
LIGHTPATH: LIGHTWEIGHT AND SCALABLE PATH REPRESENTATION LEARNING

Sean Bin Yang

Aalborg University, Denmark
sean.byang@gmail.com

Jilin Hu *

East China Normal University, China
Aalborg University, Denmark
jlhu@dase.ecnu.edu.cn

Chenjuan Guo

East China Normal University, China
Aalborg University, Denmark
cjguo@dase.ecnu.edu.cn

Bin Yang

East China Normal University, China
Aalborg University, Denmark
byang@dase.ecnu.edu.cn

Christian S. Jensen

Aalborg University, Denmark
csj@cs.aau.dk

ABSTRACT

Movement paths are used widely in intelligent transportation and smart city applications. To serve such applications, path representation learning aims to provide compact representations of paths that enable efficient and accurate operations when used for different downstream tasks such as path ranking and travel cost estimation. In many cases, it is attractive that the path representation learning is lightweight and scalable; in resource-limited environments and under green computing limitations, it is essential. Yet, existing path representation learning studies focus on accuracy and pay at most secondary attention to resource consumption and scalability. We propose a lightweight and scalable path representation learning framework, termed *LightPath*, that aims to reduce resource consumption and achieve scalability without affecting accuracy, thus enabling broader applicability. More specifically, we first propose a sparse auto-encoder that ensures that the framework achieves good scalability with respect to path length. Next, we propose a relational reasoning framework to enable faster training of more robust sparse path encoders. We also propose global-local knowledge distillation to further reduce the size and improve the performance of sparse path encoders. Finally, we report extensive experiments on two real-world datasets to offer insight into the efficiency, scalability, and effectiveness of the proposed framework.

Keywords Path representation learning · Lightweight · Self-supervised learning

1 Introduction

Motivated in part by an increasing number of intelligent transportation and smart city services that operate on movement paths, path representation learning (PRL) has received remarkable attention [1, 2, 3]. Path representation learning aims to learn a generic path representation (PR) vector (ref. \mathcal{R}^d in Figure 1) that can be utilized in a range of different downstream tasks. This is in contrast to task-specific path representation learning performed by supervised methods that yield representations that work well on task-labeled data but work poorly in other tasks. For example, in Figure 1 *Lightpath* takes as input a path p and returns a generic PR that can support a variety of tasks, e.g., travel time estimation and path ranking.

In fact, a variety of intelligent transportation services involve paths, e.g., travel cost estimation [4, 5, 6, 7, 8, 9], trajectory analysis [10, 11, 12, 13, 14, 15, 16, 17], and path ranking [18, 19, 4, 20, 21]. Path representations that are both accurate and compact, thus facilitating efficient operations, are in high demand as they hold the potential to significantly improve the services that use them. Indeed, recent path representation learning methods, in particular deep learning based methods, demonstrate impressive and state-of-the-art performance on a wide variety of downstream tasks.

*Corresponding Author: Jilin Hu (jlhu@dase.ecnu.edu.cn)

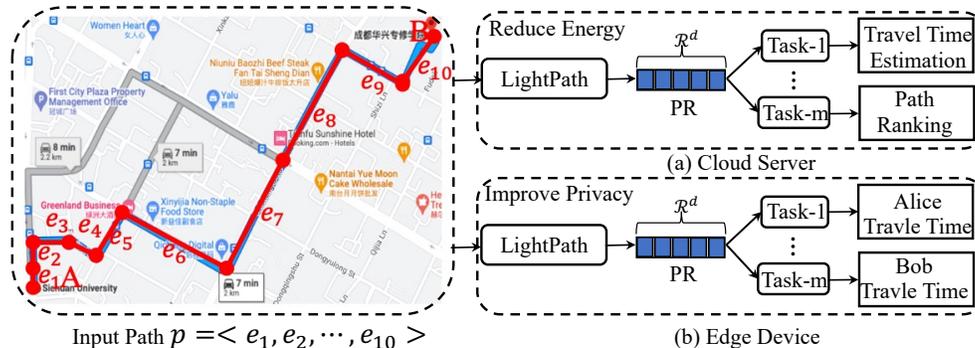


Figure 1: Intuition of the Lightweight Path Representation Learning Problem

Table 1: Model Parameter Size with Varying Encoder Layers

| Encoder Layers L | 12 | 24 | 48 | 96 |
|-----------------------|-------|-------|--------|--------|
| Parameters (Millions) | 29.85 | 55.07 | 105.51 | 206.40 |

However, existing path representation learning methods focus on accuracy improvement and pay at best secondary attention to scalability and resource usage. The resulting models often include large numbers of layers and parameters, driving up computational costs, power consumption, and memory consumption, especially for long paths. Although path encoders with many parameters may achieve good accuracy, they have two limitations. First, using large path encoders in the cloud consumes substantial energy, which is not eco-friendly (cf. Fig 1(a)). Second, increasingly many users enjoy personalized services, e.g., personalized travel time estimation based on their own trajectories. Due to privacy concerns, such personalized services often require the path encoder to be deployed in resource-limited edge environments, such as on individual users’ mobile phones (cf. Fig 1(b)), without having to upload their trajectories to the cloud. More generally, it is sensible to enable lightweight path representation learning that works in resource-limited environments.

Next, existing path representation methods suffer from two limitations.

Poor scalability w.r.t. path length Since a path is a sequence of road-network edges, path representation learning benefits from models that are good at capturing sequential relationships, such as the Transformer [22]. However, a Transformer-based method [23] employs a self-attention mechanism, where one edge attends to all other edges in a path in each attention, resulting in quadratic complexity, $\mathcal{O}(N^2)$, in the path length N . This results in poor scalability to long paths with many edges. Figure 2 gives an example of the scalability w.r.t. path length N , covering both

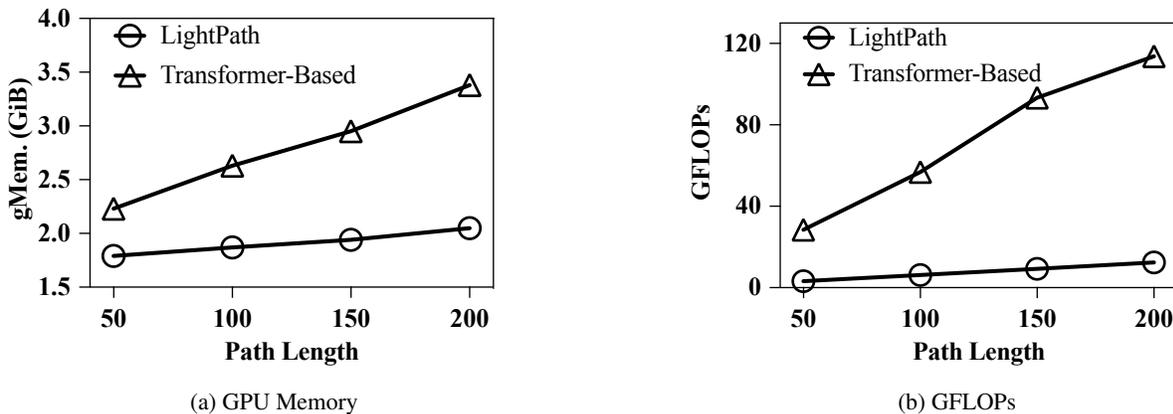


Figure 2: Scalability w.r.t. Path Length.

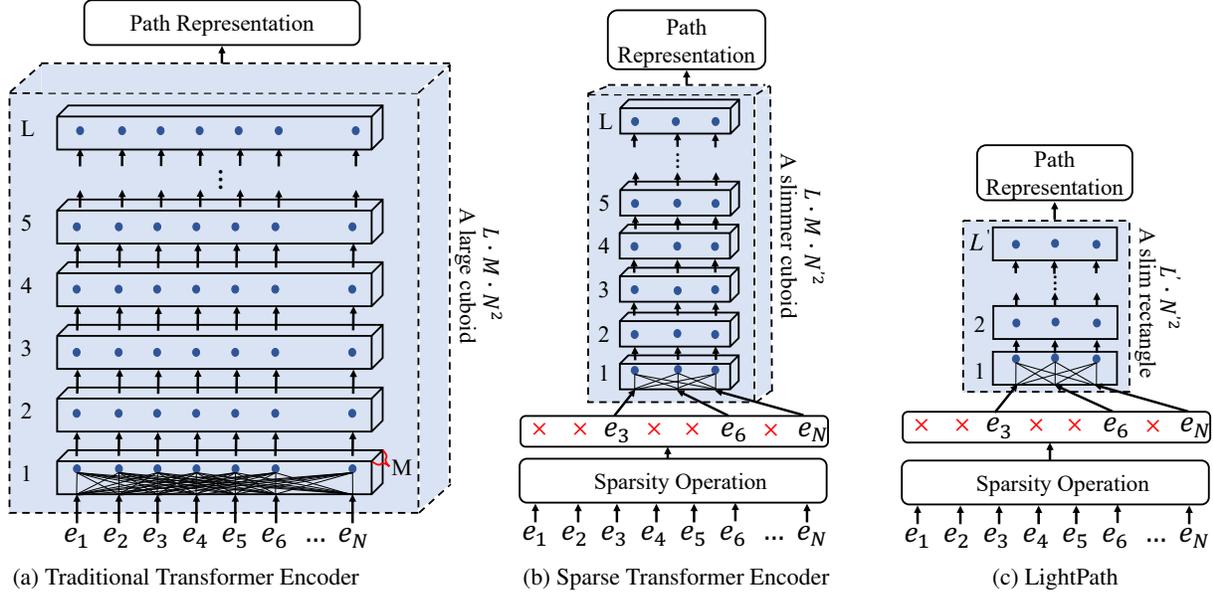


Figure 3: Encoder Architectures: (a) A traditional transformer encoder with L layers and M heads, takes as input a path (N -Length) and has complexity $\mathcal{O}(L \cdot M \cdot N^2)$; (b) A sparse transformer encoder takes as input a sparse path (i.e., reducing path length from N to N'), resulting in $\mathcal{O}(L \cdot M \cdot N'^2)$ complexity; (c) *LightPath* further compresses the traditional transformer in terms of layers and heads, yielding complexity $\mathcal{O}(L' \cdot N'^2)$, making it more scalable and lightweight than a traditional transformer encoder.

memory consumption and computational cost, in terms of GPU memory (gMem.) and Giga floating point operations per second (GFLOPs). We observe that when the path length N increases from 50 to 200 edges, the Transformer-based method performs poorly. A method that scales better w.r.t. N is desirable.

Very large model size. Many existing PRL models have large numbers of parameters, which restricts their use in resource-limited environments. For example, in a Transformer-based method [23], where the Transformer stacks L transformer layers, each layer employs multi-head (i.e., M heads) attentions. Thus, the Transformer functions like a large cuboid, with a complexity of $\mathcal{O}(L \cdot M \cdot N^2)$, as shown in Figure 3a. For example, Table 1 shows the numbers of parameters of Transformer-based path encoders when varying the number of layers among 12, 24, 48, and 96 while fixing the number heads at 8 per layer and the feature dimension of the encoder at 512. We observe that the model parameters grow dramatically when the number of encoder layers increase, preventing the models from being deployed in resource-limited environments.

Moreover, models with large amounts of parameters also suffer from high storage and computational costs, which is not eco-friendly. More specifically, as shown in Figure 2a, for path length $N = 200$, the Transformer-based model consumes almost 3.4GiB GPU memory.

