

HKUST SPD - INSTITUTIONAL REPOSITORY

Title Unified route representation learning for multi-modal transportation recommendation with spatiotemporal pre-training

Authors Liu, Hao; Han, Jindong; Fu, Yanjie; Li, Yanyan; Chen, Kai; Xiong, Hui

Source The VLDB Journal, 27 May 2022

Version Accepted Version

DOI 10.1007/s00778-022-00748-y

Publisher Springer

Copyright This version of the article has been accepted for publication, after peer review (when applicable) and is subject to Springer Nature's AM terms of use, but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at:
<http://dx.doi.org/10.1007/s00778-022-00748-y>

This version is available at HKUST SPD - Institutional Repository (<https://repository.ust.hk/ir>)

If it is the author's pre-published version, changes introduced as a result of publishing processes such as copy-editing and formatting may not be reflected in this document. For a definitive version of this work, please refer to the published version.

Unified Route Representation Learning for Multi-Modal Transportation Recommendation with Spatiotemporal Pre-Training

Hao Liu · Jindong Han · Yanjie Fu · Yanyan Li · Kai Chen · Hui Xiong

Received: date / Accepted: date

Abstract Multi-modal transportation recommendation aims to provide the most appropriate travel route with various transportation modes according to certain criteria. After analyzing large-scale navigation data, we find that route representations exhibit two patterns: spatio-temporal autocorrelations within transportation networks and the semantic coherence of route sequences. However, there are few studies that consider both patterns when developing multi-modal transportation systems. To this end, in this paper, we study multi-modal transportation recommendation with unified route representation learning by exploiting both spatio-temporal dependencies in transportation networks and the semantic co-

herence of historical routes. Specifically, we first transform the multi-modal transportation network into time-dependent multi-view transportation graphs and devise a graph-based contextual encoder to impute the missing traffic condition in transportation networks by leveraging various contextual factors. Then we propose a hierarchical multi-task route representation learning (HMTRL) framework for recommendations, including (1) a spatiotemporal graph neural network module to capture the spatial and temporal autocorrelation, (2) a coherent-aware attentive route representation learning module to explicitly model route coherence from historical routes, and (3) a hierarchical multi-task learning module to differentiate route representations for different transport modes by incorporating multiple auxiliary tasks equipped in different network layers. Moreover, to improve the model generalization capability, we further propose spatiotemporal pre-training strategies to exploit rich self-supervision signals hidden in transportation networks and historical trajectories. Finally, extensive experimental results on two large-scale real-world datasets demonstrate the effectiveness of the proposed system against eight baselines.

Keywords Multi-modal transportation · Route representation · Recommendation system · Hierarchical multi-task learning · Self-supervised learning

Hao Liu

The Thrust of Artificial Intelligence, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China and the Department of Computer Science and Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR, China.
E-mail: liuh@ust.hk

Jindong Han

The Thrust of Artificial Intelligence, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China
E-mail: jhanao@connect.ust.hk

Yanjie Fu

University of Central Florida, USA
E-mail: yanjie.fu@ucf.edu

Yanyan Li

Baidu Research, Beijing, China
E-mail: liyanyanliyanyan@baidu.com

Kai Chen

The Department of Computer Science and Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR, China.
E-mail: kaichen@cse.ust.hk

Hui Xiong

The Thrust of Artificial Intelligence, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China and the Department of Computer Science and Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR, China.
E-mail: xionghui@ust.hk

1 Introduction

The increasing prevalence of various transport modes (*e.g.*, car, bus, shared-bike, ride-sharing, etc.) and the rapidly expanding transportation networks (*e.g.*, road network, bus network, pedestrian network, etc.) have provided overwhelming alternatives for travelers to reach a destination. In recent years, multi-modal transportation recommendation has become an emerging routing service in many navigation and ride-hailing applications, such as Baidu Maps [1], Here [2],

and Didi Chuxing [3]. The target of multi-modal transportation recommendation is to help users find the most appropriate route from one place to another, by jointly considering one or more transport modes on a constrained transportation network. Therefore, accurate and intelligent multi-modal transportation recommendation can significantly help reduce the traveler’s decision cost and ultimately improve the user experience.

Existing studies on multi-modal transportation recommendation mainly fall into two categories. (1) *Searching based multi-modal route recommendation* aims to retrieve the shortest path on the transportation network, with a pre-defined distance metric (e.g., geographical distance, travel time, etc.). Most methods in this category [4,5] focus on extending graph search algorithms (e.g., Dijkstra’s algorithm, Bellman–Ford and contraction hierarchies [6]) to the multi-modal transportation network [7]. Such approaches are highly dependent on the pre-defined metric and overlook latent factors hidden in the data (e.g., mode and route preferences under different situational contexts [8]). (2) *Learning based transport mode recommendation* has partially addressed the problem by inferring coarse-grained transport mode preferences based on supervised or unsupervised machine learning techniques. A common routine in such methods [8] is to explicitly extract features (e.g., distance, estimated time of arrival (ETA)) from user historical data, such as GPS trajectories [9] and in-app clicks [10]. Such methods make recommendations based on empirically defined features, thus highly rely on the comprehensiveness of feature engineering. More recent studies have applied deep learning [11] and network embedding [12] for transport mode recommendation. However, such methods focus on learning coarse-grained vertex representations (e.g., origin-destination and user pair) or forecasting future travel costs (e.g., ETA), and are not capable of route-specific multi-modal transportation recommendation.

Indeed, the recent emergence of representation learning and multi-task learning techniques provides great potentials to overcome the above limitations. In this paper, we investigate the multi-modal transportation recommendation problem via the unified multi-task route representation learning, by exploiting both spatiotemporal dependencies from transportation networks and the semantic coherence from historical routes. However, three non-trivial challenges arise in achieving this goal. (1) **Spatiotemporal autocorrelation.** The multi-modal transportation network of various transport modes can be abstracted as a dynamic graph (e.g., bus line maybe created or removed, traffic condition is time varying). The dynamic graph contains rich structural and contextual information in both vertices (e.g., the degree, if it has a traffic light) and edges (e.g., distance, ETA, average speed). The first challenge is how to capture the spatial and temporal autocorrelation in the dynamic transportation

network. (2) **Route coherence representation.** After studying many routes traveled by users, we identify another important dependency in route representation learning, which we call route coherence. We analogize a route with a sentence, where each hub (e.g., road intersections, bus stations, etc.) and link (e.g., road segments, bus lines, etc.) correspond to a word. In this way, the representation of each hub and link in the route should not only depend on the transportation network but also semantically consistent with the whole route. Besides, the route sequence is of arbitrary length, and the importance of each vertex may vary. How to learn fix-length route representations by incorporating semantic information in historical routes is another challenge. (3) **Transport mode differentiation.** In the real-world, a route may be shared or partially shared by various transport modes. For example, given a bicycle route planned by navigation apps, it is with a high probability we can also travel by walk, and vice versa. The last challenge is how to differentiate the unified representation for various transport modes for recommendations.

To tackle the above challenges, we did some preliminary work [13], we proposed a *Hierarchical Multi-Task Route representation Learning* (HMTRL) framework for multi-modal transportation recommendations. Specifically, we first discretize the multi-modal transportation network into a set of graph snapshots over time and construct multi-view graphs, including (1) the *hub-centric graph* which regards transportation hubs as vertices, and (2) the *link-centric graph* which regards transportation links as vertices. After that, we propose the *spatiotemporal graph neural network* module which includes a graph convolution network layer that captures the non-linear spatial autocorrelation from multi-view graphs and a recurrent neural network (RNN) layer that captures the temporal autocorrelation across multiple graph snapshots. Moreover, a *coherent-aware attentive route representation learning* module is introduced, including (1) a bi-directional RNN layer that integrates the relatedness of historical routes into the representation of hubs and links, and (2) a self-attentive layer that projects the route sequence into a fixed-length representation with explicit quantifying the contribution of each hub and link. Finally, we propose the *hierarchical multi-task learning* module to learn mode-specific representations, and equip multiple correlated auxiliary tasks in different network layers to guide the optimization of representations for final recommendations. By incorporating structural dependencies in multi-view transportation graphs and route coherence in historical routes under various supervision signals, the mode-specific route representation enables more accurate route-level multi-modal transportation comparison and recommendation. Extensive experiments on two large-scale real-world datasets from one of the world’s largest navigation apps demonstrate HMTRL achieves the best performance compared with eight baselines.

In this paper, we further improve our HMTRL framework, based on the following three observations. First, the traffic conditions (e.g., traffic speed, ETA) in each time period are highly sparse, which may degrade the recommendation performance. Take Beijing for example, in each hour, only 66% road segments are covered by at least one trajectory at day time, and this coverage ratio is less than 37% at night. The missing traffic conditions may induce biased inputs and lead to unsatisfied recommendation results. Second, there are rich structural and contextual correlations between different hubs and links in transportation networks, which can be utilized to improve the recommendation performance. Moreover, only a small portion of historical trajectories are with transport mode labels, which may result in over-fitting problem on labeled data. In this paper, we further propose HMTRL⁺, which extends HMTRL for more accurate and generalizable multi-modal transportation recommendations. Comparing with [13], we further made the following four major contributions:

- We devise a dedicated graph-based contextual encoder to impute missing traffic conditions throughout a city. The contextual auto-encoder alleviates the data scarcity problem by jointly incorporating historical traffic patterns and partial real-time traffic conditions via a graph message passing component.
- We introduce spatiotemporal pre-training strategies to exploit various self-supervised signals for multi-modal transportation recommendation generalization. In particular, attribute prediction tasks exploit rich semantic information in transportation networks and the trajectory contrastive learning task extracts transferable knowledge for robust route representation.
- We provide a systematic complexity analysis of each component of HMTRL⁺.
- We evaluate the effectiveness and efficiency of the proposed method on two large-scale real-world datasets. The results show HMTRL⁺ achieves the best recommendation performance compared with HMTRL as well as eight baselines.

2 Preliminaries

2.1 Definitions and Problem Statement

Consider a set of transport modes $\mathcal{M} = \{m_1, m_2, \dots, m_k\}$, where each mode corresponds to a transportation network (e.g., road network, bus line network) that supports vehicle or pedestrian movement. Generally, the transportation network of each transport mode is composed of a set of *hubs* (e.g., road intersection, bus or metro station) and a set of *links* (e.g., road segment, bus line). We formally define the multi-modal transportation network based on transportation networks of each transport mode.

Definition 1 Multi-Modal Transportation Network (MMTN).

The multi-modal transportation network integrates multiple mode-specific transportation networks into a unified attributed directed graph $\mathcal{G} = (V, E, A^V, A^E, M)$, where V is the set of hubs, E is the set of links, M is a mapping function indicates the supported transport modes of each hub and link, A^V and A^E are respectively hub and link features, such as number of bus lines across the hub, spherical distance of the road segment, and ETA of the bus line.

We use $M(v_i)$ and $M(e_{ij})$ to denote the supported transport modes of each hub $v_i \in V$ and $e_{ij} \in E$. Note each hub and link may support more than one transport modes (e.g., walk and bicycle). We say two hubs are *adjacent* to each other if and only if there is a link connecting them, two links are *adjacent* if and only if a user can transfer from one to another by one hub. Without loss of generality, we constraint a user can only transfer to other links or transport modes in a hub, and a link is the smallest movement unit in the transportation network, e.g., a road segment between two adjacent road intersections, and a bus line between two adjacent bus stations.

