.
- Flashback [59] models sparse user mobility footprints by doing flashbacks on hidden states in RNNs. Especially, it explicitly employs the spatio-temporal contexts to search past hidden states with high predictive power. In our experiments, we take the GRU cell as the recurrent component in Flashback for a fair comparison.
- DeepMove [55] presents an attention-based RNN to encode human recent trajectories. Furthermore, it employs another RNN to learn user long-term preferences from historical trajectories.
- VANext [6] proposes a novel variational attention mechanism to explore human periodic regularities. In addition, it employs a simple convolutional neural network rather than RNN to capture human long-term interests.
- PLSPL [10] is a unified framework that jointly learns users' long- and short-term interests for next POI prediction.
- MobTCast [2] is a Transformer-based approach that considers multiple semantic contexts behind check-ins to enhance the understanding of human mobility. Note that we remove the Social Context Extractor in MobCast as the social relationships are not available in our context.
- $\beta$ -VAE [36] is a widely used representation learning method that is able to separate latent factors into different space by using an adjustable hyperparameter  $\beta$  to the original VAE objective. In this study, we use GRU as the encoder and decoder network structure in  $\beta$ -VAE to model the temporal semantics.
- SML [33] attempts to understand human mobility in a self-supervised learning manner. Especially, it leverages heuristic strategy to enumerate massive different views of original sparse trajectories for contrastive estimation.

### 5.1.3 Metrics

To evaluate the performance of our proposed SSDL, we follow most of the previous studies [6], [10], [33] and select three commonly used metrics to compare with the baselines. We first use the ACC@K to evaluate the recommendation performance. In this paper, we report the different testing results of  $K = 1, 5, 10$ . Additionally, we report area under the ROC curve (AUC) and mean average precision (MAP) metrics that are frequently used in classification tasks.

### 5.1.4 Implementation Details

We implement our SSDL and baselines in Python. All methods are based on the Torch library and accelerated by one NVIDIA GTX 1080 GPU. Besides, we choose Adam [60] to train all deep learning methods. In disentanglement learning, the learning rate is initialized as 0.01. We set the

coefficient  $\alpha$  of KL terms to 1. Besides,  $\beta$  and  $\gamma$  are fixed to be 1 and 0.1. In task learning, the learning rate is initialized with  $5e-4$ . The dropout rate is set as 0.5, and the batch size is 32. The hidden size of the self-attention network is set to 300. In addition, we set dimension of  $z^s$  to 256, while  $z_{1:n}^r$  to 32.

## 5.2 Performance Comparisons

Table 2 reports the performance of different approaches on the datasets of four cities, where the best achievement is highlighted with **bold** and the second best is marked with underline. Specifically, we have the following observations.

We can find that ST-RNN does not provide us with competitive achievements compared to GRU although it considers the spatial and temporal constraints. The plausible reason is that the sparsity issue of check-in data heavily affects the distillation of semantics contexts such as geographical distance. Meanwhile, relying on simple spatio-temporal features and regarding the next POI as the solo supervision usually results in an inference bias or uncertainty problem due to the boundary of available training datasets. To mitigate the data sparsity issue, HST-LSTM which combines spatial and temporal factors with a gate mechanism is able to boost the capture of human mobility patterns by a large margin. Furthermore, HST-LSTM models the periodicity of consecutive check-ins in an end-to-end manner, which brings an encouraging prospect for us to learn the complex distribution behind historical trajectories. Compared with HST-LSTM, which directly adds spatio-temporal factors to hidden states, Flashback achieves competitive performance because it explicitly uses a rich spatio-temporal context to search for past hidden states with high predictive power to predict the next POI.

As for DeepMove and VANext, they both attempt to correlate a certain user's recent trajectory and historical trajectory to accurately discover individual periodicity. Our experiments show that they achieve higher gains than models (e.g., ST-RNN) that only consider the past few check-ins. Furthermore, VANext, the first variational inference approach to model human trajectories using a prior assumption, outperforms DeepMove due to the relief of the inherent uncertainty of user mobility. The paradigm of PLSPL is similar to DeepMove, but it operates an attention mechanism to evaluate the importance of each POI in a user's historical check-ins, aiming at exploring the tastes of different users. We can find that PLSPL performs better than DeepMove. In addition, MobTCast is a Transformer-based approach that uses self-attention to study the interactive signals between POIs in a given trajectory, as well as multiple semantic contexts, such as category and temporal semantics. We obtain similar performance results compared to PLSPL, indicating that considering multiple semantic contexts does help to discover users' future check-in intentions.

$\beta$ -VAE is a popular disentanglement learning method that also obtains promising results, which suggests that employing the latent variables produced by variational Bayesian does help in understanding the inherent generative factors underlying human mobility. As for SML, it is the first self-supervised learning solution for the next POI prediction, achieving the best gains on AUC among the

TABLE 2  
Performance comparisons on four cities.