Proposed Solution. To tackle the above limitations, we propose *LightPath*, a lightweight and scalable path representation learning approach. To address the first limitation, we first propose a sparse auto-encoder targeting good scalability, w.r.t., path length. In particular, the introduction of sparseness reduces the path length from N to N' by removing edges and returning a sparse path of length N' . The sparse path is fed into a Transformer-based encoder, which reduces the complexity from $\mathcal{O}(L \cdot M \cdot N^2)$ to $\mathcal{O}(L \cdot M \cdot N'^2)$. As shown in Figure 3b, this reduces a huge cuboid to a slimmer cuboid. To avoid information loss due to the removed edges, we connect the encoder with a decoder, with the aim of reconstructing the full path. This enables scalable yet effective unsupervised training.

To further improve the training of the sparse encoder, we add an additional training scheme based on relational reasoning. In particular, for each path p_i , we construct two distinct sparse path views, denoted as p_i^1 and p_i^2 , using different reduction ratios, e.g., removing 40% and 80%, respectively. Then, we propose a dual sparse path encoder, including the original main encoder, and an additional, auxiliary encoder. The dual sparse path encoder encodes the two path views. Thus, we achieve four path presentations PR_i^1 , PR_i^2 , $\hat{P}R_i^1$, and $\hat{P}R_i^2$ for path p_i according to the two path views and the two sparse path encoders, where PR and $\hat{P}R$ denote the representations from the main and the auxiliary

encoders, respectively. Finally, given two path representations, we train a relational reasoning network to determine whether the two path representations are from the same ‘‘relation.’’ If they are from the same path, we consider them as positive relations; otherwise, they are negative relations.

To address the second limitation, we propose a global-local knowledge distillation framework that aims to reduce the model size of the main path encoder, which not only enables use in resource-limited environments but also improves accuracy. To this end, we consider the main path encoder as a teacher, and we create a lightweight sparse encoder with fewer layers and one head as a student, further reducing a slimmer cuboid to a slim rectangle (cf. Figure (3c)). The global knowledge distillation tries to push the lightweight student to mimic the teacher from a global semantic level (i.e., path representation level), while the local knowledge distillation can push the lightweight student to mimic the edge correlations from the teacher, thus building a lightweight encoder while maintaining or even improving the accuracy of downstream tasks.

To the best of our knowledge, this is the first study that systematically targets lightweight and scalable path representation learning. The study makes four main contributions.

- *Sparse Auto-encoder.* We propose a unified sparse auto-encoder framework that provides *LightPath* with good scalability. w.r.t. path length.
- *Relational Reasoning.* We introduce relational reasoning to enable efficient sparse auto-encoder training. Specifically, we propose two types of relational reasoning objectives for accurate and efficient path representation learning. These two objectives regularize each other and yield a more effective path encoder.
- *Global-local Knowledge Distillation.* We propose a novel global-local knowledge distillation framework that enables a lightweight student sparse encoder to mimic a larger teacher sparse encoder from global and local perspectives.
- *Extensive Experiments.* We report on extensive experiments on two large-scale, real-world datasets with two downstream tasks. The results offer evidence of the efficiency and scalability of the proposed framework as compared with nine baselines.

2 Preliminaries

We first cover important concepts that underlie the paper’s proposal and then state the problem addressed.

2.1 Definitions

Definition 1 Road Network. A road network is defined as a graph $\mathbf{G} = (V, E)$, where V is a set of vertices v_i that represents road intersections and $E \subseteq V \times V$ represents a set of edges $e_i = (v_j, v_k)$ that denotes road segments.

Definition 2 Path. A path $p = \langle e_1, e_2, e_3, \dots, e_N \rangle$ is a sequence of connected edges, where $e_i \in E$ denotes an edge in path and two adjacent edges share a vertex. Next, $p.N$ denotes the length of path, i.e., the number of edges in p . We let $p.\Phi = [1, 2, 3, \dots, N]$ denote a sequence of orders of the edges in p .

Definition 3 Sparse Path. A sparse path $p' = \langle e_i \rangle_{i \in p'.\Omega}$ contains a subset of the edges in path p , where $p'.\Omega$ is a sub-sequence of $p.\Phi$.

Example. Given a path $p = \langle e_1, e_3, e_4, e_6, e_7 \rangle$ and $p.\Phi = [1, 2, 3, 4, 5]$ then path $p' = \langle e_1, e_4, e_7 \rangle$, where $p'.\Omega = [1, 3, 5]$, is one of the sparse paths for p .

Definition 4 Edge Representation. The edge representation of an edge in a road network graph is a vector in \mathbb{R}^{d_k} , where d_k is the dimensionality of the vector. For simplicity, we reuse e_i to denote an edge representation.

Definition 5 Path Representation. The path representation PR of a path p is a vector in \mathbb{R}^d .

2.2 Problem Definition

Given a set of paths $\mathbb{P} = \{p_i\}_{i=1}^{|\mathbb{P}|}$ in a road network G , scalable and efficient path representation learning aims to learn a function $SPE_\theta(\cdot)$ that can generate a generic path representation for each path $p_i \in \mathbb{P}$ without relying on labels. It can be formulated as follows.

$$PR = SPE_\theta(p_i) : \mathbb{R}^{N \times d_k} \rightarrow \mathbb{R}^d, \quad (1)$$

where PR is learned path representation, θ represents the learnable parameters for the sparse path encoder, N is the path length and d_k and d are the feature dimensions for an edge and a final path representation, respectively.

2.3 Downstream Tasks

A downstream task is a task that consumes a path representation. We consider travel time estimation and path ranking score estimation. In particular, we formulate task estimation as a regression problem and define the corresponding regression model as:

$$Reg_{task_k(\psi)}(PR_i) : \mathbb{R}^d \rightarrow \mathbb{R}, \quad (2)$$

where $task_k(\psi)$ is a learnable parameter for task k and PR_i is the learned path representation of p_i .

3 Sparse Path Encoder

3.1 Transformer based Encoder

We first introduce the Transformer based encoder due to its parallelism pipeline and effectiveness for long sequence modeling. Thus, given a sequence of edge representations $\mathbf{X}_p = \langle e_1, e_2, e_3, \dots, e_N \rangle$ for a path p . Transformer based encoder takes as input \mathbf{X}_p and returns the encoded edge representations $\mathbf{Z}_p = \langle z_1, z_2, z_3, \dots, z_N \rangle$ that capture the correlation of different edges. Especially, instead of employing a single attention function, we define multi-head attention that linearly projects the queries, keys and values into M subspaces with different, learned linear projections to d_k, d_k and d_v dimensions, respectively. Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. Then, we formulate it as:

$$\mathbf{Z}_p = \text{MultiHead}(\mathbf{X}_p) = \text{Concat}(\text{head}_1, \dots, \text{head}_M) \cdot \mathbf{W}^O, \quad (3)$$

$$\text{head}_i(\cdot) = \text{softmax} \left(\left(\mathbf{X}_p \mathbf{W}_i^Q \right) \left(\mathbf{X}_p \mathbf{W}_i^K \right)^T / \sqrt{d_k} \right) \left(\mathbf{X}_p \mathbf{W}_i^V \right), \quad (4)$$

where $\text{Concat}(\cdot, \cdot)$ represents concatenation. $\mathbf{W}_i^Q \in \mathbb{R}^{d_{model} \times d_k}$, $\mathbf{W}_i^K \in \mathbb{R}^{d_{model} \times d_k}$, $\mathbf{W}_i^V \in \mathbb{R}^{d_{model} \times d_v}$, $\mathbf{W}^O \in \mathbb{R}^{M d_v \times d_{model}}$ are projections parameter matrices for scaled dot-product attention with respect to the learnable parameter θ in *LightPaht*. M denotes number of heads. d_{model} represents the feature dimension of final output. $\mathbf{Z}_p \in \mathbb{R}^{N \times d_k}$.

Except the attention sub-layers, each of the layers in Transformer based encoder also contains a fully connected feed-forward network (FFN), which is used to each position separately and identically. This FFN consists of two linear transformations with ReLU activation in between. Specifically, we have

$$\text{FFN}(\mathbf{Z}_p) = \max(0, \mathbf{Z}_p \mathbf{W}_1^{\text{FFN}} + \mathbf{b}_1^{\text{FFN}}) \mathbf{W}_2^{\text{FFN}} + \mathbf{b}_2^{\text{FFN}}, \quad (5)$$

where $\mathbf{W}_1^{\text{FFN}}$, $\mathbf{W}_2^{\text{FFN}}$, $\mathbf{b}_1^{\text{FFN}}$, and $\mathbf{b}_2^{\text{FFN}}$ are learnable parameters of feed-forward network, and $\text{FFN}(\mathbf{Z}_p) \in \mathbb{R}^{N \times d}$.

However, Transformer based encoder suffers from poor scalability w.r.t. path length and large mode size (ref. as to Section 1). To this end, we aim to study a sparse path encoder.

3.2 Overview

Figure 4 illustrates the sparse path encoder framework, which includes a sparsity operation, a sparse path encoder, and a path reconstruction decoder. The sparsity operation takes as input a full path and returns a sparse path with respect to a reduction ratio γ . Sparse path encoder takes as input a sparse path and learnable path representation and outputs learned path representations. Next, we introduce a path reconstruction decoder to reconstruct the path, thus ensuring the learned path representation captures the entire path information.

3.3 Sparsity Operation

A path consists of a sequence of edges $p = \langle e_1, e_2, e_3, \dots, e_N \rangle$, which are the basic processing units of different sequential models. The processing times of sequential models become longer when the path gets longer. Thus, we propose a sparsity operation, which is an approach to reduce the path length from N to N' , where N' is much less than N . For simplicity, we conduct path reduction by randomly removing a subset of edges in a path based on a high reduction ratio γ (e.g., $\gamma = 0.6$). A high reduction ratio γ (the ratio for edge removal) can significantly reduce the length of each input path, thus enabling the scalability of the path. Specifically, we construct the sparsity operation as:

$$p' = f(p, \gamma) = \langle e_j \rangle_{j \in \Omega}, \quad (6)$$

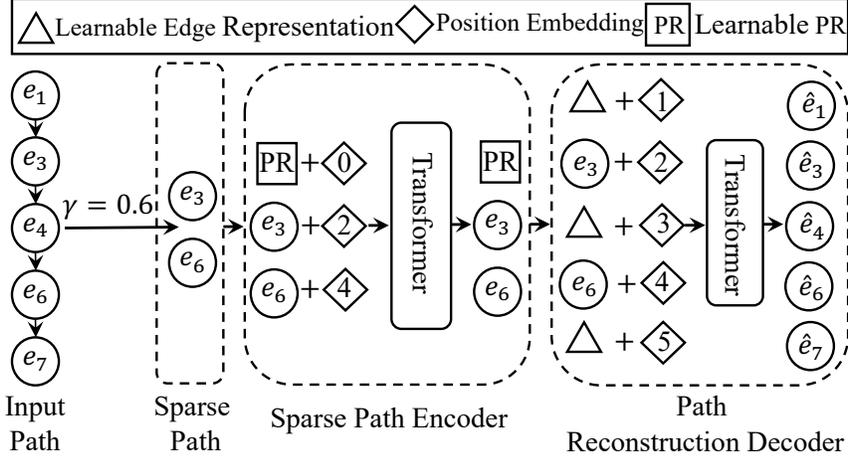


Figure 4: Sparse Auto-encoder. We remove a subset of edges from a path based on a reduction ratio γ to obtain a sparse path. We introduce a learnable path representation in front of the sparse path. And then, we fed the resulting sparse path vectors with position embeddings to a Transformer based encoder. We then introduce a learnable edge representation, denoted as a triangle, to represent the removed edges. The encoded edges in the sparse path and the removed edge representations with position embeddings are processed by a decoder that reconstructs the edges in the original path.

where p is input path. p' denotes the sparse path. For example, as shown in Figure 4, if we have a path $p = \langle e_1, e_3, e_4, e_6, e_7 \rangle$, then we conduct sparsity operation, which randomly removes a subset of edges in p , i.e., $\langle e_1, e_4, e_7 \rangle$ based on reduction ratio $\gamma = 0.6$ and achieve the sparse path $p' = \langle e_3, e_6 \rangle$ and $p'.\Omega = [2, 4]$. Thus, we can reduce path from N to N' , i.e., from 5 to 2 in this example.