Definition 2 Route. *A route is a triplet $r_i = \langle H, L, \phi \rangle$, where H is a sequence of adjacent hubs, L is a sequence of adjacent links, and ϕ is a mapping function that indicates the corresponding transport mode of each hub and link in the route.*

Different from the mapping function M in \mathcal{G} , $\phi(v_i)$ and $\phi(e_{ij})$ identify the unique transport mode in the corresponding route. In this work, we restrict a route start and terminate at a hub, a route may consist of one or more transport modes.

Definition 3 Routing Query. *A routing query is defined as a triplet $q = \langle o, d, t \rangle$, where o and d are origin and destination locations represented by a pair of longitude and latitude, and t is the departure time.*

Since the origin o and the destination d are arbitrary locations, we project them to nearby hubs for recommendation. We say a route r_i is *feasible* for q if the route start from o and terminate at d .

Problem 1 Multi-Modal Transportation Recommendation.

Given a MMTN \mathcal{G} , a routing query q and a set of feasible routes Γ for q , our problem is to recommend the most appropriate route $r_i \in \Gamma$ based on the conditional probability $\hat{y}_i \leftarrow \mathcal{F}(r_i|q, \Gamma, \mathcal{G})$, where \mathcal{F} is the unified mapping function we aim to learn.

To reduce the computational complexity, we derive the route candidate set Γ (typically less than 20 candidates) based on existing routing engines [14, 8]. To guarantee the utility of recommendations, we restrict the maximum number of mode transfer in each route candidate to three.

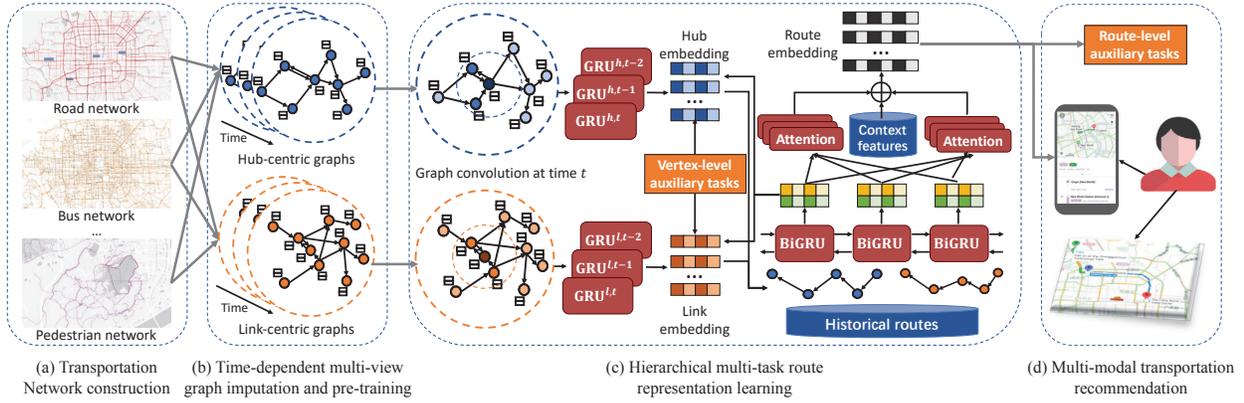


Fig. 1 An overview of unified route representation learning for multi-modal transportation recommendation.

2.2 Framework Overview

Figure 1 shows an overview of our approach, where the inputs are the multi-modal transportation network, historical routes and its corresponding routing query, and context features such as weather condition; the output is the recommended route. Overall, there are five tasks in our approach, (1) the construction of time-dependent multi-view transportation graphs, (2) the imputation of missing traffic conditions, (3) the unified route representation learning, (4) the hierarchical multi-task learning for mode-specific representation generation and route recommendation, and, (5) the spatiotemporal pre-training for model generalization. To be specific, in the first task, we transform the multi-modal transportation network to a set of time-dependent multi-view graphs from both the hub-centric perspective and the link-centric perspective. In the second task, we complement missing traffic conditions based on MMTN via a graph-based contextual encoder. In the third task, the unified route representation is obtained via (1) the joint spatiotemporal autocorrelation modeling of the hub-centric graph and the link-centric graph, and (2) the coherent-aware attentive route representation learning by exploiting historical routes. In the fourth task, we differentiate representations for various transport modes via multiple implicit tasks and boost the recommendation performance by incorporating multiple related auxiliary tasks in different neural network layers. In the final task, we improve the robustness of the recommendation system by exploiting rich self-supervision signals in MMTN and historical trajectories.

3 Constructing Time-Dependent Multi-View Transportation Graphs

First of all, we construct time-dependent multi-view transportation graphs to characterize dynamic structural and contextual information in MMTN.

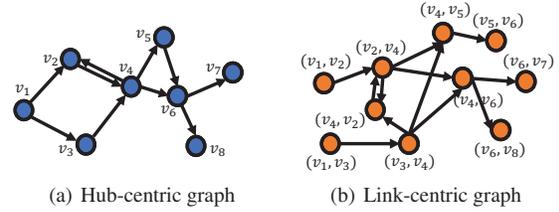


Fig. 2 An example of multi-view transportation graphs.

The hub-centric view. The hub-centric graph is a direct mapping of MMTN, where vertices and edges are respectively transportation hubs and links. Specifically, we first discretize the time-evolving graph into a sequence of snapshots, denoted by $\mathcal{G}^h = [\mathcal{G}^{h,t_1}, \mathcal{G}^{h,t_2}, \dots, \mathcal{G}^{h,t_n}]$, where \mathcal{G}^{h,t_i} is the hub-centric graph at time t_i . Figure 2(a) gives an illustrative example snapshot of time-dependent hub-centric graph. For each vertex in the graph, we attach corresponding hub features, including both time-invariant features (*e.g.*, degree, if have a traffic light) and dynamic features (*e.g.*, traffic volume). For two adjacent hubs v_i and v_j , we construct the corresponding adjacency weight based on a Gaussian kernel [15],

$$c_{ij}^h = \exp\left(-\frac{\text{dist}(v_i, v_j)^2}{\delta^2}\right), \quad (1)$$

where $\text{dist}(v_i, v_j)$ denote the spherical distance [16] between v_i and v_j , and δ is the standard deviation of spherical distances. In consequence, c_{ij}^h demonstrates the geographical distance distribution among adjacent hubs.

The link-centric view. The link-centric graph flips vertices and edges in MMTN to preserve structural and contextual information in transportation links. Similar to the hub-centric graph, we discretize the time-evolving link-centric graph into a sequence of snapshots, which is denoted by $\mathcal{G}^l = [\mathcal{G}^{l,t_1}, \mathcal{G}^{l,t_2}, \dots, \mathcal{G}^{l,t_n}]$, where \mathcal{G}^{l,t_i} indicates the link-centric graph at time t_i . Figure 2(b) shows an illustrative

example of the time-dependent link-centric graph at a specific time slice. For each vertex (*i.e.*, link) in the graph, we attach corresponding time-invariant features (*e.g.*, distance, road level) and dynamic features (*e.g.*, average speed, ETA). Consider two links $e_i = (v_1, v_2)$ and $e_j = (v_3, v_4)$, we set the adjacency constraint as

$$c_{ij}^l = \begin{cases} 1, & v_2 = v_3 \text{ and } i \neq j \\ 0, & \text{otherwise} \end{cases}. \quad (2)$$

There is a directed edge from e_i to e_j if and only if $v_2 = v_3$ in the corresponding MMTN. In other word, an edge in the link-centric graph forms a 2-hop route in the MMTN. Different with the hub-centric graph, $c_{ij}^l \in \{0, 1\}$ preserves the connectivity information.

In the same time slice, the link-centric graph is an edge-to-vertex dual of the hub-centric graph. We can re-construct the MMTN from either the hub-centric graph or the link-centric graph [17]. The time-dependent multi-view graphs therefore preserves the temporal dynamics and the structural integrity of MMTN for subsequent graph representation learning.

4 Graph-based Missing Traffic Condition Inference

In urban transportation networks, the real-time traffic condition is spatiotemporally sparse due to the unbalanced distribution of vehicle trajectories. Imputing the missing real-time traffic conditions for hubs and links can further improve the recommendation performance by providing more side information for the recommendation system. However, simply imputing missing traffic conditions by default values or averaging neighbor traffic conditions may introduce biased and noisy inputs and lead to undesired performance degradation. Recent years we have witnessed the powerful capability of graph learning for handling data sparsity in trajectory data mining [18–22]. To this end, in this work, we propose to infer missing traffic conditions via a graph-based contextual encoder, based on various contextual factors.

Let $\mathbf{D} \in \mathbb{R}^{N \times K}$ be a traffic condition matrix consisting of N road segments and K dynamic traffic condition features (*e.g.*, traffic speed and ETA) at time slot t . We define an indication matrix $\mathbf{M} \in \{0, 1\}^{N \times K}$, where $\mathbf{M}[n, k] = 0$ if the traffic condition of $\mathbf{D}[n, k]$ is missing and $\mathbf{M}[n, k] = 1$ otherwise. We formulate the traffic condition inference as a regression task on transportation graphs, where the goal is to estimate the missing values $\mathbf{D}[n, k]$ where $\mathbf{M}[n, k] = 0$.

Contextual feature construction. Intuitively, the traffic condition is correlated with various urban factors, which can be leveraged to improve the real-time traffic condition imputation accuracy. Specifically, consider a hub-centric/link-centric transportation graph \mathcal{G}^t at time t , we construct four

categories of contextual features for traffic condition imputation. First, we associate each node on the graph with a set of attribute features \mathbf{X}^s (*e.g.*, road level, speed limit), to encode rich physical information of links and hubs. Second, we extract temporal features \mathbf{X}^e to quantify the temporal patterns of traffic patterns, such as time of day, day of week, holidays. Moreover, as the current traffic condition can be affected by previous traffic patterns, we extract historical traffic conditions $\mathbf{X}^p \in \mathbb{R}^{N \times K \times T}$ from previous T time steps as input features. In addition, to handle the missing data, we simply replace the missing values in \mathbf{X}^p by using the average traffic conditions calculated based on corresponding hubs/links and time slot. Finally, as the traffic condition shows strong periodicity [23], we further construct a long-term traffic condition matrix $\mathbf{X}^l \in \mathbb{R}^{N \times K}$ to augment dynamic traffic information. In particular, we calculate the averaged traffic condition at time slot t in a day based on historical trajectory data in a longer time period (*e.g.*, 3 months). The accumulated long-term traffic data in a similar time slot can provide additional hints for current traffic conditions and can be adopted to alleviate the data sparsity problem.

Graph-based inference with node dropout. Besides contextual features, the available real-time traffic condition of neighbor nodes can also provide strong contextual signals for missing traffic condition imputation. For instance, a traffic jam on a road segment can propagate to adjacent road segments and lead to subsequent traffic congestion. We construct a graph-based encoder with node dropout to harness the information hidden in the partially observed real-time traffic conditions.

Let \mathbf{X}^u denotes the combination of four categories of contextual features and traffic condition matrix

$$\mathbf{X}^u = \mathbf{X}^s \parallel \mathbf{X}^p \parallel \mathbf{X}^l \parallel \mathbf{X}^e \parallel \mathbf{D}, \quad (3)$$

where \parallel is the concatenation operation. We capture the spatial correlation between different nodes via the graph convolution operation

$$\mathbf{x}'_i = \sigma\left(\sum_{j \in \mathcal{N}_i} c_{ij} \mathbf{W}_o \mathbf{x}_j^u\right) \parallel \mathbf{x}_i^u, \quad (4)$$

where \mathbf{x}_j^u is the features of node j , \mathbf{x}'_i is the updated vertex representation, σ is a non-linear activation function, c_{ij} is the adjacency weight, $\mathbf{W}_o \in \mathbb{R}^{d \times d}$ is learnable weighted matrix shared by all vertices in \mathcal{G}^t , and \mathcal{N}_i is the set of neighboring vertices of v_i in \mathcal{G}^t . We stack l layers to model long-range spatial dependencies. Finally, we use a linear regression model to estimate the missing traffic conditions.