Method	Tokyo					New York				
	ACC@1	ACC@5	ACC@10	AUC	MAP	ACC@1	ACC@5	ACC@10	AUC	MAP
GRU	13.11	27.88	34.28	88.01	7.36	15.37	31.73	36.10	81.40	8.59
ST-RNN	13.38	29.20	36.45	89.82	7.41	13.50	32.86	40.05	81.91	8.20
HST-LSTM	18.70	39.14	46.47	90.66	9.82	17.48	42.77	50.82	86.25	8.53
Flashback	18.23	39.42	46.66	90.40	10.47	22.22	49.52	57.11	87.74	13.59
DeepMove	19.92	40.61	48.25	90.47	12.21	21.56	45.09	52.17	87.30	13.06
VANext	20.21	44.49	52.63	91.30	12.36	22.54	51.26	58.78	89.30	14.02
PLSPL	20.19	43.64	52.45	91.37	12.86	23.02	53.33	63.34	89.21	14.83
MobTCast	19.58	43.41	51.95	89.95	11.67	22.37	54.31	64.18	88.59	14.03
$\beta$ -VAE	20.10	44.78	53.89	91.35	12.75	22.26	50.71	58.68	89.38	14.07
SML	20.25	44.70	53.58	91.48	12.51	22.62	52.16	60.18	90.17	14.74
<b>SSDL</b>	<b>22.93</b>	<b>46.80</b>	<b>55.31</b>	<b>92.39</b>	<b>14.99</b>	<b>25.07</b>	<b>56.78</b>	<b>65.72</b>	<b>90.72</b>	<b>16.60</b>

Method	Los Angeles					Houston				
	ACC@1	ACC@5	ACC@10	AUC	MAP	ACC@1	ACC@5	ACC@10	AUC	MAP
GRU	10.11	19.05	22.67	78.07	4.97	10.74	18.01	21.33	80.47	5.99
ST-RNN	10.01	19.30	23.64	80.48	5.11	11.33	20.21	24.81	82.34	6.88
HST-LSTM	12.01	23.97	29.04	82.57	5.29	13.41	22.86	27.21	82.58	6.58
Flashback	13.81	25.60	30.32	83.74	7.65	14.37	24.52	28.70	84.54	8.79
DeepMove	13.31	25.73	30.35	82.41	7.26	14.13	24.59	29.03	83.29	8.46
VANext	14.36	27.91	33.16	86.22	7.73	14.88	26.78	31.44	86.06	8.31
PLSPL	14.92	28.26	33.86	84.34	7.97	16.06	29.22	34.67	86.11	9.74
MobTCast	14.30	28.62	33.38	83.41	7.65	15.40	28.27	32.87	83.87	8.66
$\beta$ -VAE	14.39	27.43	32.96	85.95	7.72	15.00	27.09	31.89	86.10	8.23
SML	14.77	28.12	33.35	86.38	7.86	15.15	27.46	32.31	86.48	9.17
<b>SSDL</b>	<b>15.94</b>	<b>31.02</b>	<b>36.80</b>	<b>86.98</b>	<b>8.72</b>	<b>16.91</b>	<b>30.92</b>	<b>36.56</b>	<b>87.30</b>	<b>10.70</b>

baselines. The reason is that it primarily seeks to produce massive synthetic trajectories for data augmentation and leverage contrastive learning to study the diversity of human moving intents behind existing historical check-ins.

In general, our proposed SSDL significantly outperforms the compared approaches by a relatively large margin across the four cities. For instance, SSDL respectively yields 8.57% and 11.94% averaged improvement over the best baseline regarding ACC@1 and MAP. This observation demonstrates the superiority of the self-supervised disentanglement learning paradigm in our SSDL.

### 5.3 Ablation Study

To evaluate the contribution of different components of SSDL, we devise several variants regarding SSDL from two aspects, i.e., POI embedding and trajectory embedding. First, we select 7 popular embedding methods to scrutinize the efficacy of our POI embedding.

- **One-Hot** [55] is the simplest method that maps each POI to a unique vector without any semantic information.
- **Random** uses a dense matrix sampled from a Gaussian distribution to represent the POIs.
- **Word2vec** is a popular embedding technique in NLP, aiming at exploring the surrounding context of a given word. Also, it has successfully applied in POI embedding [14], [31]. We implement the skip-gram model for POI embedding.
- **Causal** [6] is a variant of word2vec that treats the previous footprints of the current POI as its semantic context to incorporate human practical transitional behaviors.
- **Deepwalk** [17] is a data augmentation method that builds a POI graph to integrate users' historical visiting interests and geographical proximity and then leverages the skip-gram technique for POI embedding.
- **GraphAE** is popular method for node embedding. In this paper, we treat each POI as a node and build the same graph as Deepwalk for POI embedding.
- **GraphVAE** is a variant of GraphAE, taking the advantage of VAE for POI embedding.