3.4 Learnable Path Representation

We use Transformer as our path encoder since it processes the input edges parallelly with respect to the self-attention mechanism. In contrast, the recurrent neural network (RNN) family is inefficient due to its recurrent nature. To avoid achieving path representation through extra aggregation function [4], we add a super extra learnable path representation representative PR in front of each sparse path. Moreover, PR is attached to position 0 for every path, thus enabling it to capture global information of the path during the training procedure. Thus, we update the p' as:

$$p' = \langle PR \rangle + \langle e_j \rangle_{j \in \Omega} = \langle e_k \rangle_{k \in \Omega'}, \quad (7)$$

where $e_0 = PR$ denotes a virtual edge and $\Omega' = [0, \Omega]$.

To preserve the sequential information of the path, we add learnable position embedding into the sparse path representations based on order information in Ω' . Specifically, we have:

$$\mathbf{X}'_p = \text{Concat} \langle x_k \rangle_{k \in \Omega'}, \text{ where } x_k = e_k + pos_k, \quad (8)$$

where pos_k represents the learnable position embedding for edges in the sparse path and \mathbf{X}'_p represents the sparse path edge representation after concatenation.

Take Figure 4 as an example, we first construct $p' = \langle PR, e_3, e_6 \rangle$. Then, we add corresponding position embedding to the edge vectors of p' , i.e., positions 0, 2, and 4, where the added position embeddings can help the Transformer encoder to be aware of the input order instead of treating them as a set of unordered path edges. Meanwhile, they enable the learned path representation PR to capture global-level semantics in the sense that edges might play a different role in a road network. The intuition is that the super learnable path representation representative PR can attend attention with other edges, which captures global-level semantic. In contrast, the learnable edge representation aims to construct a full path set and reconstruct the specific edge in input path.

3.5 Transformer Path Encoder

To achieve better performance, we usually stack multiple Transformer layers, each consisting of two sub-layers: multi-head attention and position-wise feed-forward network mentioned above (ref. as to Section 3.1). Motivated

by [24], we employ a residual connection around each sub-layers, followed by layer normalization [25]. The stacked transformer model can be formulated as:

$$Z'_p = \text{LayerNorm}(\mathbf{X}'_p + \text{MultiHead}(\mathbf{X}'_p)), \quad (9)$$

$$PR = \text{LayerNorm}(Z'_p + \text{FFN}(Z'_p)), \quad (10)$$

where LayerNorm represents layer normalization and PR is learned path representation.

Remarkably, our path encoder only takes as input a small subset of edges (e.g., 60%) of the full path edges, which means we ignore the removed edges and just consider unremoved edges during the encoder stage to enable the path scalability. Path scalability enables us to train our path encoders concerning different lengths of path effectively and reduce the corresponding computational cost and memory usage.

3.6 Path Reconstruction Decoder

To capture the global information of the full path, we further introduce a lightweight path decoder to reconstruct the removed edges in a path. As shown in Figure 4, we first complement the encoded path edges and path representation with a shared, learnable vector that represents the presence of a removed edge based on the original index of each edge in a path. Then, we add the position embedding vectors to all edge representation, which enables the learnable path representation vector to capture the global information of the entire path. Next, the path decoder takes as input the full set of representations, including (1) path representation, (2) encoded unremoved edges, and (3) removed edges. We select a more lightweight decoder structure, which has less number of Transformer layers. Since the path decoder is only used to perform path reconstruction, the architecture of our path decoder can be more flexible and independent of the path encoder. Thus, the decoder is much shallower than the encoder, e.g., one layer for the decoder and 12 layers for the encoder, which significantly reduces training time. We reconstruct the input path by predicting the removed edges to ensure the learned path representation contains complete information about the entire path. We employ mean squared error (MSE) as our reconstruction loss function and compute MSE between the reconstructed and initial edge representations in the edge level. We only employ MSE loss on removed edges, which can be formulated as follows:

$$\mathcal{L}_{rec} = \frac{1}{N} \sum_{i=1}^N (e_i - \hat{e}_i)^2, \quad (11)$$

where e_i and \hat{e}_i are the initial and predicted masked edge representation, respectively. N represents the number of edges for each input path.

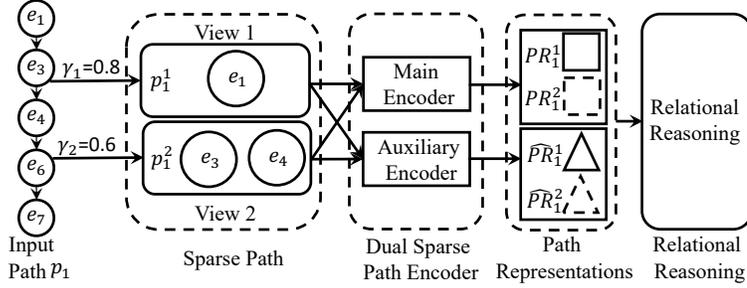
4 Relational Reasoning Path Representation Learning

4.1 Overview

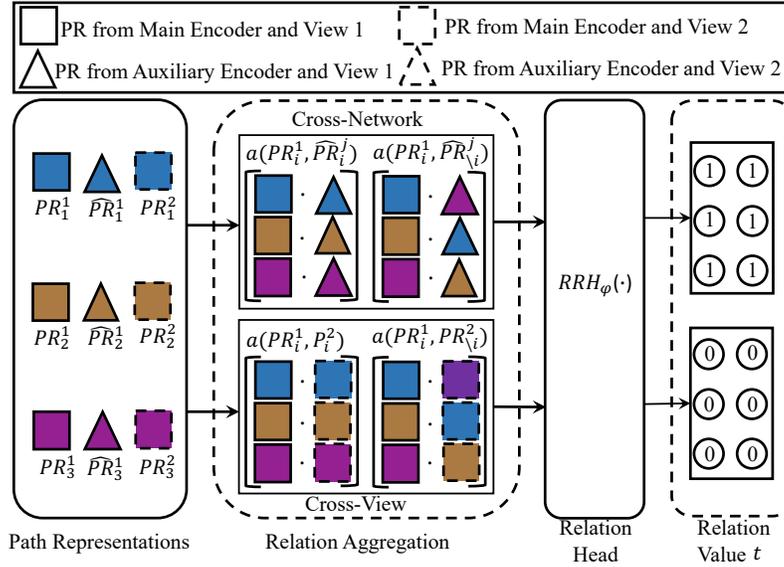
To further enhance sparse auto-encoder (cf. Section 3) training, we propose a novel self-supervised relational reasoning (RR) framework, as shown in Figure 5. The intuition behind this is that we train a relation head $RRH_\varphi(\cdot)$ to discriminate how path representations relate to themselves (same class) and other paths (different class). In particular, this framework consists of path representation construction (cf. Figure 5a) and relational reasoning (cf. Figure 5b), which includes cross-network relational reasoning and cross-view relational reasoning. To train our dual sparse auto-encoder, we first generate two path views, denotes as p_1^1 and p_1^2 , based on two different reduction ratios γ_1 and γ_2 . After this, by processing these two path views via the main encoder and the auxiliary encoder of the sparse path encoder, we construct different paths on multiple views in the representation space. Finally, we employ relational reasoning to enable efficient path representation learning.

4.2 Dual Sparse Path Encoder

In this section, we introduce our dual sparse path encoder (SPE) that is employed to generate different path representations based on different path views. As shown in Figure 5a, given a path p_1 , we first generate sparse paths in terms of two different reduction ratios γ_1 and γ_2 . We consider them as different path views, i.e., path view 1 and path view 2. Then, our dual sparse path encoder, including a main encoder and an auxiliary encoder, takes as input two different path views (i.e., p_1^1 and p_1^2) and returns different path representations. Specifically, each encoder takes as input two different path views and returns two different path representations, where solid and dotted \square denote the path representations returned from main encoder based on path view 1 and path view 2, respectively, i.e., PR_1^1 and PR_1^2 . In contrast, solid and dotted \triangle represent the path representations achieved from auxiliary encoder based on both path views, respectively,



(a) Path Representation Construction Using Dual Sparse Path Encoders



(b) Relational Reasoning

Figure 5: Illustration of RR Training: (a) Given an input path p_1 , we construct two path views (i.e., p_1^1 and p_1^2) through two reduction ratios γ_1 and γ_2 , based on which a main encoder and an auxiliary encoder are employed to generate path representations for each view (i.e., PR_1^1 , PR_1^2 , \hat{PR}_1^1 , and \hat{PR}_1^2). (b) After getting corresponding path representations for paths in a minibatch, a relational reasoning path representation learning scheme, which utilizes both cross-network and cross-view relational reasoning modules, is deployed. In particular, for both modules, an aggregation function a joins positives (representations of the same paths, e.g., $a(PR_1^1, \hat{PR}_1^1)$, $a(PR_1^1, PR_1^2)$) and negatives (randomly paired representations, e.g., $a(PR_1^1, \hat{PR}_3^1)$, $a(PR_1^1, PR_3^2)$) through a commutative operator. Then relation head module $RRH_\varphi(\cdot)$ estimates the relation score y , which must be 1 for positive and 0 for negatives. Both cross-network and cross-view objectives are optimized minimizing the binary cross-entropy (BCE) between prediction and target relation value t . In this example, $i \in [1, 2, 3]$ denotes the number of paths in the minibatch and $j \in [1, 2]$ represents the number of views.

i.e., \hat{PR}_1^1 and \hat{PR}_1^2 . To this end, we construct four different path representations for a given path, which promote our design of cross-network relational reasoning and cross-view relational reasoning in turn. Finally, we formulate it as:

$$PR_i^j = SPE_\theta(p_i^j, \gamma), \hat{PR}_i^j = SPE_{\hat{\theta}}(p_i^j, \gamma), \quad (12)$$

where PR_i^j and \hat{PR}_i^j are path representations obtained from the main encoder and the auxiliary encoder, respectively. p_i denotes the i -th path in the path set. $j \in [1, 2]$ denotes the path views. θ and $\hat{\theta}$ are the parameters for the main encoder and auxiliary encoder.