To fully utilize the information hidden in observed real-time traffic conditions, we devise a node dropout strategy for the imputation task. Specifically, in each iteration, we randomly mask a set of nodes on the transportation graph

Algorithm 1: Missing Traffic Condition Imputation

Input: time slot t , traffic condition matrix \mathbf{D} , indication matrix \mathbf{M} , transportation graph \mathcal{G}^t , number of layers L , node dropout rate p , neighborhood function \mathcal{N} , activation function σ

Output: output estimated traffic condition matrix $\hat{\mathbf{D}}$

- 1 $\mathbf{X}^s, \mathbf{X}^p, \mathbf{X}^l, \mathbf{X}^e \leftarrow \text{FeatureAugmentation}(t, \mathcal{G}^t)$;
- 2 **if** $\text{training} == \text{True}$ **then**
- 3 | $\mathbf{D}' \leftarrow \text{NodeDrop}(\mathbf{D}, \mathbf{M}, p)$
- 4 **else**
- 5 | $\mathbf{D}' = \mathbf{D} \odot \mathbf{M}$
- 6 $\mathbf{X}^u \leftarrow \mathbf{X}^s \parallel \mathbf{X}^p \parallel \mathbf{X}^l \parallel \mathbf{X}^e \parallel \mathbf{D}'$;
- 7 $\mathbf{X}^0 \leftarrow \mathbf{X}^u$;
- 8 **for** $l = 1$ **to** L **do**
- 9 | **for** $\mathbf{x}_i^{l-1} \in \mathbf{X}^{l-1}$ **do**
- 10 | | $\mathbf{x}_i^l = \sigma((\sum_{j \in \mathcal{N}_i} c_{ij} \mathbf{W}_s \mathbf{x}_j^{l-1}) \parallel \mathbf{x}_i^{l-1})$
- 11 $\hat{\mathbf{D}} \leftarrow \text{LinearRegression}(\mathbf{X}^{l-1})$;
- 12 **return** $\hat{\mathbf{D}}$

with a pre-defined dropout rate p . We construct a mask matrix \mathbf{M}_{drop} , where each entry $\mathbf{M}_{drop}[n, k]$ is set to 1 with probability p and 0 otherwise. We update the indication matrix \mathbf{M}' as follows

$$\mathbf{M}' = \mathbf{M} \odot \mathbf{M}_{drop}, \quad (5)$$

Then, we randomly mask nodes with real-time traffic conditions by using node dropout,

$$\mathbf{D}' = \mathbf{D} \odot \mathbf{M}'. \quad (6)$$

In the training phase, we use the masked traffic conditions \mathbf{D}' as features in (3), and predict the rest traffic conditions.

The graph based contextual encoder aims to optimize the mean square error between the ground-truth traffic conditions and predictions

$$\mathcal{L}_m = \frac{1}{N_m} \sum_{n=1}^N \sum_{k=1}^K (\mathbf{M}[n, k] - \mathbf{M}'[n, k]) (\hat{\mathbf{D}}[n, k] - \mathbf{D}[n, k])^2. \quad (7)$$

where $\hat{\mathbf{D}}[n, k]$ is the estimated traffic condition using the graph based contextual encoder, N_m denotes the number of non-zeros entries in $(\mathbf{M} - \mathbf{M}')$. Note we impute various traffic conditions, *e.g.*, traffic speed and traffic volume. The traffic condition imputation task follows the multi-task learning paradigm with hard parameter sharing. Overall, the complete computation of missing traffic condition inference is summarized in Algorithm 1.

5 Hierarchical Multi-Task Route Representation Learning

Based on time-dependent multi-view transportation graphs, we obtain mode-specific route representation and make multi-

modal transportation recommendations with the following intuitions.

Intuition 1: Graph autocorrelation preservation. The time-dependent multi-view transportation graphs contain rich structural and contextual information that varies over time. The vertex (*i.e.*, hub and link) in each graph at different time slices is both spatially and temporally autocorrelated with other vertices, and can contribute to the overall route representation. Therefore, the model should be able to collaboratively learn the spatial and temporal autocorrelation for both hubs and links.

Intuition 2: Coherent-aware route representation. The routes are arbitrary-length sequences and different hubs and links are playing different important roles in different routes. The fixed-length route representation should pay different attention to hubs and links in the sequence to distill salient features of each route. Besides, each hub and link is semantically coherent with its historical routes. Therefore, the representation of hubs and links should reflect a higher relevance with historical routes it involved in.

Intuition 3: Multi-modal route representation differentiation. A transportation hub or link may be shared by various transport modes. Correlating tasks such as link ETA prediction and route preference inference offer potential auxiliary signals to help differentiate mode-specific representations. In consequence, the proposed method should be capable of integrating various auxiliary tasks in different granularity (*e.g.*, vertex level and route level) for route representation differentiation and recommendation.

5.1 Spatiotemporal Autocorrelation Modeling

We first introduce the *spatiotemporal graph neural network* module to capture spatiotemporal autocorrelations based on time-dependent multi-view transportation graphs.

Modeling spatial autocorrelation. We employ Graph Neural Network (GNN) [24] to capture spatial autocorrelation at each time step. By iteratively aggregating and transforming neighbor representations [25], GNN obtains locally smoothed representations where spatially adjacent hubs and links tend to be close in the latent space.

Formally, consider a transportation graph \mathcal{G}^t at time t , let \mathbf{x}_i denotes the d -dimensional representation of vertex $v_i \in \mathcal{G}^t$, we define the graph convolution operation ($GConv$) as

$$\mathbf{x}'_i = GConv(\mathbf{x}_i) = \sigma((\sum_{j \in \mathcal{N}_i} c_{ij} \mathbf{W}_s \mathbf{x}_j) \parallel \mathbf{x}_i), \quad (8)$$

where \mathbf{x}'_i is the updated vertex representation, σ is a non-linear activation function, c_{ij} is the adjacency weight, $\mathbf{W}_s \in \mathcal{R}^{d \times d}$ is learnable weighted matrix shared by all vertices in \mathcal{G}^t , \parallel is the concatenation operation, and \mathcal{N}_i is the set of neighbor vertices of v_i in \mathcal{G}^t . Note that we can repeat l times

graph convolution operations to capture l -hop spatial dependencies. We update representations of $v_i^{p,t} \in \mathcal{G}^{p,t}$ and $v_j^{l,t} \in \mathcal{G}^{l,t}$ by $\mathbf{x}_i^{p,t} = GConv(\mathbf{x}_i^{p,t})$ and $\mathbf{x}_j^{l,t} = GConv(\mathbf{x}_j^{l,t})$, respectively.

Modeling temporal autocorrelation. The representations of hubs and links are not only correlated with neighboring vertices in \mathcal{G}^p and \mathcal{G}^l , but also influenced by their status in previous time periods. We extend GNN by Gated Recurrent Unit (GRU) [26], a simple yet effective variant of RNN, for temporal autocorrelation modeling. Consider representations of v_i in previous T steps $(\mathbf{x}_i^{t-T}, \mathbf{x}_i^{t-T+1}, \dots, \mathbf{x}_i^t)$, where \mathbf{x}_i^t is the output of the graph convolution operation at time t . We denote the status of v_i at time step $t-1$ and t as \mathbf{h}_i^{t-1} and \mathbf{h}_i^t , respectively. The GRU operation is defined as

$$\mathbf{h}_i^t = GRU(\mathbf{h}_i^{t-1}, \mathbf{x}_i^t) = (1 - \mathbf{z}_i^t) \circ \mathbf{h}_i^{t-1} + \mathbf{z}_i^t \circ \tilde{\mathbf{h}}_i^t, \quad (9)$$

where \mathbf{z}_i^t and $\tilde{\mathbf{h}}_i^t$ are defined as

$$\begin{cases} \mathbf{r}_i^t = \sigma(\mathbf{W}_r[\mathbf{h}_i^{t-1} \parallel \mathbf{x}_i^t] + \mathbf{b}_r) \\ \mathbf{z}_i^t = \sigma(\mathbf{W}_z[\mathbf{h}_i^{t-1} \parallel \mathbf{x}_i^t] + \mathbf{b}_z) \\ \tilde{\mathbf{h}}_i^t = \tanh(\mathbf{W}_{\tilde{h}}[\mathbf{r}_i^t \circ \mathbf{h}_i^{t-1} \parallel \mathbf{x}_i^t] + \mathbf{b}_{\tilde{h}}) \end{cases}, \quad (10)$$

where $\mathbf{W}_r, \mathbf{W}_z, \mathbf{W}_{\tilde{h}}, \mathbf{b}_r, \mathbf{b}_z, \mathbf{b}_{\tilde{h}}$ are learnable parameters, \parallel is the concatenation operation, and \circ denotes Hadamard product. The hidden state \mathbf{h}_i^t reflects both the spatial and temporal autocorrelation of vertex v_i in corresponding time-dependent graphs. For each hub and link in corresponding views, we respectively derive $\mathbf{h}_i^{h,t} = GRU(\mathbf{h}_i^{h,t-1}, \mathbf{x}_i^{h,t})$ and $\mathbf{h}_i^{l,t} = GRU(\mathbf{h}_i^{l,t-1}, \mathbf{x}_i^{l,t})$ for route representation learning.

5.2 Route Representation Learning

Then we present the *coherent-aware attentive route representation learning* module, including (1) the Bi-directional RNN based route coherence modeling block first incorporates route coherence constraints in historical routes to hub and link representations, and (2) the self-attentive route representation learning block further projects arbitrary-length routes into fixed-length representation vectors by automatically learning the importance of each hub and link in the corresponding route.

Bi-directional RNN based route coherence modeling.

The insight of route coherence modeling is to incorporate the relatedness of prefix and suffix sub-routes into the current hub and link representations. Figure 3(a) shows an illustrative example of the prefix sub-route coherence on a road network, where the orange arrows form a historical route traveled by a user, the yellow arrow is the current link, and the green arrows are candidate links. Given a historical route $[e_1, e_4, e_9, e_{12}]$, consider e_9 as the current link, there is another candidate link e_7 for prefix sub-route $[e_1, e_4]$ and a

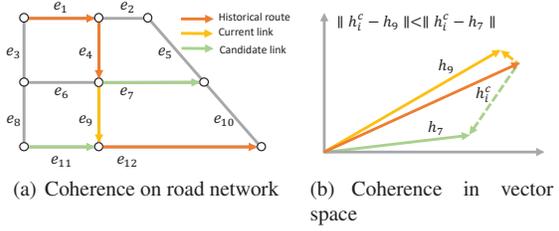


Fig. 3 An illustrative example of route coherence modeling.

candidate link e_{11} for suffix sub-route $[e_{12}]$. Based on the historical route, e_9 is more relevant with the prefix sub-route $[e_1, e_4]$ and the suffix sub-route $[e_{12}]$. Therefore, the representation of e_9 should reflect not only graph dynamics in the MMTN, but also the historical route dependency. We adopt the *Bi-directional GRU (BiGRU)* operation to integrate the route coherence dependency into both hub and link representations from both forward direction and backward direction.