Fig. 2 reports the performance of SSDL using different POI embedding methods. We can observe that our embedding method achieves the best results on the vast majority of metrics across the four cities, which indicates its higher effectiveness in capturing multiple human interests behind historical trajectory data.

Next, we turn to investigate the effectiveness of devised components in SSDL. Herein, we conduct the experiments with three SSDL variants. The details are shown as follows:

- **SSDL-Base** is a basic model that removes both graph-based embedding and mutual information regularization of SSDL. Instead, we use the word2vec technique for POI embedding.
- **SSDL w/o G** only removes the graph-based embedding and use the word2vec for POI embedding.
- **SSDL w/o H** only removes the mutual information regularization of SSDL.

Fig. 3 illustrates the performance of variants on four cities. First, we can find that removing any modules would bring significant performance degradation, suggesting that both modules in our SSDL benefit to enhance POI prediction. Second, SSDL-Base performs worse than SSDL w/o H across all cities, demonstrating that considering multiple common interests behind historical check-in data are useful to discover human mobility patterns. Third, SSDL w/o G outperforming SSDL-Base proves that our self-supervised disentanglement learning is an effective module to provide promising representations for task inference.

### 5.4 Disentanglement Interpretability

In this part, we focus on studying the disentangled representations from the interpretability aspect. We first investigate whether  $z^s$  and  $z_{1:n}^t$  can be well extracted from original trajectories and reflect human time-invariant periodicity/habits and time-varying interests, respectively. To this end, we randomly sample eight different users' trajectories and change their orders to generate several groups of trajectories. Then, we use the TSNE toolkit [61] to visualize the distribution time-invariant representations. We can find that

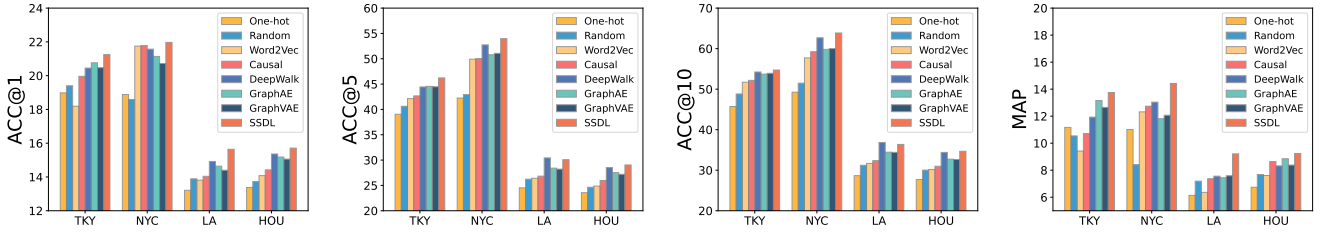


Fig. 2. Effect of POI embedding.

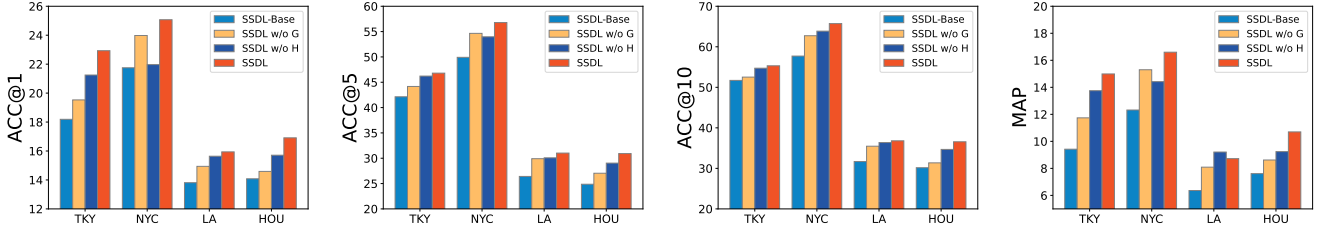


Fig. 3. Effects of different components in SSDL.

the representations of  $z^s$  produced by SSDL are grouped well, demonstrating that it can successfully separate the time-invariant factors to uncover the inherent preference of users that are not influenced by temporal factors. For  $z_{1:n}^r$ , we visualize the distribution of the last states  $z_n^r$  for simplicity, we can find they are entangled, indicating that they are really affected by the temporal factors. Therefore, we conclude that  $z^s$  and  $z_{1:n}^r$  indeed play well the roles of time-invariant and time-varying representations, respectively.

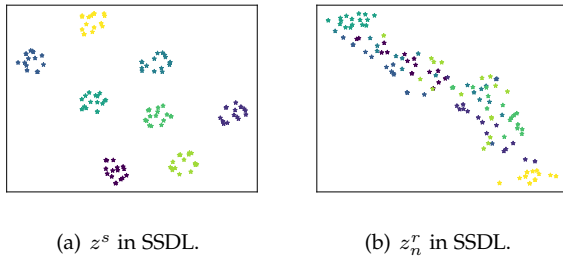


Fig. 4. The visualization of latent representations.