4.3 Relational Reasoning

4.3.1 Cross-Network Relational Reasoning

In *LightPath*, we employ a dual sparse path encoder, which includes main and auxiliary encoder, as shown in Figure 5a. We first construct path representations through sparsity operation based on different reduction ratios γ_1 and γ_2 . Given a set of path $\{p_1, p_2, \dots, p_K\}$, we can have a set of path representations $\{PR_1^1, PR_2^1, \dots, PR_K^1\}$ from main encoder and $\{\hat{PR}_1^1, \hat{PR}_2^1, \dots, \hat{PR}_K^1\}$ or $\{\hat{PR}_1^2, \hat{PR}_2^2, \dots, \hat{PR}_K^2\}$ from auxiliary encoder by using path representation construction. Then we employ relation aggregation $a(\cdot)$ that joins the positive path representation relations $\langle PR_i^1, \hat{PR}_i^1 \rangle$ or $\langle PR_i^1, \hat{PR}_i^2 \rangle$ and the negative path representation relations $\langle PR_i^1, \hat{PR}_{\setminus i}^1 \rangle$, where i denotes the i -th path sample and $\setminus i \neq i$ represents randomly selected path representations in a minibatch. Take Figure 5b as an example, where $K = 3$. we join $\langle PR_1^1, \hat{PR}_1^1 \rangle$ as a positive relation pair (representation from same path), and $\langle PR_1^1, \hat{PR}_2^1 \rangle$ as a negative relation pair (representation from different paths) through aggregation function a . Next, the relational head $RRH_\varphi(\cdot)$, which is non-linear function approximator parameterized by φ , takes as input representation relation pairs of cross-network and returns a relation score y . Finally, we formulate the cross-network relational reasoning task as a binary classification task, where we use binary cross-entropy loss to train our sparse path encoder, which is given as follows.

$$\mathcal{L}_{cn} = \underset{\theta, \varphi}{\operatorname{argmin}} \sum_{i=1}^K \sum_{j=1}^2 \mathcal{L} \left(RRH_\varphi \left(a \left(PR_i^j, \hat{PR}_i^j \right) \right), t = 1 \right) + \mathcal{L} \left(RRH_\varphi \left(a \left(PR_i^j, \hat{PR}_{\setminus i}^j \right) \right), t = 0 \right), \quad (13)$$

where K is the the number of path samples in the minibatch. $a(\cdot, \cdot)$ is an aggregation function. \mathcal{L} is a loss function between relation score and a target relation value. t is a target relation values.

The intuition behind this is to discriminate path presentations of same path and different paths, which are from different views across dual sparse path encoder and are able to distill the knowledge from historical observations, as well as stabilizing the main encoder training. To realize this, we adopt Siamese architecture for our dual sparse path encoder, where the auxiliary encoder does not directly receive the gradient during the training procedure. In contrast, we update its parameters by leveraging the momentum updating principle:

$$\hat{\theta}^t = m \times \hat{\theta}^{(t-1)} + (1 - m) \times \theta^t, \quad (14)$$

where m is momentum parameters. θ and $\hat{\theta}$ are the parameters of the main encoder and the auxiliary encoder.

4.3.2 Cross-View Relational Reasoning

To enhance the learning ability of our *LightPath*, we further consider the ties between two views within main encoder, which acts as a strong requarization to enhance the learning ability of our methods. We do not have to include such relational reasoning within the auxiliary encoder because it will not directly compute gradient during training, and our goal is to train main encoder. Figure 5b shows the design of our cross-view relational reasoning, which contains two similar representations from two views based on γ_1 and γ_2 . The intuition of cross-view relational reasoning is to preserve the relation between two views of the same path and discriminate them from the view of other paths.

Similar with cross-network, given a set of paths $\{p_1, p_2, \dots, p_K\}$. We first achieve two set of path representations in terms of two path views based on main encoder, i.e., $\{PR_1^1, PR_2^1, \dots, PR_k^1\}$ and $\{PR_1^2, PR_2^2, \dots, PR_K^2\}$. Then, we join the positive relation pairs (e.g., $\langle PR_i^1, PR_i^2 \rangle$) and negative relation pairs (e.g., $\langle PR_i^1, PR_{\setminus i}^2 \rangle$) through aggregation function. For example, as shown in Figure 5b, there are 3 paths in the set. Thus, we can denote $\langle PR_1^1, PR_1^2 \rangle$ as a positive pair and $\langle PR_1^1, PR_3^2 \rangle$ as a negative pair. Then, we further employ relational head $RRH_\varphi(\cdot)$, which takes as input a positive pair and a negative pairs from different views, to achieve the corresponding relation score y for the cross-view relational reasoning. Last, we formulate the cross-view relational reasoning loss to discriminate how different views of a path is related to themselves and other paths. In this phase, the complete learning objective can be specified as:

$$\mathcal{L}_{cv} = \underset{\theta, \varphi}{\operatorname{argmin}} \sum_{i=1}^K \mathcal{L} \left(RRH_\varphi \left(a \left(PR_i^1, PR_i^2 \right) \right), t = 1 \right) + \mathcal{L} \left(RRH_\varphi \left(a \left(PR_i^1, PR_{\setminus i}^2 \right) \right), t = 0 \right), \quad (15)$$

where K is the the number of path samples in the minibatch.

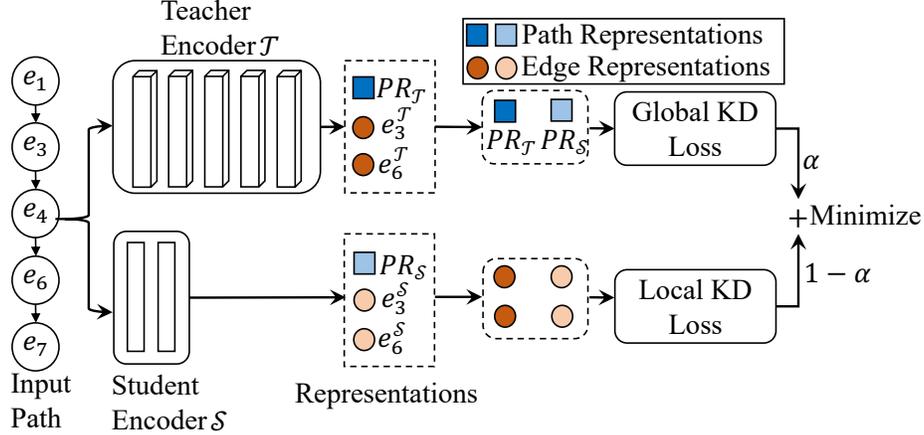


Figure 6: Illustration of GLKD. Given an input path, we formulate our GLKD as a weighted sum of global path representation knowledge distillation (GPRKD) loss and local edge representation knowledge distillation (LERKD) loss.

4.3.3 Objective for RR

To train our dual path encoder end-to-end and efficient learn path representations for downstream tasks, we jointly leverage both the cross-network and cross-view relation reasoning loss. Specifically, the overall objective function is formulated as Eq. 16.

$$\min_{\theta, \varphi} \mathcal{L}_{RR} = \mathcal{L}_{cn} + \mathcal{L}_{cv} \quad (16)$$

4.4 LightPath Training

To train our sparse path encoder and learn path representations for downstream tasks, we jointly minimize the reconstruction and RR loss. Specifically, the overall objective function is defined as:

$$\mathcal{L} = \mathcal{L}_{rec} + \mathcal{L}_{RR} \quad (17)$$

5 Global Local Knowledge Distillation (GLKD)

So far, we realize our *LightPath* through sparse auto-encoder and relational reasoning and transform it from large cuboid (cf. Figure 3a) to a slim cuboid (cf. Figure 3b). To enable the *LightPath* that can be deployed on resource-limited mobile devices, we introduce our global-local knowledge distillation (GLKD) to further reduce the size of the sparse auto-encoder, as shown in Figure 6. We first train a large cuboid teacher encoder with multiple transformer layers and heads (cf. Figure 3b) based on path reconstruction (cf. Section 3) and relational reasoning (cf. Section 4). Then, we employ a small rectangle student encoder (cf. Figure 3c), which has less layers and heads, to mimic a large teacher model and use the teacher’s knowledge to obtain similar or superior performance based on GLKD. Specifically, GLKD constructs a local knowledge distillation by matching the representations of correlated edges. On such a basis, the global term distills the knowledge from teacher to student that enabling the informative and powerful path representation for the student model.

5.1 Global-path Representation Distillation

Given a path $p_i = \langle e_1, e_2, e_3, \dots, e_N \rangle$, where N is the number of edges in a path. We define $PR^T(p_i)$ and $PR^S(p_i)$ represent the path representations achieved from teacher encoder \mathcal{T}_θ and student encoder \mathcal{S}_θ . The intuition of global path representation knowledge distillation is to let the student encoder mimic the global properties captured by a large cuboid teacher encoder. And thus, the goal of global path representation knowledge distillation is to put closer the path representation from teacher encoder and student encoder in the latent space. We formalize this problem as minimizing a latent space distance representation pairs in terms of the large cuboid teacher encoder and the rectangle student encoder.

The formulation of the objective function is given as follows:

$$\min_{\theta} \mathcal{L}_{global} \left(PR^T(p_i), PR^S(p_i) \right) = \left\| sp(PR^T(p_i)/t) - sp(PR^S(p_i)/t) \right\|^2, \quad (18)$$

where $sp(\cdot)$ is exponential function. t denotes the temperature. Using a higher value for t produces a softer probability distributions over path representations.

5.2 Local-edge Correlation Distillation

The goal of local-edge structure distillation is to preserve the local similarity of the edge correlations in a path. In particular, it is expected that the representation of the same edge in a path represented by the teacher encoder and the student encoder should be close to each other. The intuition is that a rectangle student encoder mimics the edge correlations in a path captured by a large cuboid teacher encoder. Using a similarity measurement, we formulate the local-edge structure distillation problem as minimizing the latent space distance of edge representations from the teacher encoder and then student encoder.

In specific, given a path $p = \langle e_1, e_2, e_3, \dots, e_N \rangle$ in a road network, where N is the number of edges in a path. Through applying an L -layers Transformer encoder (i.e., teacher encoder \mathcal{T}_{θ}) and L' -layers Transformer encoder (i.e., student encoder \mathcal{S}_{θ}) upon sparse path p' , where $L \ll L'$, the edge representation that captures spatial dependencies are derived as follows.

$$F^T(e_i)_{i=1}^{N'} = \mathcal{T}_{\theta}(p), F^S(e_i)_{i=1}^{N'} = \mathcal{S}_{\theta}(p), \quad (19)$$

where $F^T(e_i)_{i=1}^{N'}$ and $F^S(e_i)_{i=1}^{N'}$ represent the edge representation with respect to the teacher and student encoder, respectively.

In this phase, the goal of learning is to reduce the latent space distance between same edge pair from the teacher and student encoder, respectively. To this end, we aim to minimize the following objective functions between edge representation pairs in terms of the parameters of the student encoder.

$$\min_{\theta} \mathcal{L}_{local} \left(\mathbf{F}^T(e_i), \mathbf{F}^S(e_i) \right) = \frac{1}{n} \sum_{i=1}^n \left\| sp(\mathbf{F}^T(e_i)/t) - sp(\mathbf{F}^S(e_i)/t) \right\|^2, \quad (20)$$

where $sp(\cdot)$ represents exponential function. t denotes the temperature. Using a higher value for t produces a softer probability distributions over edges.

5.3 Objective for GLKD

To train our global and local knowledge distillation in an end-to-end fashion, we jointly leverage both the global and local knowledge distillation loss. Specifically, the overall objective function to minimize is defined in Eq. 21.

$$\min_{\theta} \mathcal{L}_{GLKD} = \alpha * \mathcal{L}_{global} + (1 - \alpha) * \mathcal{L}_{local}, \quad (21)$$

where α is balancing factor.

6 Experiments

6.1 Experimental Setup

6.1.1 Datasets

We conduct experiments on two real-world datasets and one synthetic dataset to enable fair comparisons with existing studies. Based on two real-world datasets, we report results for two downstream tasks: travel time estimation [26, 6], and path ranking [19, 27, 18]. Due to the lack of large amounts of long paths in the real-world datasets, we construct one synthetic dataset that contains paths with lengths of 100, 150, and 200 to verify the efficiency and scalability of *LightPath*.