Specifically, we reuse the GRU operation in Equation (9) for hub and link representation update. Formally, for vertex v_i , consider its prefix sub-route $[\dots, v_{i-2}, v_{i-1}]$ and suffix sub-route $[v_{i+1}, v_{i+2}, \dots]$, we obtain the forward coherent dependency and backward coherent dependency of v_i by $\vec{\mathbf{h}}_i^c = GRU(\vec{\mathbf{h}}_{i-1}, \mathbf{x}_i)$ and $\overleftarrow{\mathbf{h}}_i^c = GRU(\overleftarrow{\mathbf{h}}_{i+1}, \mathbf{x}_i)$, and define the *BiGRU* operation as

$$\mathbf{h}_i^c = BiGRU(\vec{\mathbf{h}}_i^c, \overleftarrow{\mathbf{h}}_i^c) = \mathbf{W}_c[\vec{\mathbf{h}}_i^c \parallel \overleftarrow{\mathbf{h}}_i^c], \quad (11)$$

where $\mathbf{W}_c \in \mathcal{R}^{2d \times d}$ is the learnable parameter projecting the concatenated representation to d -dimensional vector. Take Figure 3(b) for example, denote the representation of the historical route as \mathbf{h}_i^c , and the representations of e_7 and e_9 as $\mathbf{h}_7, \mathbf{h}_9$, the route coherence modeling forces $dist(\mathbf{h}_i^c, \mathbf{h}_9) < dist(\mathbf{h}_i^c, \mathbf{h}_7)$, where $dist(\cdot, \cdot)$ is a distance function in the latent vector space, e.g., the Euclidean Distance. Specifically, we aim to simultaneously minimize the distance between \mathbf{h}_i^c and \mathbf{h}_9 and maximize the distance between \mathbf{h}_i^c and \mathbf{h}_7 . Here we leverage *triplet loss* to achieve this goal, the details are elaborated in Section 5.4. The updated vertex representation incorporates both prefix and suffix sub-route information and is more informative for multi-modal transportation recommendations.

Self-attentive route representation learning. There are still two problems to obtain unified route representation learning: (1) the length of each route may vary, and (2) the importance of each hub and link in the route may be different. Simply averaging representations of hubs and links can not capture the diversified importance of each hub and link, while the RNN in Equation (11) suffers from the gradient vanishing problem [27]. Inspired by the recent success of the attention mechanism [28] on modeling weighted dependencies of long sentences. We analogize multi-modal routes as

sentences and employ a self-attention mechanism to transform arbitrary-length routes to fixed-length route representation vectors, with explicit quantifying the importance of both hubs and links in each route.

Given a hub or route sequence of n vertices, we devise K independent self-attentive operations to stabilize the learning process. Specifically, we define the k -th attentive score of v_i as

$$\alpha_{i,k} = \frac{\exp(\mathbf{W}_{a,k} \tanh(\mathbf{W}_{b,k} \mathbf{h}_i))}{\sum_{j=1}^n \exp(\mathbf{W}_{a,k} \tanh(\mathbf{W}_{b,k} \mathbf{h}_j))}, \quad (12)$$

where \mathbf{h}_i and \mathbf{h}_j are representations of v_i and v_j , $\mathbf{W}_{a,k}$ and $\mathbf{W}_{b,k}$ are learnable weights in the k -th attentive operation. Then, we derive the sequence representation by

$$\mathbf{h}' = \parallel_{k=1}^K \left(\sum_{i=1}^n \alpha_{i,k} \mathbf{W}_{r,k} \mathbf{h}_i \right), \quad (13)$$

where \parallel is the vector concatenation operation and $\mathbf{W}_{r,k} \in \mathcal{R}^{d \times d}$ is the learnable parameter corresponding to k -th self-attentive operation. Based on Equation (13), we derive the corresponding hub sequence representation $\mathbf{h}^{r,h}$ and link sequence representation $\mathbf{h}^{r,l}$, and derive the unified route representation as

$$\mathbf{h}^r = \mathbf{h}^{r,h} \parallel \mathbf{h}^{r,l}. \quad (14)$$

5.3 Hierarchical Multi-Task Learning

Finally we introduce the *hierarchical multi-task learning* module for multi-modal route representation differentiation and recommendation optimization. By jointly learning multiple related tasks, multi-task learning shares common knowledge in each task and, therefore, improves the generality of the model [29]. Incorporating auxiliary tasks in different granularity has been proved beneficial in many tasks such as document parsing and synonym prediction [30,31]. In HMTRL, we introduce various auxiliary tasks as complement supervision signals, where different tasks are equipped at different neural network layers.

Specifically, we employ the hard parameter sharing [32] in HMTRL, where different tasks are sharing part of the model but have individual output layers. In HMTRL, the learning tasks can be categorized into two classes, (1) the *Vertex-level MTL* that corresponds to representation learning of vertices (*i.e.*, hubs and links) in time-dependent multi-view graphs, (2) the *Route-level MTL* that corresponds to route representation optimization and recommendation.

Vertex-level MTL. Let $\{\mathcal{T}^{v,i}\}_{i=1}^{\tau_1}$ denote a set of auxiliary vertex tasks, where each task $\mathcal{T}^{v,i}$ corresponds to a set of labels $\{y_j^i\}_{j=1}^{n_i}$ if any, where $y_j^i \in \mathcal{R}$. We first introduce transport mode differentiation tasks to generate mode-specific representations for hubs and links. Specifically, for

each transport mode $m_i \in \mathcal{M}$, we define a corresponding task \mathcal{T}^{m_i} to obtain the mode-specific representation after the spatiotemporal graph neural network, $\mathbf{h}_j^{m_i} \leftarrow \mathcal{F}^{m_i}(\mathbf{h}_j)$, where \mathcal{F}^{m_i} is a mode-specific function implemented by a fully connected multi-layer neural network. Note that because not all hubs and links are feasible for all transport modes (*e.g.*, a bus link does not support walk), we mask infeasible transport modes in optimization. Different from other auxiliary tasks, transport mode differentiation tasks do not have direct supervision signals. Instead, all such tasks are optimized based on higher-level task feedbacks (*e.g.*, the link type classification, the route distance prediction and multi-modal transportation recommendation) via the back-propagation.

Thereafter, we extract various vertex attributes as supervision signals and facilitate multiple auxiliary task specific layers. Concretely, we integrate regression tasks including the distance prediction and forecasting ETA of the next time step, and integrate classification tasks including hub type (road intersection, bus station, etc.) and link type (road segment, bus line, etc.) inference.

Route-level MTL. Similarly, we define a set of auxiliary route tasks $\{\mathcal{T}^{r,i}\}_{i=1}^{\tau_2}$. Specifically, we first integrate the route coherence modeling task, by leveraging the intermediate prefix sub-route and suffix-route representation derived by Equation (11). Rather than set explicit labels, we optimize the representation in a self-supervised manner, to allow the vertex representation \mathbf{h}_i in the latent space to approximate the corresponding sub-route representation \mathbf{h}_i^c more closely. After that, we incorporate various route related regression tasks, including route distance prediction and future ETA prediction. Besides, for each route r_j , we facilitate the transport mode prediction task by applying a multi-class classifier.

Finally, we define the main recommendation task. Consider the route representation \mathbf{h}_i^r , we define the output layer as

$$\hat{y}_i = \sigma(\mathbf{w}_{main}[\mathbf{h}_i^r \parallel \mathbf{x}_i^{context}] + b_{main}), \quad (15)$$

where \hat{y}_i is the estimated travel likelihood of route r_i , σ is a non-linear activation function, \mathbf{w}_{main} are the learnable parameters of the main task, and b_{main} is the bias. Similar to [8], we also concatenate a context vector $\mathbf{x}_i^{context}$ to incorporate the situational context, including features such as weather condition and time periods. To facilitate the readability and reproducibility, we provide detailed settings of auxiliary tasks in the supplementary material.

5.4 Optimization

In HMTRL, we optimize both the main task as well as auxiliary tasks in different layers jointly. For the main task and

auxiliary classification tasks, we employ the *cross-entropy loss* for optimization. For regression tasks such as distance and ETA prediction, the objective is to minimize the *mean square error loss*. Please refer to supplementary material for detailed auxiliary losses. Additionally, we introduce the *triplet loss* for the optimization of route coherence,

$$\mathcal{L}_c = -\frac{1}{nk} \sum_{i=1}^n \sum_{j=1}^k \max\{(\|\mathbf{h}_i^c - \mathbf{h}_i\|_2 - \|\mathbf{h}_i^c - \mathbf{h}_j\|_2 + \gamma), 0\}, \quad (16)$$

where \mathbf{h}_j is the representation of the negative sample, γ is a margin constant between positive pair $(\mathbf{h}_i^c, \mathbf{h}_i)$ and negative pair $(\mathbf{h}_i^c, \mathbf{h}_j)$. We draw adjacent vertices $v_j \in \mathcal{N}(v_i)$ in the corresponding time-dependent transportation graph as negative samples, and force the representation of the vertex in the route \mathbf{h}_i is closer to the coherent state \mathbf{h}_i^c than negative samples \mathbf{h}_j .

Overall, we aim to optimize the following objective,

$$\mathcal{L} = \mathcal{L}_{main} + \beta_1 \sum_{i=1}^{\tau_1} \mathcal{L}_i^v + \beta_2 \sum_{i=1}^{\tau_2} \mathcal{L}_i^r, \quad (17)$$

where \mathcal{L}_i^v and \mathcal{L}_i^r are auxiliary vertex and route tasks, β_1 and β_2 are hyper-parameters control the importance of auxiliary tasks. We employ Adam optimizer [33] for training with an exponential decay.

6 Spatiotemporal Pre-Training

The multi-modal transportation network and historical trajectories contain rich contextual and functional information. Inspired by the success of deep neural network pre-training techniques for improving generalization capability in computer vision and nature language processing, we propose spatiotemporal pre-training to learn robust and transferable model parameters by introducing additional self-supervised signals. Previous works [34,35] have examined that self-supervised pre-trained models are more robust to adversarial examples and noises, and can be easily transfer to enhance the performance of unseen downstream tasks [36,35]. The key insight of our method is to exploit the transportation network and massive unlabeled historical trajectories for self-supervised representation learning. Specifically, in this work, our approach comprises two parts: masked attribute prediction and trajectory contrastive learning.

6.1 Masked Attribute Prediction

In masked attribute prediction, we aim to capture the semantic knowledge and attribute patterns underlying the multi-view transportation graphs. The intuition behind masked attribute prediction is to fully exploit easily obtained (*i.e.*, with

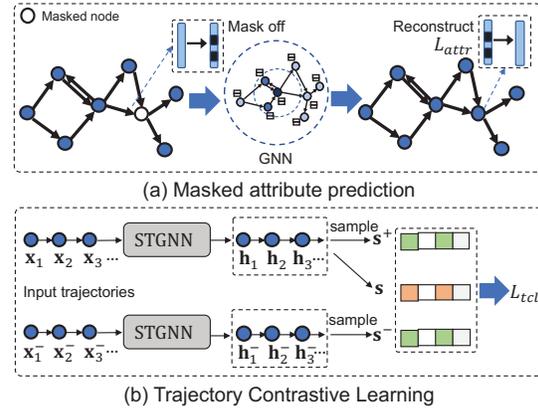


Fig. 4 Illustration of spatiotemporal pre-training tasks, including (1) the attribute prediction, and (2) the trajectory contrastive learning.

little extra labeling cost) attributes on the graph to learn better feature representations. For the transportation graph at time step t , we first mask target attributes and then predict the masked features based on the observed attributes and graph structure. Specifically, we randomly remove available attributes on the graph by replacing the masked features with default values (*e.g.*, zero). After that, we adopt the graph convolutional network introduced in Section 5.1 to generate the hidden representation \mathbf{x}'_i for node v_i in the corresponding graph. Finally, a linear model $f_{decoder}(\cdot)$ is devised to reconstruct the masked attributes of node i by taking \mathbf{x}'_i as input. In this task, we aim to minimize the Mean Square Error (MSE) between the predicted and masked attributes,

$$\mathcal{L}_{attr} = \|f_{decoder}(\mathbf{x}'_i) \odot \mathbf{m}_i - \mathbf{x}_i \odot \mathbf{m}_i\|_2^2 \quad (18)$$

where \mathbf{x}_i denotes input attributes of node v_i , \mathbf{m}_i is an indication vector, we denote $\mathbf{m}_i[j] = 1$ if the j -th input feature is masked and $\mathbf{m}_i[j] = 0$ otherwise.