Besides, we also study the impact of our data augmentation approaches from a visualization perspective. We visualize the  $z^s$  distribution of randomly sampled trajectories of eight different users after task training. As shown in Fig. 5(a), we can find  $\beta$ -VAE can only separate the representations with a small margin. Fig. 5(b) presents the results of SSDL without any data augmentations, and Fig. 5(c) shows the results of SSDL that has no augmentation for time-invariant factors. Compare to Fig. 5(d), we can clearly find that both augmentations used in SSDL can significantly help us distinguish the trajectory representations of different users. This observation further suggests that different user movement patterns can be well refined by our SSDL.

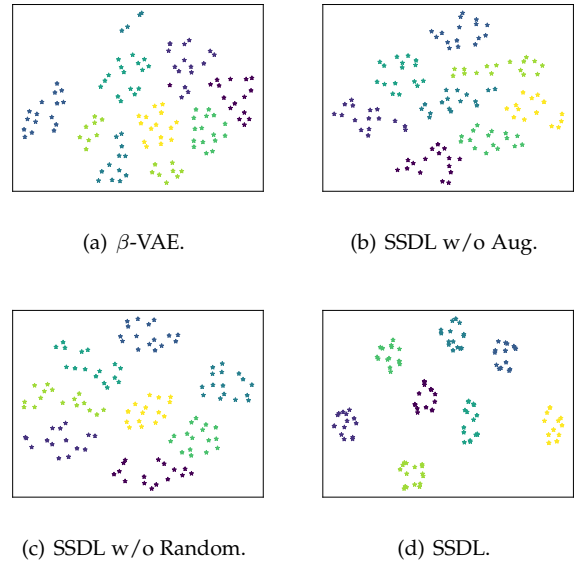


Fig. 5. The impact of augmentations on  $z^s$ .

5.5 Sensitivity Analysis

Finally, we investigate the impact of significant hyperparameters in our SSDL to evaluate the model’s robustness.

- *Weight coefficients.* The objective of our representation learning (cf. 21) contains three coefficients, which would determine the optimization procedure of each relative term. To this end, we generate different combinations of coefficients to investigate their impacts. The results of ACC@1 are shown in Fig. 6. We observe that  $\gamma=0.1$  obtains better performance than  $\gamma=1$  in general. We also find that the larger  $\beta$  helps to improve the accuracy of prediction since  $\beta$  represents the importance of mutual information between latent variables and trajectories. Finally, the weight coefficient  $\alpha$  cannot be too large, otherwise it

would constrain the performance.

- *Dimension of  $z^s$ .* Fig. 7 shows the performance variations of SSDL at different sizes of  $z^s$ . We find that the larger dimension of  $z^s$  does not give us promising results. Hence, for efficiency reasons, we set its dimensionality to 256.
- *Dimension of  $z_{1:n}^r$ .* Fig. 8 shows how different dimension of  $z_{1:n}^r$  would influence the performance of SSDL. The performance decreases when the dimension of  $z_{1:n}^r$  in each time step is larger than 32 and stays stable when the dimension increases. To obtain best performance, we set the dimension of  $z_{1:n}^r$  to 32 in our experiments.
- *Embedding size.* Embedding size is one of the critical factors affecting task prediction performance. Fig. 9 presents the effect of the embedding size. We can observe that the performance of SSDL climbs as the embedding size increases, and degrades or stays stable when the embedding size is larger than 256. In our experiments, we set the embedding size to 256.

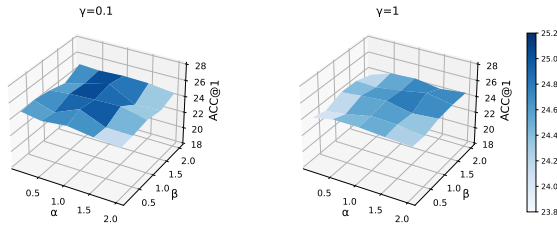


Fig. 6. The influence of weight coefficients in New York dataset.

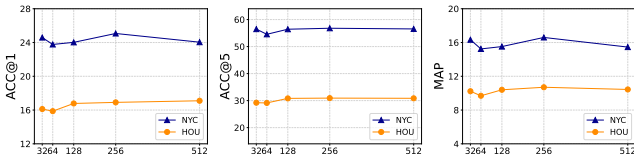


Fig. 7. The influence of the dimension of  $z^s$ .

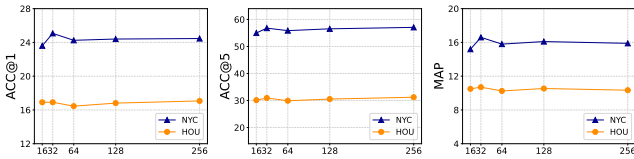


Fig. 8. The influence of the dimension of  $z_{1:n}^r$ .