Aalborg, Denmark: We collect the road network of Aalborg from OpenStreetMap² that contains 10,017 nodes and 11,597 edges. Specifically, this dataset contains 180 million GPS records from 183 vehicles sampled at 1 Hz over a two-year period from 2017 to 2018. After map matching [28], we obtain 39,160 paths with length 50.

²<https://www.openstreetmap.org>

Chengdu, China: This dataset was collected from Chengdu, China, on October and November 2016. We obtain the corresponding road network from OpenStreetMap. The network contains 6,632 nodes and 17,038 edges. The GPS data was sampled at about 1/4-1/2 Hz. We obtain 50,000 paths through map matching with lengths 50.

Synthetic: Due to the lack of large amounts of long paths in the real-world datasets, we construct one synthetic dataset to verify the efficiency and scalability of *LightPath*. In particular, we first randomly pick 500 nodes in road network of Aalborg dataset, and then expand each node to a path by random walking until the path length reach the threshold (i.e., 100, 150, 200), which we refer as to generation process. Subsequently, we iterate the generation process 10 times to construct 5,000 paths for each path length.

6.1.2 Downstream Tasks

We report the results on two downstream tasks:

Path Travel Time Estimation: We obtain travel time (in seconds) for each path from the trajectory. We aim to utilize a regression model to predict the travel time based on the learned path representations. We employ Mean Absolute Error(MAE), Mean Absolute Relative Error(MARE), and Mean Absolute Percentage Error(MAPE) to evaluate the performance of travel time estimations. The smaller values of these metrics, the better performance we achieve.

Path Ranking: Each path is assigned a ranking score in the range $[0, 1]$, which is obtained from historical trajectories by following the existing studies [18, 27]. More specifically, we take the path that is used by a driver in the historical trajectories as the trajectory path, which is denoted as the top ranked path. Then, we generate multiple paths connecting the same source and destination via path finding algorithms [29]. Finally, we calculate the similarity between a generated path and the trajectory path as a ranking score. The higher ranking score indicates a generated path is more similar to the trajectory path, and the trajectory path itself has a score of 1 and ranks the highest. To measure the path ranking, we apply MAE, the Kendall rank correlation coefficient (τ), and the Spearman’s rank correlation coefficient (ρ), which are widely used metrics in path ranking, to evaluate the effectiveness of path ranking.

6.1.3 Models for Downstream Tasks

For all unsupervised learning methods, we first achieve the corresponding d dimensionality path representation and then we build a regression model that takes as input a path representation and output estimated the travel time and path ranking, respectively. In particular, we select ensemble model *Gradient Boosting Regressor*(GBR) [30] as our prediction model since they are regression problems.

6.1.4 Baselines

We compare *LightPath* with 9 baselines, including 6 unsupervised learning-based methods and 3 supervised learning-based methods. The details of these baseline methods are summarized as follows:

- **Node2vec** [31] is an unsupervised node representation model that learn node representation in a graph. We achieve the path representation by aggregating the node representations of the nodes in a path.
- **MoCo** [32] is a momentum contrast for unsupervised visual representation learning. Here we use momentum contrast to learn path representations.
- **Toast** [23] first uses auxiliary traffic context information to learn road segment representation based on the skip-gram model and then utilizes a stacked transformer encoder layer to train trajectory representation through route recovery and trajectory discrimination tasks. We use the same schema to learn path representations.
- **t2vec** [33] is a trajectory representation learning method for similarity computation based on the encoder-decoder framework, which is trained to reconstruct the original trajectory. We use a sequence of edges in a path to represent a trajectory.
- **NeuTraj** [34] is a method that revised the structure of LSTM to learn representations of the grid in the process of training their framework. To support our task with it, we replace the grid with edges in their framework.
- **PIM** [4] is an unsupervised path representation learning approach that first generates negative samples using curriculum learning and then employs global and local mutual information maximization to learn path representations.
- **HMTRL** [19] is a supervised path representation learning framework for multi-modal transportation recommendation.
- **PathRank** [18] is a supervised path representation learning model based on GRUs, which treats departure time and driver ID as additional information.

- **CompactETA** [35] aims to estimate travel time based on a real time predictor and an asynchronous updater.
- **HierETA** [36] is a supervised multi-view trajectory representation method to estimate the travel time.
- **LightPath-Sup** is a supervised version of our *LightPath*, where we train it in a supervised manner.

For all unsupervised learning methods, we first use unlabeled training data (e.g., 30K unlabeled Aalborg dataset and 50K unlabeled Chengdu dataset) to train path encoders to obtain path representations. Then a regression model takes as input path representations and returns the estimated travel time and path ranking score using a limited labeled training dataset, e.g., the 12K labeled Aalborg dataset. In comparison, for all supervised learning methods, we directly train path encoders using a limited labeled training dataset, e.g., the 12K labeled Aalborg dataset and Chengdu dataset.

6.1.5 Implementation Details

We employ an asymmetrical sparse auto-encoder architecture and randomly initialize all learnable parameters with uniform distributions. In particular, we adopt Siamese architecture, where we update the parameters of the auxiliary encoder based on the momentum updating principle based on the main encoder and we set the momentum parameter $m = 0.99$. We employ node2vec [31] to embed each edge to 128-dimensional vectors and set the dimension for path representation to 128. For a fair comparison, we set the path representation dimensionality of all baseline methods as 128. We select concatenate as the relation aggregation function $a(\cdot, \cdot)$. We use the AdamW optimizer with a cosine decay learning rate schedule over 400 epochs, with a warm-up period of 40 epochs. We set the base learning rate to $1e-3$ and betas as $(0.9, 0.95)$. We vary γ from 0.1, 0.3, 0.5, 0.7, 0.9 to study the effect of path scalability and efficiency for the *LightPath*. In addition, we consider four different path lengths, i.e., 50, 100, 150, and 200, to study the effectiveness, efficiency, and scalability of the *LightPath*. We then evaluate our *LightPath* as well as all baselines on a powerful Linux server with 40 Intel(R) Xeon(R) W-2155 CPU @ 3.30GHz and two TITAN RTX GPU cards. Finally, all algorithms are implemented in PyTorch 1.11.0.

6.2 Experimental Results

6.2.1 Overall Performance

Table 2 shows the overall performance of our *LightPath* and all the compared baselines on both datasets in terms of different evaluation metrics. Especially, we select 30K unlabeled paths on Aalborg and Chengdu, respectively, but we only have 12K labeled paths for both datasets. Thus, we use 30K unlabeled paths to train path encoder for

Table 2: Overall Accuracy on Travel Time Estimation and Ranking Score Estimation

| Method | Aalborg | | | | | | Chengdu | | | | | |
|----------------------|------------------------|-------------|--------------|--------------|-------------|-------------|------------------------|-------------|--------------|--------------|-------------|-------------|
| | Travel Time Estimation | | | Path Ranking | | | Travel Time Estimation | | | Path Ranking | | |
| | MAE | MARE | MAPE | MAE | τ | ρ | MAE | MARE | MAPE | MAE | τ | ρ |
| <i>Node2vec</i> | 154.07 | 0.20 | 25.22 | 0.24 | 0.59 | 0.64 | 267.28 | 0.23 | 26.30 | 0.15 | 0.74 | 0.77 |
| <i>MoCo</i> | 146.29 | 0.19 | 21.60 | 0.25 | 0.53 | 0.57 | 237.14 | 0.20 | 23.13 | 0.15 | 0.77 | 0.81 |
| <i>Toast</i> | 137.27 | 0.17 | 20.43 | 0.24 | 0.59 | 0.63 | 240.57 | 0.21 | 23.50 | 0.11 | 0.65 | 0.68 |
| <i>t2vec</i> | 147.24 | 0.19 | 22.13 | 0.25 | 0.52 | 0.56 | 242.96 | 0.21 | 23.65 | 0.14 | 0.77 | 0.82 |
| <i>NeuTraj</i> | 117.06 | 0.15 | 18.09 | 0.25 | 0.60 | 0.64 | 232.96 | 0.20 | 22.73 | 0.12 | 0.79 | 0.83 |
| <i>PIM</i> | 102.09 | 0.14 | 14.92 | 0.20 | 0.63 | 0.67 | 223.34 | 0.19 | 21.69 | 0.12 | 0.80 | 0.84 |
| <i>HMTRL</i> | 101.81 | 0.13 | 14.51 | 0.17 | 0.68 | 0.72 | 218.94 | 0.19 | 21.22 | 0.09 | 0.83 | 0.84 |
| <i>PathRank</i> | 115.37 | 0.15 | 16.41 | 0.21 | 0.64 | 0.68 | 229.85 | 0.20 | 22.53 | 0.11 | 0.81 | 0.82 |
| <i>CompactETA</i> | 106.47 | 0.15 | 16.22 | 0.17 | 0.67 | 0.70 | 236.28 | 0.20 | 23.13 | 0.11 | 0.79 | 0.80 |
| <i>HierETA</i> | 88.95 | 0.12 | 14.23 | 0.15 | 0.71 | 0.74 | 215.39 | 0.19 | 21.66 | 0.09 | 0.84 | 0.85 |
| <i>LightPath-Sup</i> | 105.51 | 0.15 | 16.35 | 0.14 | 0.68 | 0.72 | 218.67 | 0.19 | 21.36 | 0.13 | 0.76 | 0.79 |
| <i>LightPath</i> | 85.76 | 0.11 | 12.12 | 0.13 | 0.73 | 0.77 | 212.61 | 0.18 | 20.75 | 0.07 | 0.87 | 0.88 |

unsupervised-based methods. However, supervised approaches can only use the 12K labeled paths. Overall, *LightPath* outperforms all the baselines on these two tasks for both datasets, which demonstrates the advance of our model. Specifically, we can make the following observations. Graph representation learning based approach *Node2vec* is much worse than *LightPath*. This is because *Node2vec* fails to capture spatial dependencies in a path. In contrast, *LightPath* considers the spatial dependencies through the self-attention mechanism, thus achieving better performance.