In this way, we obtain a neural encoder to preserve latent characteristics and correlations of nodes in multi-view transportation graphs. The pre-trained network can be utilized as a better low-level neural network initialization of the HMTRL framework for downstream route recommendation task.

6.2 Trajectory Contrastive Learning

Prior works [37,38] have shown that GPS trajectories can reflect important functional patterns of transportation network, *i.e.*, functional distributions. The key insight of trajectory contrastive learning is to enforce hub and link representations in the same trajectory to have higher correlation with each other. In this way, we can capture intra-trajectories dependencies from massive unlabeled trajectories to improve the effectiveness of multi-modal route recommendation. In

this work, we achieve this goal by contrasting the trajectory and its local parts based on mutual information maximization [39].

Formally, consider an entire trajectory sequence $e_{1:N} = \{e_1, e_2, \dots, e_N\}$ and a randomly clipped sub-trajectory $e_{i:j} = \{e_i, \dots, e_j\}$ from $e_{1:N}$. We first employ map-matching algorithm [40] to project raw trajectories to corresponding time-ordered node (*i.e.*, hub and link) sequences $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, where each node is associated with a set of hub or link features. Similar to the masked attribute prediction task, we first adopt the graph convolution network to derive node representation \mathbf{h}_i for each node v_i . For a trajectory $e_{1:N}$, we can learn representations $\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N\}$ from scratch or reuse the learned representation from the masked attribute prediction task. To maximize the mutual information between the entire trajectory and the sub-trajectory, we aggregate sequence of node representations $\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N\}$ and $\{\mathbf{h}_i, \dots, \mathbf{h}_j\}$ into unified trajectory-level representations \mathbf{s} and \mathbf{s}_0 . For efficiency concern, we obtain trajectory representations via average pooling

$$\mathbf{s} = \text{MEAN}(\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N\}) = \sigma\left(\frac{1}{N} \sum_{i=1}^N \mathbf{h}_i\right), \quad (19)$$

where σ is the sigmoid activation function. Note other advanced pooling operations can also be applied.

Mutual information measures how much the uncertainty reduced in a random variable Y when knowing another variable X , it can be defined as follows

$$I(X, Y) = H(X) - H(X | Y). \quad (20)$$

Since optimizing the mutual information is usually intractable, we leverage contrastive loss [39, 41], a lower bound of $I(X, Y)$ that works well in practice, to attain the goal of maximizing mutual information. The contrastive loss is defined as

$$\mathcal{L}_{tcl} = -\log \frac{\exp(f(\mathbf{s}, \mathbf{s}_0)/\tau)}{\sum_{i=0}^M \exp(f(\mathbf{s}, \mathbf{s}_i)/\tau)}, \quad (21)$$

where τ denotes the temperature coefficient, M is the number of negative sub-trajectories sampled from other trajectories, $f(\cdot, \cdot)$ is implemented as a linear scoring function $f(\mathbf{s}, \mathbf{s}_i) = \sigma(\mathbf{s}_i^\top W_c \mathbf{s})$, indicating the likelihood of \mathbf{s}_i belong to the target trajectory. By using the above contrastive loss, we can increase the similarity between the entire trajectory and sub-trajectory, and meanwhile, the similarities between different trajectories are decreased. Therefore, the model can effectively maximize the mutual information between local and global parts of a specific trajectory.

In this way, two distant nodes in a same trajectory are enforced to be similar by preserving the mutual information. In addition, rich implicit information (*e.g.*, functional patterns) of unlabeled trajectories can be incorporated into our model for route recommendation.

6.3 Learning Strategy and Discussion

In this work, our pre-training module works sequentially by first performing masked attribute prediction and then trajectory contrastive learning. When the pre-training is done, we can fine-tune the entire model based on pre-trained low-level model parameters on transportation recommendation task in an end-to-end manner. The spatiotemporal pre-training can learn both contextual information and sequential patterns by utilizing both transportation network structure and unlabeled historical trajectory data.

Both the hierarchical multi-task learning module and the pre-training module exploit semantic and auxiliary information in the transportation network as well as historical trajectories, where the hierarchical multi-task learning module follows the joint training paradigm, and the pre-training module follows the pre-train and fine-tune paradigm. In practice, the pre-training can leverage large-scale unlabeled trajectories and further improve the transportation recommendation effectiveness.

7 Complexity Analysis

In this section, we analyze the computational complexity of each component of HMTRL⁺, including the traffic condition inference, the hierarchical multi-task route representation learning, and the spatiotemporal pre-training.

Complexity of traffic condition inference. Let $|V^h|$, $|E^h|$, $|V^l|$, and $|E^l|$ denote the number of nodes and edges in hub-centric and link-centric graphs, respectively. We adopt sparse matrix multiplication to efficiently implement graph convolution operation. For each time slot, the complexity of the inference component is

$$\mathcal{T}_{infer} = \mathcal{O}(F^2l(|E^h| + |E^l|)), \quad (22)$$

where F and l respectively represent the number of input features and stacked convolutional layers. Here we derive $|V| = |V^h \cup V^l|$ and $|E| = |E^h \cup E^l|$. The complexity of traffic condition inference can be written as $\mathcal{T}_{infer} = \mathcal{O}(F^2l|E|)$, *i.e.*, linear in the number of graph edges.

Complexity of HMTRL. For each user query, we generate a candidate route set Γ based on existing routing engine and rank these routes by using HMTRL. The computational cost comes from three modules, *i.e.*, spatiotemporal graph neural network, route representation learning and hierarchical multi-task learning. Same as traffic condition inference, the complexity of spatial autocorrelation modeling module at each time step is $\mathcal{O}(F^2l|E|)$. Additionally, GRU block is used in temporal autocorrelation modeling module. For each time step, the GRU operation has complexity $\mathcal{O}(F^2|V|)$. Assume we adopt previous T step features as model input,

the total complexity of the spatiotemporal graph neural network is

$$\mathcal{T}_{st} = \mathcal{O}(TF^2(|E| + |V|)). \quad (23)$$

Here we set the input length $T = 2$ to reduce the computational cost. Then we analyze the complexity of route representation learning, which comprises of two blocks, route coherence modeling block and self-attention block. On the one hand, the complexity of route coherence modeling block (*i.e.*, bi-directional GRU) is $\mathcal{O}(|\Gamma|R_{max}F^2)$, where R_{max} denote the maximum length of candidate routes, $|\Gamma|$ is the number of routes in candidate set. On the other hand, the complexity of the self-attention block (*i.e.*, multi-head self-attentive operation) is $\mathcal{O}(|\Gamma|KR_{max}^2F)$, where K denote the number of self-attention operations in our model. Therefore, the total complexity of route representation learning can be written as

$$\mathcal{T}_{route} = \mathcal{O}(|\Gamma|R_{max}F(KR_{max} + F)). \quad (24)$$

Moreover, the complexity of the hierarchical multi-task learning module is

$$\mathcal{T}_{mt} = \mathcal{O}(F|V|\tau_1 + F|\Gamma|\tau_2), \quad (25)$$

where τ_1 and τ_2 are the number of vertex-level and route-level auxiliary tasks. To summarize, the overall complexity of HMTRL is the combination of the above three modules, which can be written as follows

$$\mathcal{T}_{hmtrl} = \mathcal{O}(TF^2(|E| + |V|) + F(|V|\tau_1 + |\Gamma|\tau_2) + |\Gamma|R_{max}F(KR_{max} + F)). \quad (26)$$

Since the real-world transportation networks are usually extremely large, the number of vertices and edges of transportation graph is the major bottleneck in HMTRL. However, once we obtain all the node representations by using spatiotemporal autocorrelation modeling component at a specific time slot, we can directly reuse these node representations to answer user queries. Therefore, we only need to run the low-level spatiotemporal component once in the same time period, and the computational complexity can be significantly reduced. In practice, the proposed model is efficient and applicable on large-scale transportation networks.

Complexity of spatiotemporal pre-training. We pre-train low-level layers of HMTRL, *i.e.*, the spatiotemporal graph neural network. The complexity of spatiotemporal pre-training component is similar with Equation 23. Besides, the complexity of masked attribute prediction is $\mathcal{O}(F^2|V|)$. The complexity of trajectory contrastive learning is $\mathcal{O}(F^2M)$, where M denotes the number of negative trajectories. Therefore, the overall complexity of spatiotemporal pre-training can be written as

$$\mathcal{T}_{pre} = \mathcal{O}(TF^2(|E| + |V|) + F^2(|V| + M)). \quad (27)$$

Table 1 Statistics of datasets.

Data description	BEIJING	SHANGHAI
# of routing queries & trajectories	2,804,274	2,101,028
# of road intersections	334,421	333,163
# of road segments	420,889	426,247
# of bus lines	22,364	25,652
# of bus stations	9,651	11,587

8 Experiments

8.1 Data Description

We conduct experiments on two real-world datasets, BEIJING and SHANGHAI. Both datasets are provided by one of the world’s largest navigation applications in the world. The datasets include: (1) transportation networks of car, bus, cycle and walk, (2) routing query data extracted from user in-app logs, (3) historical trajectory data collected from user navigation events, (4) context data including weather conditions and user demographic attributes. The raw data of BEIJING and SHANGHAI are 4.13 TB and 4.36 TB, respectively. Both datasets are ranged from August 1, 2019 to October 30, 2019. The Minimum Boundary Rectangle (MBR) of BEIJING and SHANGHAI are (116.21, 39.76), (116.56, 40.03) and (121.35, 31.12), (121.65, 31.38). The statistics of each dataset are summarized in Table 1. We chronologically order each data set, take the first 80% as training set, the following 10% for validation and the rest 10% for testing.

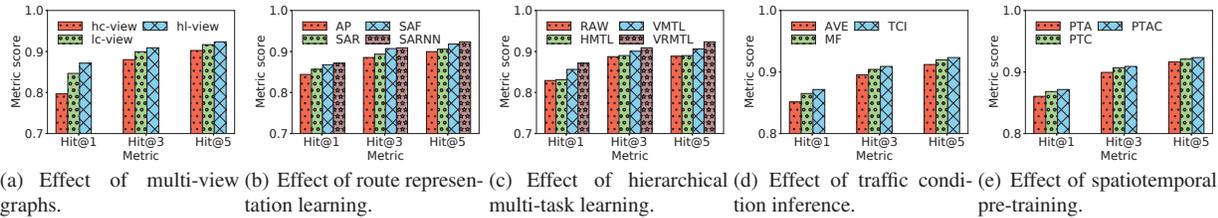
We use the route that user traveled in the real world as the ground truth. First, we match routing queries and historical trajectories based on anonymized user ID and timestamp, so that each routing query corresponds to a GPS trajectory from the origin to destination. Therefore, each matched record indicates a real-world trip after a routing query, reflecting the practical user preference. For each record, we use unselected routes generated by Baidu Maps as negative samples to recover the real context when users issue routing queries. For example, given origin and destination as well as some candidate routes r_1, r_2, r_3 , if user traveled from origin to destination via route r_1 in the real-world trajectory dataset, we set the label of route r_1 to 1, and the other two routes r_2 and r_3 to 0.