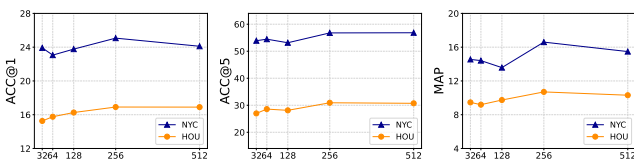


Fig. 9. The influence of embedding size.

## 6 CONCLUSION

In this paper, we present a self-supervised disentanglement learning framework, namely SSDL, to understand human mobility for tackling the next POI prediction problem. In contrast to existing sequential dynamics learning paradigms, SSDL mainly concentrates on disentangling the time-invariant and time-varying factors underlying massive sequential trajectories, which provides us an interpretable perspective to become familiar with human complex mobility patterns. Meanwhile, we present two practical trajectory augmentation strategies to relieve the sparsity issue of check-in data, which also enables the disentanglement of latent representations. Besides, we introduce a flexible graph structure learning method to incorporate multiple heterogeneous collaborative signals from historical check-ins. We believe that several other associated contexts such as social relations and textual data are also easily incorporated into our graph learning. Finally, our extensive experiments on four datasets demonstrate the superiority of SSDL compared to state-of-the-art baselines. As our future work, we plan to investigate the possible more intricate prior assumption during representation learning.

## ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China (Grant No.62102326 and No.62072077), the Key Research and Development Project of Sichuan Province under Grant 2022YFG0314, National Science Foundation SWIFT grant 2030249, and Guanghua Talent Project.

## REFERENCES

- [1] J. Feng, Y. Li, Z. Yang, M. Zhang, H. Wang, H. Cao, and D. Jin, "User identity linkage via co-attentive neural network from heterogeneous mobility data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 2, pp. 954–968, 2022.
- [2] H. Xue, F. Salim, Y. Ren, and N. Oliver, "Mobtcast: Leveraging auxiliary trajectory forecasting for human mobility prediction," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [3] Z. Luo and C. Miao, "Rlmob: Deep reinforcement learning for successive mobility prediction," in *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, 2022, pp. 648–656.
- [4] P. Zhao, A. Luo, Y. Liu, F. Zhuang, J. Xu, Z. Li, V. S. Sheng, and X. Zhou, "Where to go next: A spatio-temporal gated network for next poi recommendation," *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [5] S. Wang, J. Cao, and P. Yu, "Deep learning for spatio-temporal data mining: A survey," *IEEE transactions on knowledge and data engineering*, 2020.
- [6] Q. Gao, F. Zhou, G. Trajcevski, K. Zhang, T. Zhong, and F. Zhang, "Predicting human mobility via variational attention," in *The World Wide Web Conference*, 2019, pp. 2750–2756.
- [7] H. Zang, D. Han, X. Li, Z. Wan, and M. Wang, "Cha: Categorical hierarchy-based attention for next poi recommendation," *ACM Transactions on Information Systems (TOIS)*, vol. 40, no. 1, pp. 1–22, 2021.
- [8] W. Mathew, R. Raposo, and B. Martins, "Predicting future locations with hidden markov models," in *Proceedings of the 2012 ACM conference on ubiquitous computing*, 2012, pp. 911–918.
- [9] D. Massimo and F. Ricci, "Harnessing a generalised user behaviour model for next-poi recommendation," in *Proceedings of the 12th ACM Conference on Recommender Systems*, 2018, pp. 402–406.
- [10] Y. Wu, K. Li, G. Zhao, and Q. Xueming, "Personalized long- and short-term preference learning for next poi recommendation," *IEEE Transactions on Knowledge and Data Engineering*, 2020.