Although *MoCo* considers the dependencies among edges in a path, this method still performs worse. The main reason is that *MoCo* can leverage the spatial dependencies, but it converges very slow since it needs large amounts of

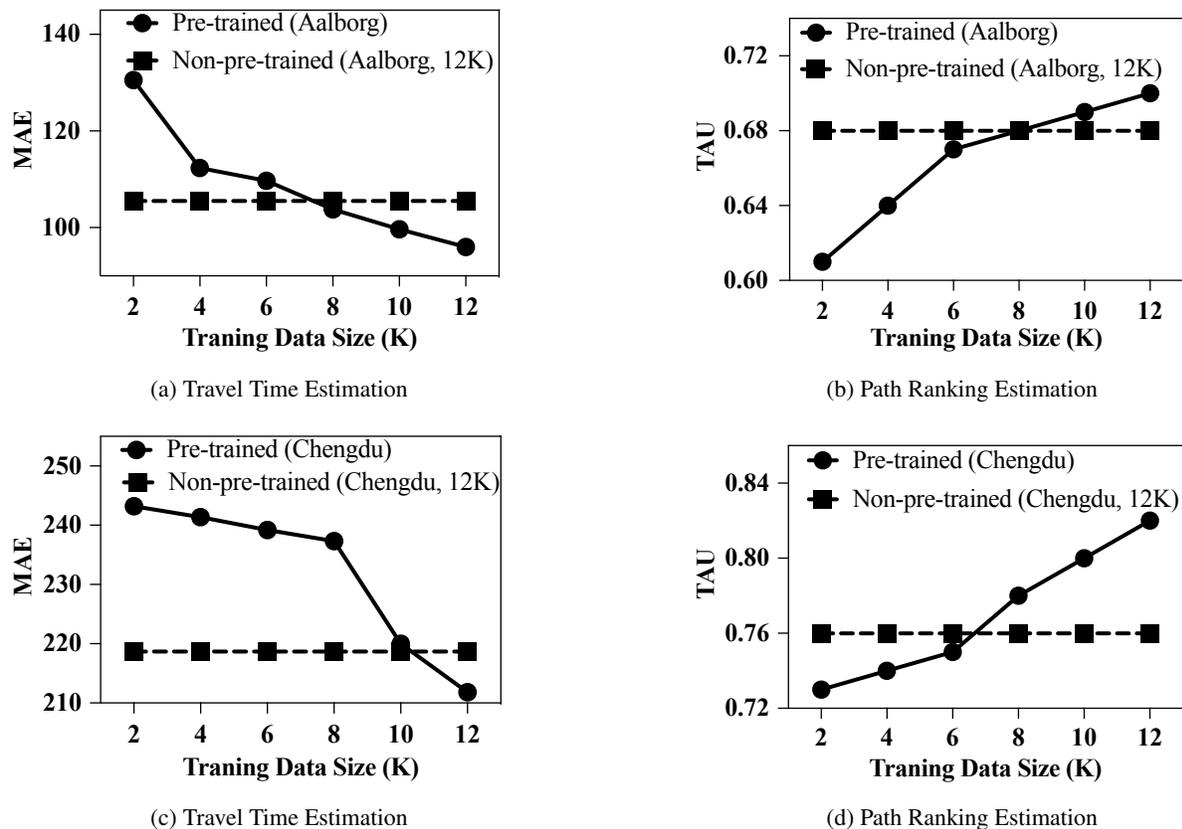


Figure 7: Effects of Pre-training.

negative samples to enable training. *LightPath* also outperforms *t2vec* and *NeuTraj*, which both are first design to learn trajectory representation for trajectory similarity computation. This suggests that random drops on some edges and not reconstruct these edges in a path resulting in spatial information missing, thus achieving the worse performance on downstream tasks. *PIM* consistently outperforms all other unsupervised baselines, which demonstrates the effectiveness of representation learning. The main reason is that *PIM* is designed for path representation learning. However, *PIM* is InfoNCE based method and has high computation complexity, making it hard to deploy on resource-limited edge devices. *HMTRL*, *PathRank* and *LightPath-Sup* are three supervised learning methods that achieve relatively worse performance due to the lack of labeled training data. Since labeling data is very time-consuming and expensive. We consider a scenario where labeled data is limited in this paper.

6.2.2 Using *LightPath* as Pre-training Methods

Model pre-training aims to create generic representations that can then be fine-tuned for multiple downstream tasks using a limited labeled dataset. Especially for many personalized services (e.g., personalized travel time estimation or path ranking estimation) in transportation applications, a user can use his/her trajectory to fine-tune the pre-trained model and then achieve the personalized service, where we do not have many trajectories from the user. In this experiment, we evaluate the effect of Pre-training. We employ *LightPath* as a pre-training method for the supervised method *LightPath-Sup*. Specifically, we first train *LightPath* in an unsupervised fashion, and then we use the learned transformer path encoder to initialize the transformer in *LightPath-Sup*. Here, it takes as input a sequence of edge representations and predicts the travel time and path ranking score. Figure 7 illustrates the performance of *LightPath-Sup* w. and w/o pre-training over two downstream tasks on both datasets. When employing non-pre-trained *LightPath-Sup*, we train it using 12K labeled training paths. We notice that (1) when employing pre-training, we can obtain the same performance with no-pre-trained *LightPath-Sup* using less labeled data. For example, *LightPath-Sup* w. pre-training only needs 8K, and 10K labeled training paths for the Aalborg and Chengdu, respectively, to achieve the same performance of *LightPath-Sup* w/o pre-training with 12k labeled samples on the task of travel time estimation. (2) *LightPath-Sup* w.

Table 3: Effect of Variants of *LightPath*

| Method | Aalborg | | | | | | Chengdu | | | | | |
|-------------------------|------------------------|-------------|--------------|--------------|-------------|-------------|------------------------|-------------|--------------|--------------|-------------|-------------|
| | Travel Time Estimation | | | Path Ranking | | | Travel Time Estimation | | | Path Ranking | | |
| | MAE | MARE | MAPE | MAE | τ | ρ | MAE | MARE | MAPE | MAE | τ | ρ |
| <i>w/o RR</i> | 90.90 | 0.12 | 13.85 | 0.17 | 0.66 | 0.70 | 224.31 | 0.19 | 21.90 | 0.14 | 0.76 | 0.79 |
| <i>w/o Rec.</i> | 103.45 | 0.14 | 15.76 | 0.15 | 0.65 | 0.69 | 229.24 | 0.20 | 22.36 | 0.16 | 0.69 | 0.73 |
| <i>w/o ME</i> | 91.57 | 0.12 | 13.09 | 0.16 | 0.68 | 0.72 | 223.81 | 0.19 | 21.86 | 0.13 | 0.78 | 0.81 |
| <i>w/o CN</i> | 93.17 | 0.12 | 13.35 | 0.15 | 0.68 | 0.73 | 217.14 | 0.18 | 21.29 | 0.08 | 0.80 | 0.83 |
| <i>w/o CV</i> | 89.84 | 0.12 | 13.51 | 0.15 | 0.68 | 0.72 | 215.59 | 0.18 | 21.20 | 0.09 | 0.81 | 0.83 |
| <i>LightPath</i> | 85.76 | 0.11 | 12.12 | 0.13 | 0.73 | 0.77 | 212.61 | 0.18 | 20.75 | 0.07 | 0.87 | 0.88 |

Table 4: Effect of KD, Global Loss and Local Loss

| Method | Aalborg | | | | | | Chengdu | | | | | |
|-------------------------|------------------------|-------------|--------------|--------------|-------------|-------------|------------------------|-------------|--------------|--------------|-------------|-------------|
| | Travel Time Estimation | | | Path Ranking | | | Travel Time Estimation | | | Path Ranking | | |
| | MAE | MARE | MAPE | MAE | τ | ρ | MAE | MARE | MAPE | MAE | τ | ρ |
| <i>w/o KD</i> | 87.77 | 0.11 | 12.94 | 0.14 | 0.70 | 0.74 | 213.26 | 0.18 | 20.97 | 0.08 | 0.84 | 0.86 |
| <i>w/o Global</i> | 90.24 | 0.12 | 13.31 | 0.18 | 0.67 | 0.71 | 220.32 | 0.19 | 21.52 | 0.09 | 0.79 | 0.81 |
| <i>w/o Local</i> | 89.23 | 0.12 | 12.78 | 0.16 | 0.69 | 0.73 | 215.03 | 0.19 | 21.02 | 0.08 | 0.82 | 0.84 |
| <i>LightPath</i> | 85.76 | 0.11 | 12.12 | 0.13 | 0.73 | 0.77 | 212.61 | 0.18 | 20.75 | 0.07 | 0.87 | 0.88 |

pre-training achieves higher performance than *LightPath-Sup* w/o pre-training. We observe similar results on the task of path ranking, demonstrating that *LightPath* can be used as a pre-training method to enhance supervised methods.

6.2.3 Ablation Studies

To verify the effectiveness of different components in *LightPath*, we conduct ablation studies on *LightPath*: a) effect of variants of *LightPath*, specifically reconstruction (Rec) loss, relational reasoning (RR) loss, cross-network loss and cross-view loss; b) effect of global-local knowledge distillation.

a) Effect of variants of *LightPath*, we consider five variants of *LightPath*: 1) *w/o RR*; 2) *w/o Rec.*; 3) *w/o ME*; 4) *w/o CN*; 5) *w/o CV*. In *w/o RR*, we only consider path reconstruction loss and use main encoder; In *w/o Rec.*, we only consider relational reasoning loss; In *w/o ME*, we consider both path reconstruction and relational reasoning losses, but we do not consider Siamese architectures in dual path encoder; In *w/o CN*, we remove the cross-network loss in RR; And in *w/o CV*, we remove cross-view loss in RR. Table 3 reports the results on both dataset. We can observe that (1) *LightPath w/o Rec.* achieves the worst performance because the learned *PR* only capture information from sparse path while ignoring the removed edges, which verifies the importance of path reconstruction decoder; (2) *LightPath w/o RR* also achieves the poor performance, which implies the effectiveness of self-supervised relational reasoning. (3) We observe that the performance of *LightPath* degrades without cross-network and cross-view loss on both datasets, which further demonstrates the effectiveness of our relational reasoning loss. (4) We notice that *LightPath* achieves the best performance. This result implies that all the proposed modules contribute positively to the final performance, which validates that *LightPath* takes advantage of all designed components.

b) Effect of KD, global KD loss, local KD loss: We further study the effect of global-local knowledge distillation. We compared our framework with three variants: 1) *w/o KD*, which denotes the performance of the teacher model; 2) *w/o global KD loss*, which removes global loss from global-local knowledge distillation; and 3) *w/o local KD loss*, which removes local loss from global-local knowledge distillation. As shown in Table 4, compared with *KD*, *LightPath* achieves a better performance, which verifies that the teacher model can improve the performance of the student model. Both global and local loss can improve the performance of the learned path representation of the student model. In specific, global loss makes more contributions to the learned path representations. As a result, removing global loss degrades performance significantly.

6.2.4 Parameter Sensitivity Analysis

We proceed to study three important hyper-parameters, including 1) model scalability w.r.t. reduction ratio and path length, 2) Effect of Reduction Ratio γ , 3) the parameter of temperature for global-local knowledge distillation, and 4) effect of balancing factor α .