8.2 Implementation Details

HMTRL settings. We initialize all trainable parameters randomly with the uniform distribution. We apply an embedding operation to project each categorical features to 16-dimensional embedding vectors and concatenate them with continuous features. The dimension of hidden state d is fixed to 64. We stack two layers of graph convolution to capture spatial autocorrelation, and choose LeakyReLU ($\alpha = 0.2$) as the activation function in graph convolution operation.

Table 2 Overall performance comparison using six metrics on BEIJING and SHANGHAI.

Method	BEIJING						SHANGHAI					
	Hit@1	Hit@3	Hit@5	NDCG@3	NDCG@5	NDCG@10	Hit@1	Hit@3	Hit@5	NDCG@3	NDCG@5	NDCG@10
RBT	0.1337	0.3874	0.5794	0.3515	0.4278	0.4713	0.1596	0.4121	0.5609	0.3119	0.3730	0.4569
RBD	0.3647	0.6212	0.7339	0.4801	0.5583	0.6116	0.3178	0.4437	0.6096	0.4068	0.4337	0.5051
LR	0.7188	0.8329	0.8705	0.7864	0.8018	0.8203	0.6687	0.8010	0.8423	0.7468	0.7638	0.7827
GBDT	0.7370	0.8474	0.8851	0.8021	0.8176	0.8341	0.6814	0.8083	0.8524	0.7563	0.7745	0.7950
DeepWalk	0.5213	0.6642	0.7587	0.5955	0.6344	0.6731	0.4916	0.6591	0.7591	0.5886	0.6297	0.6687
DeepFM	0.7658	0.8538	0.8853	0.8166	0.8295	0.8452	0.7068	0.8209	0.8592	0.7743	0.7901	0.8096
Hydra	0.7604	0.8508	0.8827	0.8139	0.8270	0.8434	0.7292	0.8324	0.8713	0.7854	0.8177	0.8251
MURAT	0.7892	0.8654	0.8993	0.8345	0.8467	0.8631	0.7508	0.8415	0.8889	0.8009	0.8297	0.8334
HMTRL	0.8545	0.8946	0.9184	0.8735	0.8856	0.8990	0.8115	0.8823	0.9111	0.8533	0.8652	0.8761
HMTRL ⁺	0.8716	0.9091	0.9235	0.8907	0.8983	0.9082	0.8326	0.8974	0.9163	0.8735	0.8829	0.8904

**Fig. 5** Ablation study of HMTRL⁺ on BEIJING.

We employ a sigmoid function in the final output layer. The hyper-parameters K , β_1 , β_2 , T , γ are set to 8, 0.3, 0.1, 3, 0.5, respectively. We set the learning rate $lr = 0.0001$ and the batch size 256. We fix the length of the sub-route to 6 for coherence modeling. We evaluate our model as well as all baselines on a powerful Linux server with 26 Intel Xeon Gold 5117 CPUs, 8 NVIDIA Tesla P40 GPUs, 256GB memory and 10TB disk. For a fair comparison, we carefully fine-tuned the hyper-parameters for all baselines on our datasets via grid search based on settings in their original paper. Please refer to source code¹ for more details.

8.3 Metrics

We employ Hit@ k and Normalized Discounted Cumulative Gain (NDCG@ k) [42], two widely used metrics in recommenders, to evaluate the recommendation effectiveness. In the following experiments, we report Hit@1, Hit@3, Hit@5, and NDCG@3, NDCG@5, NDCG@10.

8.4 Baselines

We compare HMTRL⁺ with two rule-based methods and six learning methods.

- **RBT** is a rule-based method that recommends the fastest route, in which we rank route candidates by ETA.
- **RBD** is another rule-based method that recommends the shortest route, in which we rank route candidates by road network distance.

- **LR** uses logistic regression [43] for recommendation. The inputs are same as raw features used in HMTRL⁺.
- **GBDT** adopts the Gradient Boosting Decision Tree for recommendation, which is widely used in both academia and industry. We implement the baseline based on XGboost [44]. The input features are the same as LR.
- **DeepWalk** [45] is a unsupervised network embedding method that learns vertex representations of a graph. We apply random walks on the MMTN to generate vertex representations, and apply average pooling on route sequences to obtain route representation. We further apply a LR layer for recommendation.
- **DeepFM** [46] is a deep recommendation model that combines the factorization machine and deep neural network to model both first-order and higher-order feature interactions. The input is the same as HMTRL⁺.
- **Hydra** [8] is the state-of-the-art multi-modal transport mode recommendation method based on multi-sourced urban data. It is fed both handcrafted features as well as pre-trained latent embedding features to a gradient boosting tree-based model. We extend it by adding a regression layer to enable multi-modal route recommendation.
- **MURAT** [47] is a novel multi-task graph representation learning framework for travel time estimation. We also use our multi-view graphs as the input and devise the the output layer to fit our recommendation task.

¹ <https://github.com/hanjindong/HMTRL-Pytorch>

8.5 Overall Performance

Table 2 shows the overall performance of our method and all the compared baselines on two datasets with respect to six evaluation metrics. Overall, HMTRL⁺ outperforms all the baselines on both datasets using all metrics, which demonstrate the advance of our model. To be specific, HMTRL⁺ achieves (10.4%, 5.0%, 2.7%) Hit@*k* and (6.7%, 6.1%, 5.2%) NDCG@*k* improvement compared with the state-of-the-art approach (MURAT) on BEIJING. Similarly, the improvement of Hit@*k* and NDCG@*k* on SHANGHAI dataset are (10.9%, 6.6%, 3.1%) and (9.1%, 6.4%, 6.8%). We also observe HMTRL⁺ consistently outperforms HMTRL in terms of all metrics, indicates the effectiveness of the traffic condition inference and spatiotemporal pre-training. Compared with HMTRL, HMTRL⁺ achieves (2.0%, 1.6%, 0.6%) and (1.9%, 1.4%, 1.0%) improvements on Hit@*k* and NDCG@*k* on the BEIJING dataset, and the improvement on SHANGHAI is consistent.

Moreover, we can make the following observations. (1) The performance of RBT is much worse than RBD. This observation indicates that travel distance is a more significant indicator than ETA for user trip decision. (2) DeepWalk achieves a better performance than rule-based methods, but performs worse than other learning-based methods. The main reason is that DeepWalk can leverage the structural information but it fails to consider contextual features. Besides, due to its unsupervised property, DeepWalk neglects the user preference signal in historical data. (3) Hydra outperforms all other non-deep learning models by incorporating fine-grained handcrafted features and high-order embedding features. However, compared with deep learning-based methods, including DeepFM and MURAT, the manually extracted features limit the recommendation capability of the model. (4) MURAT consistently outperforms all other baselines, which demonstrate the effectiveness of multi-task graph representation learning. However, MURAT neglects the information in link-centric graphs as well as the low level supervision signals, therefore performs worse than our approach.

8.6 Ablation Study

Then we conduct ablation study on HMTRL⁺.

Effect of multi-view graphs. We first examine the effectiveness of multi-view graphs by evaluating three variants of HMTRL⁺, (1) *hc-view* only uses the hub-centric graph, (2) *lc-view* uses the link-centric graph only, and (3) *hl-view* uses both graphs for recommendation. As shown in Figure 5(a), the performance of *hl-view* outperforms *hc-view* and *lc-view* by (9.6%, 3.4%, 2.5%) and (4.4%, 1.2%, 0.9%) on (Hit@1, Hit@3 and Hit@5), respectively. Moreover, the *lc-view* performs better than *hc-view*, which demonstrate the

structural and contextual information in transportation links plays a more important role for multi-modal transportation recommendation.

Effect of coherent-aware attentive route representation learning. We further construct and evaluate the following variants, (1) *AP* uses average pooling to aggregate hub and link representations, (2) *SAR* derives route representation by self-attention only, (3) *SAF* removes backward GRU in *BiGRU*, and (4) *SARNN* includes both self-attentive operation and the *BiGRU* to integrate route coherence. As shown in Figure 5(b), self-attentive based route aggregation achieves a better performance than *AP*. Moreover, by integrating the route coherence, *SARNN* achieves significant improvement compared with *SAR*. Additionally, compared with *SAF*, we observe *SARNN* achieves consistent improvement by incorporating backward sub-route coherence, demonstrate the effectiveness of bi-directional RNN.

Effect of hierarchical multi-task learning. We compare the following variants, (1) *RAW* directly learns route representation without auxiliary tasks, (2) *HMTL* only incorporates vertex-level hub-related tasks, (3) *VMTL* incorporates both vertex-level hub-related and link-related tasks, and (4) *VRMTL* integrates both vertex-level tasks and route-level tasks. As reported in Figure 5(c), we observe consistent improvement by respectively adding vertex level and route level auxiliary tasks, validate the effectiveness of different supervision signals for multi-modal transportation recommendation. In particular, *VMTL* achieves more significant improvement over *HMTL* than *HMTL* over *RAW*, indicating link related auxiliary tasks plays a more important role in multi-modal transportation recommendations.

Effect of traffic condition inference. To validate the effect of traffic condition inference, we exam (1) *AVE* uses average values to interpolate the missing traffic conditions, (2) *MF* imputes missing values by applying a collaborative matrix factorization method [48], (3) *TCI* infers traffic condition through the proposed context encoder. As reported in Figure 5(d), the proposed approach achieves the best performance in recommendation task than other imputing methods, either average interpolation or matrix factorization. Overall, real-time traffic conditions are strong signals for transportation recommendation, accurate traffic condition inference can further provide more information to improve model performance.

Effect of spatiotemporal pre-training. Finally, we verify the spatiotemporal pre-training used in our framework. Specifically, we compare three variants of HMTRL⁺: (1) *STA* excludes trajectory contrastive learning, (2) *STC* without attribute prediction, (3) *STAC* incorporates both attribute prediction and trajectory contrastive learning tasks for pre-training. As shown in Figure 5(e), we observe only performing attribute prediction during pre-training stage gives limited performance improvement, while trajectory contrastive

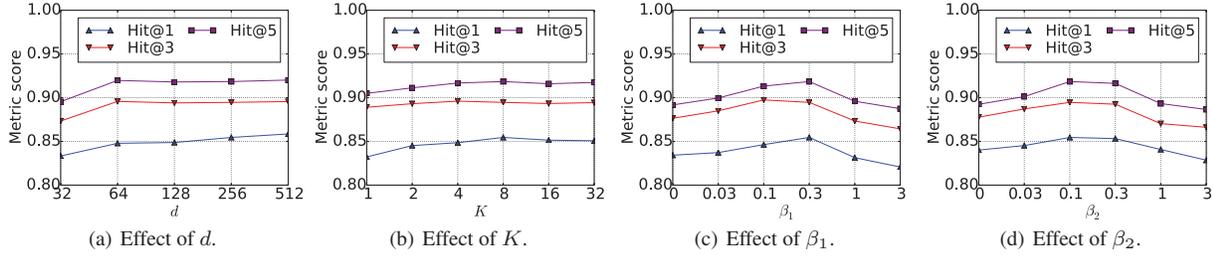


Fig. 6 Parameter sensitivities on BEIJING.