- [11] X. Rao, L. Chen, Y. Liu, S. Shang, B. Yao, and P. Han, "Graph-flashback network for next location recommendation," in *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 1463–1471.
- [12] D. Kong and F. Wu, "Hst-lstm: A hierarchical spatial-temporal long-short term memory network for location prediction," in *IJCAI*, vol. 18, no. 7, 2018, pp. 2341–2347.
- [13] Y. Luo, Q. Liu, and Z. Liu, "Stan: Spatio-temporal attention network for next location recommendation," in *Proceedings of the Web Conference 2021*, 2021, pp. 2177–2185.
- [14] S. Feng, G. Cong, B. An, and Y. M. Chee, "Poi2vec: Geographical latent representation for predicting future visitors," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [15] S. Zhao, T. Zhao, I. King, and M. R. Lyu, "Geo-teaser: Geotemporal sequential embedding rank for point-of-interest recommendation," in *Proceedings of the 26th international conference on world wide web companion*, 2017, pp. 153–162.
- [16] L. Huang, Y. Ma, Y. Liu, and K. He, "Dan-snr: A deep attentive network for social-aware next point-of-interest recommendation," *ACM Transactions on Internet Technology (TOIT)*, vol. 21, no. 1, pp. 1–27, 2020.
- [17] Q. Gao, F. Zhou, T. Zhong, G. Trajcevski, X. Yang, and T. Li, "Contextual spatio-temporal graph representation learning for reinforced human mobility mining," *Information Sciences*, 2022.
- [18] N. Lim, B. Hooi, S.-K. Ng, X. Wang, Y. L. Goh, R. Weng, and J. Varadarajan, "Stp-udgat: spatial-temporal-preference user dimensional graph attention network for next poi recommendation," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 845–854.
- [19] Y. Li, T. Chen, Y. Luo, H. Yin, and Z. Huang, "Discovering collaborative signals for next poi recommendation with iterative seq2graph augmentation," in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, Z.-H. Zhou, Ed. International Joint Conferences on Artificial Intelligence Organization, 8 2021, pp. 1491–1497, main Track. [Online]. Available: <https://doi.org/10.24963/ijcai.2021/206>
- [20] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [21] K. Cho, B. v. M. C. Gulcehre, D. Bahdanau, F. B. H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1724–1734.
- [22] Q. Liu, S. Wu, L. Wang, and T. Tan, "Predicting the next location: A recurrent model with spatial and temporal contexts," in *Thirtieth AAAI conference on artificial intelligence*, 2016.
- [23] F. Yu, L. Cui, W. Guo, X. Lu, Q. Li, and H. Lu, "A category-aware deep model for successive poi recommendation on sparse check-in data," in *Proceedings of the web conference 2020*, 2020, pp. 1264–1274.
- [24] K. Zhao, Y. Zhang, H. Yin, J. Wang, K. Zheng, X. Zhou, and C. Xing, "Discovering subsequence patterns for next poi recommendation." in *IJCAI*, 2020, pp. 3216–3222.
- [25] H. Sun, J. Xu, K. Zheng, P. Zhao, P. Chao, and X. Zhou, "Mfnf: A meta-optimized model for few-shot next poi recommendation," in *IJCAI International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence, 2021, pp. 3017–3023.
- [26] M. Zhang, Y. Yang, R. Abbas, K. Deng, J. Li, and B. Zhang, "Snpr: A serendipity-oriented next poi recommendation model," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 2568–2577.
- [27] C. Miao, Z. Luo, F. Zeng, and J. Wang, "Predicting human mobility via attentive convolutional network," in *Proceedings of the 13th International Conference on Web Search and Data Mining*, 2020, pp. 438–446.
- [28] Y. Chen, C. Long, G. Cong, and C. Li, "Context-aware deep model for joint mobility and time prediction," in *Proceedings of the 13th International Conference on Web Search and Data Mining*, 2020, pp. 106–114.
- [29] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [30] X. Liu, Y. Liu, and X. Li, "Exploring the context of locations for personalized location recommendations." in *IJCAI*, 2016, pp. 1188–1194.
- [31] Q. Gao, F. Zhou, K. Zhang, G. Trajcevski, X. Luo, and F. Zhang, "Identifying human mobility via trajectory embeddings," in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 2017, pp. 1689–1695.
- [32] S. Yang, J. Liu, and K. Zhao, "Getnext: Trajectory flow map enhanced transformer for next poi recommendation," in *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2022, pp. 1144–1153.
- [33] F. Zhou, Y. Dai, Q. Gao, P. Wang, and T. Zhong, "Self-supervised human mobility learning for next location prediction and trajectory classification," *Knowledge-Based Systems*, vol. 228, p. 107214, 2021.
- [34] F. Zhou, X. Liu, T. Zhong, and G. Trajcevski, "Metamove: On improving human mobility classification and prediction via meta-learning," *IEEE Transactions on Cybernetics*, 2021.
- [35] H. Tan, D. Yao, T. Huang, B. Wang, Q. Jing, and J. Bi, "Meta-learning enhanced neural ode for citywide next poi recommendation," in *2021 22nd IEEE International Conference on Mobile Data Management (MDM)*. IEEE, 2021, pp. 89–98.
- [36] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner, "Understanding disentangling in beta-vae," *arXiv e-prints*, pp. arXiv–1804, 2018.
- [37] R. T. Chen, X. Li, R. B. Grosse, and D. K. Duvenaud, "Isolating sources of disentanglement in variational autoencoders," *Advances in neural information processing systems*, vol. 31, 2018.
- [38] J. Bai, W. Wang, and C. P. Gomes, "Contrastively disentangled sequential variational autoencoder," *Advances in Neural Information Processing Systems*, vol. 34, pp. 10 105–10 118, 2021.
- [39] Y. Li and S. Mandt, "Disentangled sequential autoencoder," *arXiv preprint arXiv:1803.02991*, 2018.
- [40] Z. Ding, Y. Xu, W. Xu, G. Parmar, Y. Yang, M. Welling, and Z. Tu, "Guided variational autoencoder for disentanglement learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 7920–7929.
- [41] Y. Zhu, M. R. Min, A. Kadav, and H. P. Graf, "S3vae: Self-supervised sequential vae for representation disentanglement and data generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6538–6547.
- [42] C. Huang, X. Wang, X. He, and D. Yin, "Self-supervised learning for recommender system," in *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2022, pp. 3440–3443.
- [43] J. Ma, C. Zhou, H. Yang, P. Cui, X. Wang, and W. Zhu, "Disentangled self-supervision in sequential recommenders," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 483–491.
- [44] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [45] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola, "What makes for good views for contrastive learning?" *Advances in Neural Information Processing Systems*, vol. 33, pp. 6827–6839, 2020.
- [46] M. Gutmann and A. Hyvärinen, "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 297–304.
- [47] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [48] J. Jeon, S. Kang, M. Jo, S. Cho, N. Park, S. Kim, and C. Song, "Lightmove: A lightweight next-poi recommendation for taxicab rooftop advertising," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 3857–3866.
- [49] Y. Chen, X. Wang, M. Fan, J. Huang, S. Yang, and W. Zhu, "Curriculum meta-learning for next poi recommendation," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 2692–2702.
- [50] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=rJXmpikCZ>
- [51] K. Akuzawa, Y. Iwasawa, and Y. Matsuo, "Information-theoretic regularization for learning global features by sequential vae," *Machine Learning*, vol. 110, no. 8, pp. 2239–2266, 2021.
- [52] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Schölkopf, and O. Bachem, "Challenging common assumptions in the unsu-