Table 5: Model Scalability vs. Reduction Ratio (γ) and Path Length (N)(Aalborg)

| LightPath | | | | | | | | | | | | | |
|-------------------------|--------------|-------|----------------|-------|----------------|-------|----------------|-------|----------------|-------|----------------|-------|-------|
| N | $\gamma = 0$ | | $\gamma = 0.1$ | | $\gamma = 0.3$ | | $\gamma = 0.5$ | | $\gamma = 0.7$ | | $\gamma = 0.9$ | | Para. |
| | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | |
| 50 | 8.01 | 1.47 | 7.48 | 1.39 | 6.44 | 1.37 | 5.39 | 1.36 | 4.34 | 1.34 | 3.18 | 1.33 | 1.570 |
| 100 | 15.77 | 1.60 | 14.72 | 1.50 | 12.62 | 1.48 | 10.52 | 1.46 | 8.43 | 1.44 | 6.23 | 1.43 | 1.570 |
| 150 | 23.53 | 1.72 | 21.95 | 1.64 | 18.81 | 1.60 | 15.66 | 1.58 | 12.52 | 1.55 | 9.27 | 1.52 | 1.570 |
| 200 | 31.29 | 1.90 | 29.19 | 1.81 | 24.99 | 1.77 | 20.80 | 1.73 | 16.61 | 1.68 | 12.31 | 1.65 | 1.570 |
| LightPath w/o KD | | | | | | | | | | | | | |
| N | $\gamma = 0$ | | $\gamma = 0.1$ | | $\gamma = 0.3$ | | $\gamma = 0.5$ | | $\gamma = 0.7$ | | $\gamma = 0.9$ | | Para. |
| | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | GFLOPs | gMem. | |
| 50 | 33.68 | 1.78 | 30.64 | 1.70 | 22.55 | 1.61 | 18.47 | 1.53 | 12.39 | 1.47 | 5.70 | 1.41 | 5.525 |
| 100 | 66.60 | 2.53 | 60.52 | 2.37 | 48.36 | 2.11 | 36.19 | 1.86 | 24.03 | 1.72 | 11.26 | 1.58 | 5.525 |
| 150 | 99.53 | 3.44 | 90.41 | 3.23 | 72.16 | 2.76 | 53.91 | 2.35 | 35.65 | 2.03 | 16.82 | 1.82 | 5.525 |
| 200 | 132.54 | 4.74 | 120.29 | 4.30 | 95.96 | 3.53 | 71.64 | 2.94 | 47.31 | 2.43 | 22.37 | 2.10 | 5.525 |

Table 6: Effect of Reduction Ratio γ

| γ | Aalborg | | | | | |
|----------|------------------------|-------------|--------------|--------------|-------------|-------------|
| | Travel Time Estimation | | | Path Ranking | | |
| | MAE | MARE | MAPE | MAE | τ | ρ |
| 0.1 | 82.79 | 0.11 | 11.95 | 0.12 | 0.74 | 0.77 |
| 0.3 | 84.75 | 0.11 | 12.14 | 0.13 | 0.73 | 0.77 |
| 0.5 | 84.81 | 0.11 | 11.86 | 0.14 | 0.72 | 0.76 |
| 0.7 | 85.91 | 0.11 | 12.49 | 0.14 | 0.71 | 0.75 |
| 0.9 | 85.76 | 0.11 | 12.12 | 0.13 | 0.73 | 0.77 |

Table 7: Effect of Temperature t in KD

| t | Aalborg | | | | | | Chengdu | | | | | |
|-----|------------------------|-------------|--------------|--------------|-------------|-------------|------------------------|-------------|--------------|--------------|-------------|-------------|
| | Travel Time Estimation | | | Path Ranking | | | Travel Time Estimation | | | Path Ranking | | |
| | MAE | MARE | MAPE | MAE | τ | ρ | MAE | MARE | MAPE | MAE | τ | ρ |
| 1 | 90.01 | 0.12 | 13.30 | 0.16 | 0.65 | 0.69 | 225.70 | 0.20 | 22.02 | 0.09 | 0.79 | 0.82 |
| 3 | 94.11 | 0.12 | 13.54 | 0.15 | 0.68 | 0.72 | 217.07 | 0.19 | 20.77 | 0.09 | 0.81 | 0.84 |
| 5 | 90.39 | 0.12 | 12.76 | 0.15 | 0.66 | 0.70 | 216.93 | 0.19 | 21.24 | 0.08 | 0.83 | 0.86 |
| 7 | 89.64 | 0.12 | 12.76 | 0.15 | 0.70 | 0.74 | 214.88 | 0.19 | 20.98 | 0.08 | 0.86 | 0.87 |
| 9 | 85.76 | 0.11 | 12.12 | 0.13 | 0.73 | 0.77 | 212.61 | 0.18 | 20.75 | 0.07 | 0.87 | 0.88 |
| 11 | 87.15 | 0.12 | 12.43 | 0.14 | 0.70 | 0.74 | 214.17 | 0.19 | 21.00 | 0.08 | 0.83 | 0.85 |

Model Scalability In the sequel, we explore the model scalability in terms of reduction ratio and path length based on the synthetic dataset. Table 5 depicts the results for both *LightPath* and its teacher model, with varying $\gamma = 0, 0.1, 0.3, 0.5, 0.7, 0.9$. $\gamma = 0$ denotes we do not conduct sparsity operation for the input path, i.e., using a classic Transformer based encoder. We can observe that the GFLOPs and gMem. (GiB) decrease with the increase in the reduction ratio. It is because the higher value of γ is, the more edges we can remove. Second, *LightPath* has significantly reduced model complexity, w.r.t., GFLOPs and gMem.. For example, we can reduce the training GFLOPs by $2.54\times$ for the *LightPath* by increasing the reduction ratio γ from 0 to 0.9 in terms of path length 200. Moreover, *LightPath* also shows better performance (i.e., GFLOPs and gMem.) over teacher model, e.g., $1.79\times$ GFLOPs speedup with reduction ratio $\gamma = 0.9$. Third, the parameters (Para. (Millions)) of teacher model is at least $3.5\times$ of *LightPath*, which implies the effectiveness of our proposed framework. Overall, *LightPath* shows potential of scalability to support path representation learning for long paths.

Effect of Reduction Ratio γ To study the impact of reduction ratio γ in the final performance, we conduct an experiment by varying the γ from 0.1 to 0.9 on Aalborg, which is shown in Table 6. We can observe that the overall performance in both downstream tasks degrades a little when γ increases, which is reasonable as the the model has more input information. However, we can also observe the performance differences are not so significant, which suggests the effectiveness of our proposed framework. Even when a high reduction ratio is applied, the performance does not does not go down too much. Therefore, our proposed method can achieve good scalability while ensuring accuracy.

Effect of Temperature t of Knowledge Distillation To study the effect of the temperature t , we conduct a parameter study on both datasets, which is reported in Table 7. We can observe that the performance of *LightPath* varies with different temperatures. It can be figured out that the best temperature t is 9, which indicates warm temperature can mitigate the peakiness of the teacher model and results in better performance.

Effect of Balancing Factor α To study the effect of the balancing factor of global-local knowledge distillation, we conduct a parameter study on both datasets. Based on the results reported in Table 8, we observe that the performance of our model changes when varying α . We can observe that the optimal α is 0.6, which means that global and local knowledge distillation loss can contribute to the *LightPath*'s performance. When $\alpha = 0$, the global knowledge distillation loss is ignored, which yields poor performance. When $\alpha = 1.0$, the local knowledge distillation loss is ignored, and the performance also performs poorly. This confirms our conjecture that the two proposed global-local

Table 8: Effect of Balancing Factor α

| α | Aalborg | | | | | | Chengdu | | | | | |
|----------|------------------------|-------------|--------------|--------------|-------------|-------------|------------------------|-------------|--------------|--------------|-------------|-------------|
| | Travel Time Estimation | | | Path Ranking | | | Travel Time Estimation | | | Path Ranking | | |
| | MAE | MARE | MAPE | MAE | τ | ρ | MAE | MARE | MAPE | MAE | τ | ρ |
| 0 | 90.24 | 0.12 | 12.78 | 0.16 | 0.69 | 0.73 | 220.32 | 0.19 | 21.52 | 0.09 | 0.79 | 0.81 |
| 0.2 | 89.35 | 0.12 | 12.85 | 0.14 | 0.69 | 0.73 | 217.10 | 0.19 | 21.34 | 0.09 | 0.78 | 0.80 |
| 0.4 | 91.57 | 0.12 | 13.17 | 0.15 | 0.69 | 0.73 | 217.33 | 0.19 | 21.21 | 0.08 | 0.85 | 0.87 |
| 0.6 | 85.76 | 0.11 | 12.12 | 0.13 | 0.73 | 0.77 | 212.61 | 0.18 | 20.75 | 0.07 | 0.87 | 0.88 |
| 0.8 | 87.44 | 0.12 | 12.76 | 0.14 | 0.70 | 0.75 | 214.34 | 0.19 | 20.93 | 0.08 | 0.84 | 0.86 |
| 1 | 89.23 | 0.12 | 12.78 | 0.16 | 0.69 | 0.73 | 215.03 | 0.19 | 21.02 | 0.08 | 0.82 | 0.84 |

Table 9: Comparison with Whole Path Input

| | Aalborg | | |
|--|------------------------|------|-------|
| | Travel Time Estimation | | |
| | MAE | MARE | MAPE |
| <i>LightPath w/o RR</i> ($\gamma = 0.1$) | 91.97 | 0.12 | 13.53 |
| <i>LightPath w/o RR</i> ($\gamma = 0.1$) | 90.85 | 0.12 | 13.39 |
| <i>LightPath</i> ($\gamma = 0.1$) | 82.79 | 0.11 | 11.95 |

knowledge distillation losses can regularize each other and achieve better results than only optimizing one of them (i.e., $\alpha = 0.0$ or $\alpha = 1.0$).

6.3 Comparison with Whole Path Input

Compared with whole path input, we consider a variant “LightPath w/o RR” where only reconstruction loss is used and the relational reasoning (RR) loss is disabled. In this setting, we can see in the table 9 that LightPath w/o RR ($\gamma = 0$), i.e., using whole paths, outperforms LightPath w/o RR ($\gamma = 0.1$), i.e., using partial paths. This means that, when only using the reconstruction loss, using whole paths are indeed better than using partial paths. However, when using both losses, LightPath ($\gamma = 0.1$) outperforms LightPath w/o RR ($\gamma = 0$). This demonstrates that the relational reasoning loss, which employs partial paths to create different views, is indeed effective. *LightPath*,

6.4 Comparison with Fixed Interval Strategy

We then conducted additional experiments where we removed edges at fixed intervals strategy. Specifically, we set the fixed interval to 10 and removed n edges out of every 10 (for instance, we deleted 1 edge out of every 10 edges, i.e., corresponding to a removal ratio of $\gamma = 0.1$). The results on travel time estimation in Aalborg shown in Table 10, which indicate that the fixed interval edge removal strategy (ref. as to first three rows) achieves the worse performance compared with the random edge removal strategy (ref. as to last row).

6.4.1 Model Efficiency

We finally evaluate the model efficiency, including training and inference phases. Figure 8 illustrates the corresponding results. The first observation is that *LightPath* outperforms *PIM* and *Toast* in both training and inference phases. In the training phase, *LightPath* is more than $3\times$ faster than *PIM* and almost $5\times$ faster than *Toast* when path length is 200. In the testing phase, we measure the running time for each path sample. As observed, *LightPath* achieves up to at least 100% and almost 200% performance improvement compared with *PIM* and *Toast* when path length is 200.