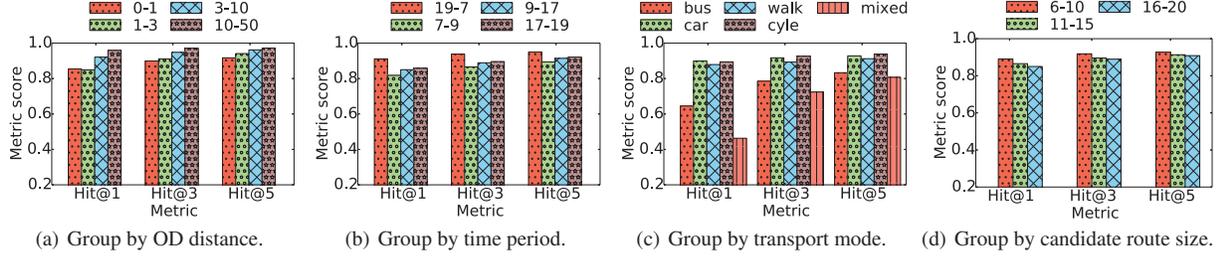


Fig. 7 Robustness check on BEIJING.

learning have significantly better performance gain than attribute prediction task. This finding confirms unlabeled trajectories introduce additional crucial information and patterns, which can be utilized to make our model more generalizable and improve the total recommendation performance.

8.7 Parameter Sensitivity

We further study the parameter sensitivity of HMTRL. Each time we vary a parameter, we set others to their default values.

First, we vary the dimension d from 32 to 512. The results are reported in Figure 6(a). As the dimension increases, the performance first increases and then remains stable. However, too large d will induce a higher training cost. Therefore, set the dimension to 64 is enough to capture representation information.

Then, we vary the number of self-attentive operations K from 1 to 32. The results are reported in Figure 6(b). We observe a performance improvement when increasing K from 1 to 8, but a slight performance degradation by further increasing K from 8 to 32. Using 8 self-attentive operations is good enough to capture diversified vertex importance for route representation learning.

After that, we vary vertex-level multi-task weight β_1 from 0 to 3. The results are reported in Figure 6(c). We observe a significant performance gain when increasing β from 0 to 0.3, and then the performance degrades when we further increase β from 0.3 to 3. Above results prove incorporating

low-level supervision signals is beneficial to the main recommendation task, but may introduce more noises with too large task weight.

Finally, to test the impact of route-level auxiliary tasks weight, we vary β_2 from 0 to 3. The results are reported in Figure 6(d). HMTRL achieves the best performance when $\beta_2 = 0.1$, and we observe a performance degradation when we increase or decrease β_2 . This is possibly because too small β_2 cannot fully take advantage of the common information in route level auxiliary tasks, whereas too large β_2 makes the auxiliary tasks dominate the optimization and weakens the importance of the main recommendation task.

8.8 Robustness Check

A robust transportation recommendation model should perform evenly well in different routing query subgroups. We evaluate the robustness of HMTRL⁺ from the following three perspectives. First, we group queries by OD pair distance, *i.e.*, less than 1Km, 1Km to 3Km, 3Km to 10Km, and more than 10Km. Second, we split routing queries by day using four time intervals, *i.e.*, the morning peak hour (7,9], the evening peak hour (17,19], and two off-peak intervals (19,7], (9,17]. Third, we group queries based on the selected transport mode, including bus, cycle, walk, car and mixed (*i.e.*, route with more than one transport modes). Finally, we split routing queries based on the size of route candidate set \mathcal{L} , *i.e.*, 6 to 10, 11 to 15, and 16 to 20. Figure 7 shows the results of HMTRL⁺ on different subgroups on BEIJING. For different OD distance intervals, we observe the performance

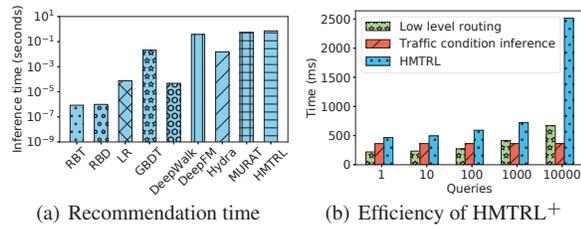


Fig. 8 Efficiency analysis.

difference is smaller than 10.9%. Besides, our model performs better on longer distance transportation recommendations, which is perhaps because routes of walk and cycle are no longer attractive for long distance trip, therefore ease the recommendation. For different time periods, we observe the difference is smaller than 10.2%. Besides, we observe a more accurate recommendation result at night, and a worse result in morning rush hour. This is possibly because the traffic condition during the morning rush hour is more complex and hard to predict. For different transport modes, we observe the performance of mixed group is notably lower than others. This may be induced by the scarcity of mixed historical routes, and the route representation with combined transport mode is more sophisticated to learn. Note that the performance of HMTRL⁺ on mixed routes still significantly higher than all baselines, please refer to supplementary material for details. The above results suggests us to pay more attention on mixed routes in the further work to obtain a better overall performance. For different size of route candidate set Γ , we observe a performance degradation when increasing size of Γ from 6-10 to 16-20. The reason perhaps is that the proposed model is confused by similar routes, as a large size will introduce more diverse candidate routes and raise the difficulty to find the appropriate route.

8.9 Efficiency Analysis

Finally, we present the efficiency of each approach, which was evaluated on a server with single NVIDIA Tesla P40 GPU. We first randomly test 1000 user queries and report the average running time of our model as well as each baseline. As shown in Figure 8(a), we observe rule-based methods are fastest, while learning-based methods run much slower. In addition, deep learning based methods, such as DeepFM and MURAT, take longer time than other models. Although HMTRL is more time-consuming, it achieves nearly 10% accuracy improvement than MURAT. We also test the online latency of HMTRL⁺, which consists of three components, low-level routing, traffic condition inference and HMTRL recommendation. The results are reported in Figure 8(b). On average, we can infer city-wide traffic conditions by using the proposed component in 362ms. Moreover, when we

vary the number of queries from 1 to 10000, the low-level routing latency increased from 220ms to 671ms, while the recommendation latency increased from 466ms to 2515ms. We observe the latency gap is relatively small for HMTRL at the beginning, and goes larger when we further increase the queries from 1000 to 10000. According to the efficiency study, we find HMTRL is the major bottleneck of our model. In the future, we will further optimize the recommendation component to reduce the total latency.

9 Related Work

Route Recommendation has become a core component in map services (*e.g.*, Google Maps, Baidu Maps) and has gradually received more research attention [49, 50]. With the ubiquity of mobile devices and location-based services, massive historical data (*e.g.*, GPS trajectory data [9] and mobile check-in data [51, 52]) has been leveraged to improve the quality of route recommendation. For example, Chen *et al.* [53], Wang *et al.* [54], and Yang *et al.* [55] leverage historical trajectories for better routing, but cannot be directly generalized to multi-modal recommendations. Liu *et al.* [14] proposes a general framework for public transportation routing. Nevertheless, it focuses on uni-modal route recommendation and fails to model the relationship between different transport modes. Recently, a few machine learning based multi-modal transportation recommendation techniques has been introduced. For instance, FAVOUR [56] proposed a probabilistic model for multi-modal route recommendation based on a series of user-provided profile and survey data. Trans2vec [12] learns network embedding of users, OD pairs for transport mode recommendation, but it cannot generalize well when there exists massive cold-start users with sparse data. Besides, Hydra [8, 57] constructed various context features and MTRecS-DLT [10] developed a convolutional neural network based model for personalized transport mode recommendation. However, the above studies ignore rich semantic information in the transportation network and historical routes, which lead to unsatisfactory multi-modal route recommendations.

Graph Representation Learning extends the convolutional neural network for capturing spatial dependencies on non-Euclidean graph structures [24, 25]. Recently, graph representation learning has been widely used in many spatiotemporal mining tasks, such as flow prediction [15, 58], region representation [59], and parking availability prediction [60]. Beyond vertex classification, a few studies investigate the classification problem of sequences on dynamic graphs [61]. However, none of the above works are dedicated to multi-modal transportation recommendations.

Multi-Task Learning is a learning paradigm that aims to improve the performance of multiple correlated tasks by sharing common information. Based on information sharing

strategy, multi-task learning can be categorized into hard parameter sharing based and soft parameter sharing based [62]. Recent studies [31, 63, 30] have successfully facilitated multiple tasks in lower neural network layers to guide the overall optimization. In this paper, we employ the hierarchical multi-task learning framework by using hard parameter sharing to integrate auxiliary tasks in different network layers.

Self-Supervised Learning aims at learning transferable and generalizable feature representations based on various auxiliary supervised signals extracted from the data itself. The paradigm of self-supervised learning in deep neural networks can be categorized into two classes: pre-training and joint training. On the one hand, pre-training first learn model parameters with self-supervised signals as pretext tasks then fine-tune the neural network based on downstream tasks [64–66]. On the other hand, joint training simultaneously train self-supervised pretext tasks together with target downstream task [67, 35]. Due to its effectiveness, self-supervised learning has been applied in many fields, such as computer vision [41], natural language processing [68], graph mining [36], and trajectory data mining [69]. In this paper, we propose the spatiotemporal pre-training to fully exploit the information hidden in the dynamic road networks and massive historical trajectories.

10 Conclusion

In this paper, we proposed HMTRL⁺, a unified route representation learning framework for multi-modal transportation recommendation. We first constructed time-dependent multi-view transportation graphs to characterize the structural and contextual information of both hubs and links. Then, we devised a graph-based contextual encoder for missing traffic condition imputation. Furthermore, we proposed a spatiotemporal graph neural network for collaborative learning of spatial and temporal autocorrelation. After that, a coherent-aware self-attentive route representation learning module is introduced to project arbitrary-length routes into fixed-length route representation vectors, with explicit modeling of route coherence from historical routes. Moreover, a hierarchical multi-task learning module is proposed to derive transport mode-specific route representations for recommendation by integrating various supervision signals in different network layers. Finally, we introduced spatiotemporal pre-training strategies to enhance the robustness of the recommendation system by exploiting various self-supervision signals in the multi-modal transportation network and unlabeled historical trajectories. Extensive experimental results on two real-world datasets demonstrated the performance of HMTRL⁺ consistently outperforms eight state-of-the-art baselines. In future work, we will further reduce the error on mixed route recommendation and optimize the framework to improve the model efficiency.

Acknowledgements This work is supported by the National Natural Science Foundation of China under Grant No. 62102110, Foshan HKUST Projects (FSUST21-FYTRI01A, FSUST21-FYTRI02A), and Hong Kong RGC TRS T41-603/20-R.