pervised learning of disentangled representations,” in *international conference on machine learning*. PMLR, 2019, pp. 4114–4124.

- [53] Z. Sun, C. Li, Y. Lei, L. Zhang, J. Zhang, and S. Liang, “Point-of-interest recommendation for users-businesses with uncertain check-ins,” *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [54] L. Zhang, Z. Sun, J. Zhang, Y. Lei, C. Li, Z. Wu, H. Kloeden, and F. Klanner, “An interactive multi-task learning framework for next poi recommendation with uncertain check-ins,” in *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, 2021, pp. 3551–3557.
- [55] J. Feng, Y. Li, C. Zhang, F. Sun, F. Meng, A. Guo, and D. Jin, “Deepmove: Predicting human mobility with attentional recurrent networks,” in *Proceedings of the 2018 world wide web conference*, 2018, pp. 1459–1468.
- [56] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [57] D. Yang, D. Zhang, V. W. Zheng, and Z. Yu, “Modeling user activity preference by leveraging user spatial temporal characteristics in lbsns,” *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 45, no. 1, pp. 129–142, 2014.
- [58] Y. Liu, W. Wei, A. Sun, and C. Miao, “Exploiting geographical neighborhood characteristics for location recommendation,” in *Proceedings of the 23rd ACM international conference on conference on information and knowledge management*, 2014, pp. 739–748.
- [59] D. Yang, B. Fankhauser, P. Rosso, and P. Cudre-Mauroux, “Location prediction over sparse user mobility traces using rnns: flashback in hidden states!” in *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, 2021, pp. 2184–2190.
- [60] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [61] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne.” *Journal of machine learning research*, vol. 9, no. 11, 2008.

## 7 APPENDIX

Herein, we provide a theoretical details of our objective. Assume each  $z_\tau$  in  $z_{1:n}$  is the entangled latent code of check-in  $l_\tau$  (i.e.,  $z_\tau = \{z_\tau^r, z_\tau^s\}$ ). We attempt to learn a set of time-varying variables  $z_{1:n}^r = \{z_1^r, z_2^r, \dots, z_n^r\}$  and a time-invariant variable  $z^s$  from a given trajectory  $l_{1:n}$ . According to the Bayes rules and Variational Inference. We have:

**Proof 2.**

$$\begin{aligned}
& \log p(l_{1:n}) \\
& \geq -KL[q(z_{1:n}) \| p(z_{1:n} | l_{1:n})] + \log p(l_{1:n}) \\
& = -\mathbb{E}_{q(z_{1:n} | l_{1:n})} [\log \frac{q(z_{1:n})}{p(z_{1:n} | l_{1:n})}] + \log p(l_{1:n}) \\
& = -\mathbb{E}_{q(z_{1:n} | l_{1:n})} [\log q(z_{1:n}) - \log p(z_{1:n} | l_{1:n})] + \log p(l_{1:n}) \\
& = \mathbb{E}_{q(z_{1:n} | l_{1:n})} [\log p(l_{1:n}) - \log q(z_{1:n}) + \log p(z_{1:n} | l_{1:n})] \\
& = \mathbb{E}_{q(z_{1:n} | l_{1:n})} [\log p(l_{1:n}) - \log q(z_{1:n}) \\
& + \log \frac{p(c_{1:n} | z_{1:n})p(z_{1:n})}{p(c_{1:n})}] \\
& = \mathbb{E}_{q(z_{1:n} | c_{1:n})} [\log p(c_{1:n}) - \log q(z_{1:n}) \\
& + \log p(c_{1:n} | z_{1:n}) + \log p(z_{1:n}) - \log p(c_{1:n})] \\
& = \mathbb{E}_{q(z_{1:n} | c_{1:n})} [\log p(c_{1:n} | z_{1:n}) - (\log q(z_{1:n}) - \log p(z_{1:n}))] \\
& = \mathbb{E}_{q(z_{1:n} | c_{1:n})} [\log p(c_{1:n} | z_{1:n}) \\
& - \mathbb{E}_{q(z_{1:n} | c_{1:n})} [\log q(z_{1:n}) - \log q(z^s)q(z_{1:n}^r)] \\
& + \log q(z^s)q(z_{1:n}^r) - \log p(z_{1:n})] \\
& = \mathbb{E}_{q(z_{1:n} | c_{1:n})} [\log p(c_{1:n} | z_{1:n}) \\
& - \mathbb{E}_{q(z_{1:n} | c_{1:n})} [\log q(z_{1:n}) - \log q(z^s)q(z_{1:n}^r)] \\
& - \mathbb{E}_{q(z_{1:n} | c_{1:n})} [\log q(z^s)q(z_{1:n}^r) - \log p(z_{1:n})]
\end{aligned}$$

Since  $z_{1:n} = (z^s, z_{1:n}^r)$ , we thus have:

$$\begin{aligned}
& \log p(c_{1:n}) \\
& \geq \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log p(c_{1:n} | z^s, z_{1:n}^r)] \\
& - \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log q(z^s, z_{1:n}^r) - \log q(z^s)q(z_{1:n}^r)] \\
& - \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log q(z^s)q(z_{1:n}^r) - \log p(z^s, z_{1:n}^r)] \\
& = \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log p(c_{1:n} | z^s, z_{1:n}^r)] \\
& - MI_q(z^s; z_{1:n}^r) \\
& - \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log q(z^s)q(z_{1:n}^r) - \log p(z^s, z_{1:n}^r)]
\end{aligned}$$

Due to the prior assumption  $p(z^s, z_{1:n}^r) = p(z^s)p(z_{1:n}^r)$ , we now have:

$$\begin{aligned}
& \log p(c_{1:n}) \\
& \geq \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log p(c_{1:n} | z^s, z_{1:n}^r)] \\
& - MI_q(z^s; z_{1:n}^r) \\
& - \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log q(z^s)q(z_{1:n}^r) - \log p(z^s)p(z_{1:n}^r)] \\
& = \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log p(c_{1:n} | z^s, z_{1:n}^r)] \\
& - MI_q(z^s; z_{1:n}^r) \\
& - \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log q(z^s) - \log p(z^s) + \log q(z_{1:n}^r) - \log p(z_{1:n}^r)] \\
& = \underbrace{\mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log p(c_{1:n} | z^s, z_{1:n}^r)]}_{1^{st} \text{ term}} \\
& - \underbrace{MI_q(z^s; z_{1:n}^r)}_{6^{st} \text{ term}} \\
& - \left[ \underbrace{\mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log \frac{q(z^s)}{p(z^s)}]}_A + \underbrace{\mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} [\log \frac{q(z_{1:n}^r)}{p(z_{1:n}^r)}]}_B \right]
\end{aligned}$$

For part  $A$ , we have:

$$\begin{aligned}
A &: \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} \left[ \log \frac{q(z^s)}{p(z^s)} \right] \\
&= -\mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} \left[ \log \frac{q(z^s | c_{1:n})}{p(z^s)} - \log \frac{q(z^s | c_{1:n})}{q(z^s)} \right] \\
&= - \left[ \underbrace{KL(q(z^s | c_{1:n}) || p(z^s))}_{2^{st} \text{ term}} - \underbrace{MI_q(z^s, c_{1:n})}_{4^{st} \text{ term}} \right]
\end{aligned}$$

For part  $B$ , we have:

$$\begin{aligned}
B &: \mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} \left[ \log \frac{q(z_{1:n}^r)}{p(z_{1:n}^r)} \right] \\
&= -\mathbb{E}_{q(z^s, z_{1:n}^r | c_{1:n})} \left[ \log \frac{q(z_{1:n}^r | c_{1:n})}{p(z_{1:n}^r)} - \log \frac{q(z_{1:n}^r | c_{1:n})}{q(z_{1:n}^r)} \right] \\
&= - \left[ \underbrace{KL(q(z_{1:n}^r | c_{1:n}) || p(z_{1:n}^r))}_{3^{st} \text{ term}} - \underbrace{MI_q(z_{1:n}^r, c_{1:n})}_{5^{st} \text{ term}} \right]
\end{aligned}$$