7 Related Work

7.1 Path Representation Learning

Path Representation Learning (PRL) aims to learn effective and informative path representations in road network that can be applied to various downstream applications, i.e., travel cost estimation, and path ranking. Existing PRL studies can be categorised as supervised learning (SL) based [27, 18, 19], unsupervised learning (UL) based [4], and weakly

Table 10: Comparison with Fixed Interval Strategy

| Strategy | γ | Aalborg | | |
|-------------------------|----------|------------------------|------|-------|
| | | Travel Time Estimation | | |
| | | MAE | MARE | MAPE |
| <i>LightPath-Fixed</i> | 1/10 | 87.77 | 0.12 | 12.95 |
| <i>LightPath-Fixed</i> | 5/10 | 89.63 | 0.12 | 12.86 |
| <i>LightPath-Fixed</i> | 9/10 | 92.87 | 0.12 | 12.23 |
| <i>LightPath-Random</i> | 0.9 | 85.76 | 0.11 | 12.12 |

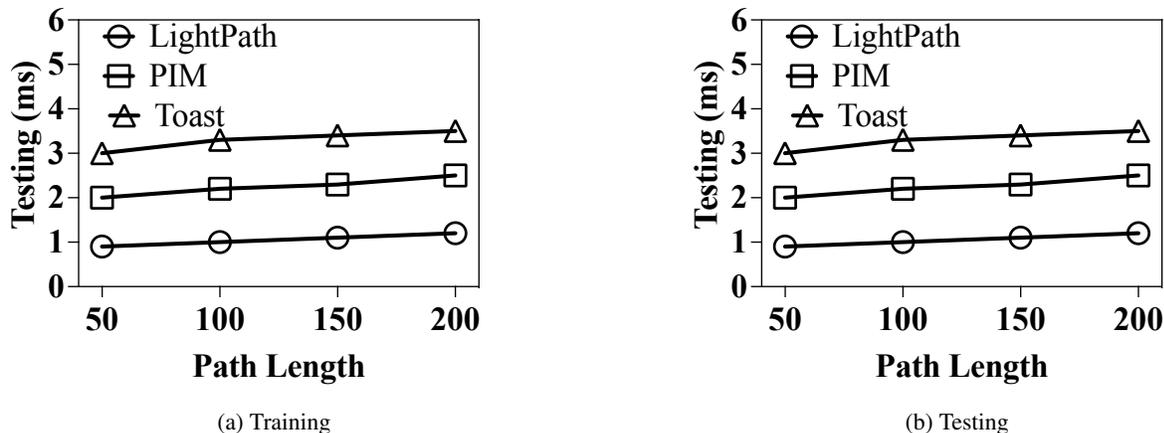


Figure 8: Model Efficiency Evaluation

supervised learning (WSL) based [26] approaches. SL based methods aim at learning a task-specific path representation with the availability of large amounts of labeled training data [19, 27, 18], which has a poor generality for other tasks. UL methods are to learn general path representation learning that does not need labeled training data and generalizes well to multiple downstream tasks [4, 26]. In contrast, WSL methods try to learn a generic temporal path representation by introducing meaningful weak labels, e.g., traffic congestion indices, that are easy and inexpensive to obtain, and are relevant to different tasks [26]. However, we aim to learn generic path representations instead of temporal path representations in this paper. Thus, we do not select WSL method as baseline method. In particular, these methods are computationally expensive and hard to deploy in resource-limited environments.

7.2 Self-supervised Learning

State-of-the-art self-supervised learning can be classified into contrastive learning-based and relation reasoning-based methods. Contrastive learning-based methods [37, 38, 39, 4, 40], especially for InfoNCE loss-based, commonly generate different views of same input data through different augmentation strategies, and then discriminate positive and negative samples. However, these methods suffer from their quadratic complexity, w.r.t. the number of data samples, given that it needs a large number of negative samples to guarantee that the mutual information lower bound is tight enough [38]. In contrast, relation reasoning-based methods [41, 42] aim to learn relation reasoning head that discriminates how entities relate to themselves and other entities, which results in linear complexity. However, existing studies construct relation reasoning between different views from the same encoder, ignoring the effect of different views between different encoders, i.e., main encoder and auxiliary encoder in Siamese encoder architecture.

8 Conclusion

We design a lightweight and scalable framework called *LightPath* for unsupervised path representation learning. In this framework, we first propose sparse auto-encoder that is able to reduce path length N to N' , where N is much larger than N' , which in turn reduces the computation complexity of the model. Then, we use path reconstruction decoder to reconstruct the input path to ensure no edges information missing. Next, we propose a novel self-supervised relational reasoning approach, which contains cross-network relational reasoning and cross-view relational reasoning loss, to

enable efficient unsupervised training. After that, we introduce global-local knowledge distillation to further reduce the size of sparse path encoder and improve the performance. Finally, extensive experiments on two real-world datasets verify the efficiency, scalability, and effectiveness of *LightPath*.

Acknowledgments

This work was supported in part by Independent Research Fund Denmark under agreements 8022-00246B and 8048-00038B, the VILLUM FONDEN under agreements 34328 and 40567.

References

- [1] Chenjuan Guo, Bin Yang, Jilin Hu, and Christian S. Jensen. Learning to route with sparse trajectory sets. In *ICDE*, pages 1073–1084, 2018.
- [2] Ziquan Fang, Lu Pan, Lu Chen, Yuntao Du, and Yunjun Gao. MDTP: A multi-source deep traffic prediction framework over spatio-temporal trajectory data. *Proc. VLDB Endow.*, 14(8):1289–1297, 2021.
- [3] Junbo Zhang, Yu Zheng, Junkai Sun, and Dekang Qi. Flow prediction in spatio-temporal networks based on multitask deep learning. *IEEE Trans. Knowl. Data Eng.*, 32(3):468–478, 2020.
- [4] Sean Bin Yang, Chenjuan Guo, Jilin Hu, Jian Tang, and Bin Yang. Unsupervised path representation learning with curriculum negative sampling. In *IJCAI*, pages 3286–3292, 2021.
- [5] Jilin Hu, Chenjuan Guo, Bin Yang, and Christian S. Jensen. Stochastic weight completion for road networks using graph convolutional networks. In *ICDE*, pages 1274–1285, 2019.
- [6] Haitao Yuan, Guoliang Li, Zhifeng Bao, and Ling Feng. Effective travel time estimation: When historical trajectories over road networks matter. In *SIGMOD*, pages 2135–2149, 2020.
- [7] Xinle Wu, Dalin Zhang, Chenjuan Guo, Chaoyang He, Bin Yang, and Christian S. Jensen. Autocts: Automated correlated time series forecasting. *Proc. VLDB Endow.*, 15(4):971–983, 2021.
- [8] Haitao Yuan, Guoliang Li, Zhifeng Bao, and Ling Feng. An effective joint prediction model for travel demands and traffic flows. In *ICDE*, pages 348–359, 2021.
- [9] Luan V. Tran, Minyoung Mun, Matthew Lim, Jonah Yamato, Nathan Huh, and Cyrus Shahabi. Deeptrans: A deep learning system for public bus travel time estimation using traffic forecasting. *Proc. VLDB Endow.*, 13(12):2957–2960, 2020.
- [10] Zeyuan Shang, Guoliang Li, and Zhifeng Bao. DITA: distributed in-memory trajectory analytics. In *SIGMOD*, pages 725–740, 2018.
- [11] Zeyuan Shang, Guoliang Li, and Zhifeng Bao. DITA: A distributed in-memory trajectory analytics system. In *SIGMOD*, pages 1681–1684, 2018.
- [12] Xiaolin Han, Reynold Cheng, Chenhao Ma, and Tobias Grubenmann. Deeptea: Effective and efficient online time-dependent trajectory outlier detection. *Proc. VLDB Endow.*, 15(7):1493–1505, 2022.
- [13] Zheng Wang, Cheng Long, Gao Cong, and Yiding Liu. Efficient and effective similar subtrajectory search with deep reinforcement learning. *Proc. VLDB Endow.*, 13(11):2312–2325, 2020.
- [14] Sheng Wang, Zhifeng Bao, J. Shane Culpepper, Timos Sellis, and Xiaolin Qin. Fast large-scale trajectory clustering. *Proc. VLDB Endow.*, 13(1):29–42, 2019.
- [15] Simon Aagaard Pedersen, Bin Yang, and Christian S. Jensen. Fast stochastic routing under time-varying uncertainty. *VLDB J.*, 29(4):819–839, 2020.
- [16] Chenjuan Guo, Christian S. Jensen, and Bin Yang. Towards total traffic awareness. *SIGMOD Record*, 43(3):18–23, 2014.
- [17] Simon Aagaard Pedersen, Bin Yang, and Christian S. Jensen. Anytime stochastic routing with hybrid learning. *Proc. VLDB Endow.*, 13(9):1555–1567, 2020.
- [18] Sean Bin Yang, Chenjuan Guo, and Bin Yang. Context-aware path ranking in road networks. *TKDE (Early Access)*, 2020.
- [19] Hao Liu, Jindong Han, Yanjie Fu, Jingbo Zhou, Xinjiang Lu, and Hui Xiong. Multi-modal transportation recommendation with unified route representation learning. *Proc. VLDB Endow.*, 14(3):342–350, 2020.

- [20] Chenjuan Guo, Bin Yang, Jilin Hu, Christian S. Jensen, and Lu Chen. Context-aware, preference-based vehicle routing. *VLDB J.*, 29(5):1149–1170, 2020.
- [21] Huiping Liu, Cheqing Jin, Bin Yang, and Aoying Zhou. Finding top-k optimal sequenced routes. In *ICDE*, pages 569–580, 2018.
- [22] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, pages 5998–6008, 2017.
- [23] Yile Chen, Xiucheng Li, Gao Cong, Zhifeng Bao, Cheng Long, Yiding Liu, Arun Kumar Chandran, and Richard Ellison. Robust road network representation learning: When traffic patterns meet traveling semantics. In *CIKM*, pages 211–220, 2021.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [25] Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *CoRR*, abs/1607.06450, 2016.
- [26] Sean Bin Yang, Chenjuan Guo, Jilin Hu, Bin Yang, Jian Tang, and Christian S. Jensen. Weakly-supervised temporal path representation learning with contrastive curriculum learning - extended version. *CoRR*, abs/2203.16110, 2022.
- [27] Sean Bin Yang and Bin Yang. Learning to rank paths in spatial networks. In *ICDE*, pages 2006–2009, 2020.
- [28] Paul Newson and John Krumm. Hidden markov map matching through noise and sparseness. In *SIGSPATIAL*, pages 336–343, 2009.
- [29] Huiping Liu, Cheqing Jin, Bin Yang, and Aoying Zhou. Finding top-k shortest paths with diversity. *IEEE Trans. Knowl. Data Eng.*, 30(3):488–502, 2018.
- [30] Sven Peter, Ferran Diego, Fred A. Hamprecht, and Boaz Nadler. Cost efficient gradient boosting. In *NIPS*, pages 1551–1561, 2017.
- [31] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *SIGKDD*, pages 855–864, 2016.
- [32] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9726–9735, 2020.
- [33] Xiucheng Li, Kaiqi Zhao, Gao Cong, Christian S. Jensen, and Wei Wei. Deep representation learning for trajectory similarity computation. In *ICDE*, pages 617–628, 2018.
- [34] Di Yao, Gao Cong, Chao Zhang, and Jingping Bi. Computing trajectory similarity in linear time: A generic seed-guided neural metric learning approach. In *ICDE*, pages 1358–1369, 2019.
- [35] Kun Fu, Fanlin Meng, Jieping Ye, and Zheng Wang. Compacteta: A fast inference system for travel time prediction. In *KDD*, pages 3337–3345, 2020.
- [36] Zebin Chen, Xiaolin Xiao, Yue-Jiao Gong, Jun Fang, Nan Ma, Hua Chai, and Zhiguang Cao. Interpreting trajectories from multiple views: A hierarchical self-attention network for estimating the time of arrival. In *KDD*, pages 2771–2779, 2022.
- [37] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018.
- [38] R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *CoRR*, abs/1808.06670, 2018.
- [39] Petar Velickovic, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon Hjelm. Deep graph infomax. In *ICLR*, 2019.
- [40] Dmitrii Babaev, Nikita Ovsov, Ivan Kireev, Mariya Ivanova, Gleb Gusev, Ivan Nazarov, and Alexander Tuzhilin. Coles: Contrastive learning for event sequences with self-supervision. In *SIGMOD*, pages 1190–1199, 2022.
- [41] Massimiliano Patacchiola and Amos J. Storkey. Self-supervised relational reasoning for representation learning. In *NeurIPS*, 2020.
- [42] Haoyi Fan, Fengbin Zhang, and Yue Gao. Self-supervised time series representation learning by inter-intra relational reasoning. *CoRR*, abs/2011.13548, 2020.