References

1. Baidu maps. https://en.wikipedia.org/wiki/Baidu_Maps, 2021. Accessed: 2021-07-01.
2. Here wego. https://en.wikipedia.org/wiki/Here_WeGo, 2021. Accessed: 2021-07-01.
3. Didi. <https://en.wikipedia.org/wiki/DiDi>, 2021. Accessed: 2021-07-01.
4. Julian Dibbelt et al. Engineering algorithms for route planning in multimodal transportation networks. *Transportation*, 2016.
5. Nilesh Borole, Dillip Rout, Nidhi Goel, P Vedagiri, and Tom V Mathew. Multimodal public transit trip planner with real-time transit data. *Procedia-Social and Behavioral Sciences*, 104:775–784, 2013.
6. Robert Geisberger, Peter Sanders, Dominik Schultes, and Christian Vetter. Exact routing in large road networks using contraction hierarchies. *Transportation Science*, 46(3):388–404, 2012.
7. Lu Liu. *Data model and algorithms for multimodal route planning with transportation networks*. PhD thesis, Technische Universität München, 2011.
8. Hao Liu, Yongxin Tong, Panpan Zhang, Xinjiang Lu, Jianguo Duan, and Hui Xiong. Hydra: A personalized and context-aware multi-modal transportation recommendation system. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2314–2324, 2019.
9. Jing Yuan, Yu Zheng, Chengyang Zhang, Wenlei Xie, Xing Xie, Guangzhong Sun, and Yan Huang. T-drive: driving directions based on taxi trajectories. In *Proceedings of the 18th SIGSPATIAL International conference on advances in geographic information systems*, pages 99–108, 2010.
10. Ayat Abedalla, Ali Fadel, Ibraheem Tuffaha, Hani Al-Omari, Mohammad Omari, Malak Abdullah, and Mahmoud Al-Ayyoub. Mtrecc-dlt: Multi-modal transport recommender system using deep learning and tree models. In *2019 Sixth International Conference on Social Networks Analysis, Management and Security (SNAMS)*, pages 274–278. IEEE, 2019.
11. Hao Zhou, Yan Zhao, Junhua Fang, Xuanhao Chen, and Kai Zeng. Hybrid route recommendation with taxi and shared bicycles. *Distributed and Parallel Databases*, pages 1–21, 2019.
12. Hao Liu, Ting Li, Renjun Hu, Yanjie Fu, Jingjing Gu, and Hui Xiong. Joint representation learning for multi-modal transportation recommendation. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*, pages 1036–1043, 2019.
13. Hao Liu, Jindong Han, Yanjie Fu, Jingbo Zhou, Xinjiang Lu, and Hui Xiong. Multi-modal transportation recommendation with unified route representation learning. *Proceedings of the VLDB Endowment*, 14(3):342–350, 2020.
14. Hao Liu, Ying Li, Yanjie Fu, Huaibo Mei, Jingbo Zhou, Xu Ma, and Hui Xiong. Polestar: An intelligent, efficient and nationwide public transportation routing engine. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2321–2329, 2020.
15. Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *6th International Conference on Learning Representations*, 2018.
16. Pengyang Wang, Yanjie Fu, Jiawei Zhang, Pengfei Wang, Yu Zheng, and Charu Aggarwal. You are how you drive: Peer and

- temporal-aware representation learning for driving behavior analysis. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2457–2466, 2018.
17. Dragoš Cvetkovic, Dragoš M Cvetković, Peter Rowlinson, Slobodan Simic, and Slobodan Simić. *Spectral generalizations of line graphs: On graphs with least eigenvalue-2*, volume 314. Cambridge University Press, 2004.
 18. Jilin Hu, Chenjuan Guo, Bin Yang, and Christian S Jensen. Stochastic weight completion for road networks using graph convolutional networks. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pages 1274–1285. IEEE, 2019.
 19. Jilin Hu, Bin Yang, Chenjuan Guo, Christian S Jensen, and Hui Xiong. Stochastic origin-destination matrix forecasting using dual-stage graph convolutional, recurrent neural networks. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, pages 1417–1428. IEEE, 2020.
 20. Chenjuan Guo, Bin Yang, Jilin Hu, and Christian Jensen. Learning to route with sparse trajectory sets. In *2018 IEEE 34th International Conference on Data Engineering (ICDE)*, pages 1073–1084. IEEE, 2018.
 21. Chenjuan Guo, Bin Yang, Jilin Hu, Christian S Jensen, and Lu Chen. Context-aware, preference-based vehicle routing. *The VLDB Journal*, 29(5):1149–1170, 2020.
 22. Bin Yang, Manohar Kaul, and Christian S Jensen. Using incomplete information for complete weight annotation of road networks. *IEEE Transactions on Knowledge and Data Engineering*, 26(5):1267–1279, 2013.
 23. Hao Liu, Qiyu Wu, Fuzhen Zhuang, Xinjiang Lu, Dejing Dou, and Hui Xiong. Community-aware multi-task transportation demand prediction. In *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2021.
 24. Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations*, 2017.
 25. Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034, 2017.
 26. Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
 27. Shuai Li, Wanqing Li, Chris Cook, Ce Zhu, and Yanbo Gao. Independently recurrent neural network (indrnn): Building a longer and deeper rnn. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5457–5466, 2018.
 28. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
 29. Yu Zhang and Qiang Yang. A survey on multi-task learning. *arXiv preprint arXiv:1707.08114*, 2017.
 30. Kazuma Hashimoto, Caiming Xiong, Yoshimasa Tsuruoka, and Richard Socher. A joint many-task model: Growing a neural network for multiple nlp tasks. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1923–1933, 2017.
 31. Hongliang Fei, Shulong Tan, and Ping Li. Hierarchical multi-task word embedding learning for synonym prediction. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 834–842, 2019.
 32. R Caruana. Multitask learning: A knowledge-based source of inductive bias. machine learning, 1997.
 33. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 34. Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning can improve model robustness and uncertainty. *Advances in Neural Information Processing Systems*, 32, 2019.
 35. Yuning You, Tianlong Chen, Zhangyang Wang, and Yang Shen. When does self-supervision help graph convolutional networks? In *International Conference on Machine Learning*, pages 10871–10880. PMLR, 2020.
 36. Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. Strategies for pre-training graph neural networks. 2020.
 37. Jing Yuan, Yu Zheng, and Xing Xie. Discovering regions of different functions in a city using human mobility and pois. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 186–194, 2012.
 38. Nicholas Jing Yuan, Yu Zheng, Xing Xie, Yingzi Wang, Kai Zheng, and Hui Xiong. Discovering urban functional zones using latent activity trajectories. *IEEE Transactions on Knowledge and Data Engineering*, 27(3):712–725, 2014.
 39. R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. 2019.
 40. Yin Lou, Chengyang Zhang, Yu Zheng, Xing Xie, Wei Wang, and Yan Huang. Map-matching for low-sampling-rate gps trajectories. In *Proceedings of the 17th ACM SIGSPATIAL international conference on advances in geographic information systems*, pages 352–361, 2009.
 41. Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
 42. Kalervo Järvelin and Jaana Kekäläinen. Cumulated gain-based evaluation of ir techniques. *ACM Transactions on Information Systems (TOIS)*, pages 422–446, 2002.
 43. David G Kleinbaum, K Dietz, M Gail, Mitchel Klein, and Mitchell Klein. *Logistic regression*. Springer, 2002.
 44. Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 785–794, 2016.
 45. Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710, 2014.
 46. Hui Feng Guo, Ruiming TANG, Yunming Ye, Zhenguo Li, and Xi-qi He. Deepfm: A factorization-machine based neural network for ctr prediction. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, pages 1725–1731, 2017.
 47. Yaguang Li, Kun Fu, Zheng Wang, Cyrus Shahabi, Jieping Ye, and Yan Liu. Multi-task representation learning for travel time estimation. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1695–1704, 2018.
 48. Jingbo Shang, Yu Zheng, Wenzhu Tong, Eric Chang, and Yong Yu. Inferring gas consumption and pollution emission of vehicles throughout a city. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1027–1036, 2014.
 49. Evangelos Kanoulas, Yang Du, Tian Xia, and Donghui Zhang. Finding fastest paths on a road network with speed patterns. In *22nd International Conference on Data Engineering*, pages 10–10, 2006.
 50. Ling-Yin Wei, Yu Zheng, and Wen-Chih Peng. Constructing popular routes from uncertain trajectories. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 195–203, 2012.
 51. Dawei Chen, Cheng Soon Ong, and Lexing Xie. Learning points and routes to recommend trajectories. In *Proceedings of the 25th*

- ACM International on Conference on Information and Knowledge Management*, pages 2227–2232, 2016.
52. Samia Shafique and Mohammed Eunus Ali. Recommending most popular travel path within a region of interest from historical trajectory data. In *Proceedings of the 5th ACM SIGSPATIAL International Workshop on Mobile Geographic Information Systems*, pages 2–11. ACM, 2016.
 53. Lisi Chen, Shuo Shang, Christian S Jensen, Bin Yao, Zhiwei Zhang, and Ling Shao. Effective and efficient reuse of past travel behavior for route recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 488–498, 2019.
 54. Jingyuan Wang, Ning Wu, Wayne Xin Zhao, Fanzhang Peng, and Xin Lin. Empowering a* search algorithms with neural networks for personalized route recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 539–547, 2019.
 55. Sean Bin Yang, Chenjuan Guo, and Bin Yang. Context-aware path ranking in road networks. *IEEE Transactions on Knowledge and Data Engineering*, 2020.
 56. Paolo Campigotto, Christian Rudloff, Maximilian Leodolter, and Dietmar Bauer. Personalized and situation-aware multimodal route recommendations: the favour algorithm. *IEEE Transactions on Intelligent Transportation Systems*, 18(1):92–102, 2016.
 57. Hao Liu, Yongxin Tong, Jindong Han, Panpan Zhang, Xinjiang Lu, and Hui Xiong. Incorporating multi-source urban data for personalized and context-aware multi-modal transportation recommendation. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1, 2020.
 58. Yuandong Wang, Hongzhi Yin, Hongxu Chen, Tianyu Wo, Jie Xu, and Kai Zheng. Origin-destination matrix prediction via graph convolution: a new perspective of passenger demand modeling. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1227–1235, 2019.
 59. Yunchao Zhang, Yanjie Fu, Pengyang Wang, Xiaolin Li, and Yu Zheng. Unifying inter-region autocorrelation and intra-region structures for spatial embedding via collective adversarial learning. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 1700–1708, 2019.
 60. Weijia Zhang, Hao Liu, Yanchi Liu, Jingbo Zhou, and Hui Xiong. Semi-supervised hierarchical recurrent graph neural network for city-wide parking availability prediction. pages 1186–1193, 2020.
 61. Jia Li, Zhichao Han, Hong Cheng, Jiao Su, Pengyun Wang, Jianfeng Zhang, and Lujia Pan. Predicting path failure in time-evolving graphs. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 1279–1289, 2019.
 62. Yu Zhang and Qiang Yang. A survey on multi-task learning. *arXiv preprint arXiv:1707.08114*, 2017.
 63. Anders Søgaard and Yoav Goldberg. Deep multi-task learning with low level tasks supervised at lower layers. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 231–235, 2016.
 64. Trieu H Trinh, Minh-Thang Luong, and Quoc V Le. Selfie: Self-supervised pretraining for image embedding. *arXiv preprint arXiv:1906.02940*, 2019.
 65. Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. Gcc: Graph contrastive coding for graph neural network pre-training. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1150–1160, 2020.
 66. Ziniu Hu, Yuxiao Dong, Kuansan Wang, Kai-Wei Chang, and Yizhou Sun. Gpt-gnn: Generative pre-training of graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1857–1867, 2020.
 67. Zhongzheng Ren and Yong Jae Lee. Cross-domain self-supervised multi-task feature learning using synthetic imagery. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 762–771, 2018.
 68. Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. 2019.
 69. Sean Bin Yang, Chenjuan Guo, Jilin Hu, Jian Tang, and Bin Yang. Unsupervised path representation learning with curriculum negative sampling. *arXiv preprint arXiv:2106.09373*, 2021.