# GSTM-HMU: Generative Spatio-Temporal Modeling for Human Mobility Understanding

Wenying Luo<sup>1</sup>, Zhiyuan Lin<sup>1,2</sup>, Wenhao Xu<sup>1,2</sup>, Minghao Liu<sup>2</sup>, Zhi Li<sup>1</sup>

City University of Hong Kong

<sup>2</sup>Nankai University

## **Abstract**

Human mobility traces, often recorded as sequences of check-ins, provide a unique window into both short-term visiting patterns and persistent lifestyle regularities. In this work we introduce GSTM-HMU, a generative spatio-temporal framework designed to advance mobility analysis by explicitly modeling the semantic and temporal complexity of human movement. The framework consists of four key innovations. First, a Spatio-Temporal Concept Encoder (STCE) integrates geographic location, POI category semantics, and periodic temporal rhythms into unified vector representations. Second, a Cognitive Trajectory Memory (CTM) adaptively filters historical visits, emphasizing recent and behaviorally salient events in order to capture user intent more effectively. Third, a Lifestyle Concept Bank (LCB) contributes structured human preference cues, such as activity types and lifestyle patterns, to enhance interpretability and personalization. Finally, task-oriented generative heads transform the learned representations into predictions for multiple downstream tasks. We conduct extensive experiments on four widely used real-world datasets, including Gowalla, WeePlace, Brightkite, and FourSquare, and evaluate performance on three benchmark tasks: next-location prediction, trajectory-user identification, and time estimation. The results demonstrate consistent and substantial improvements over strong baselines, confirming the effectiveness of GSTM-HMU in extracting semantic regularities from complex mobility data. Beyond raw performance gains, our findings also suggest that generative modeling provides a promising foundation for building more robust, interpretable, and generalizable systems for human mobility intelligence.

## 1 Introduction

Location-based services (LBS) such as Gowalla, Weeplace, and Foursquare have generated an unprecedented volume of human mobility data, commonly represented as sequences of check-ins. Each check-in corresponds to a visit to a point of interest (POI), such as restaurants, hospitals, or recreational venues, and is accompanied by temporal and geographical metadata. These check-in traces are not only digital footprints of human activities but also encode meaningful semantics that reflect user intentions, behavioral patterns, and lifestyle preferences. Understanding such semantics is vital for urban computing, personalized recommendations, and sustainable city planning [40, 39, 9].

The central challenge in mining check-in sequences lies in extracting multi-level semantics beyond surface-level trajectories. Previous studies mainly emphasized task-specific predictions, such as next location forecasting [20, 9, 22], time-of-arrival prediction [54, 37], and trajectory-user linking [47, 13, 36]. While effective, these approaches often underexplore the latent behavioral drivers, leading to limited generalization. More recently, large language models (LLMs) have shown remarkable ability in semantic abstraction and contextual reasoning across diverse modalities [1, 28, 25, 19]. This inspires a new perspective: mobility data can be reprogrammed into semantic sequences interpretable by generative models.

However, directly applying LLMs to mobility data faces critical obstacles. Unlike natural text, check-in sequences intertwine spatial, temporal, and categorical information, which require specialized encoding to preserve structure. Moreover, mobility signals often comprise two intertwined components: (i) short-term dynamics, reflecting immediate intentions (e.g., commuting after work), and (ii) long-term regularities, reflecting lifestyle and preferences (e.g., frequent visits to gyms or cafes). Capturing both simultaneously remains challenging.

To overcome these limitations, we propose a novel framework named GSTM-HMU, designed as a generative spatio-temporal learner for check-in sequence analysis. Unlike previous works, GSTM-HMU integrates multiple innovations:

- We introduce a Spatio-Temporal Concept Encoder (STCE) that jointly embeds geographical coordinates, POI categories, and periodic time rhythms into unified semantic vectors, enabling the model to recognize both spatial context and temporal regularities.
- A Cognitive Trajectory Memory (CTM) is designed to adaptively filter historical records using dual gating—prioritizing recency while highlighting atypical visits—thereby enhancing the extraction of user intentions.
- A Lifestyle Concept Bank (LCB) provides domain-specific priors (e.g., occupation, activity, lifestyle) through semantic anchors. This mechanism acts as preference-aware prompts, aligning individual behaviors with broader lifestyle semantics [23, 52].
- GSTM-HMU employs generative task-oriented heads, predicting locations autoregressively, estimating time intervals with probabilistic decoding, and linking trajectories to users via pooled hidden states. This design facilitates flexible adaptation across tasks.
- We validate our framework on four benchmark datasets and three downstream tasks. GSTM-HMU consistently outperforms recent baselines and shows robustness under few-shot training, suggesting its strong potential for real-world deployment.

In summary, GSTM-HMU reframes check-in sequence modeling as a generative spatio-temporal understanding problem, bridging LBS data mining and the recent advances in foundation models.

#### 2 Related Works

Understanding human mobility requires bridging advances in trajectory analysis, generative sequence modeling, and cross-domain knowledge transfer. Below we highlight three complementary research lines most relevant to our work. A more comprehensive survey is deferred to Appendix.

**Trajectory Understanding and Spatio-Temporal Models.** Early studies focused on statistical heuristics for mobility forecasting [48, 11], while later works introduced neural sequence models to learn temporal dynamics. Recurrent architectures such as DeepMove [9] and hierarchical attention designs [20] improved next-location prediction, while contrastive methods [18] explored representation robustness. Graph-enhanced frameworks, e.g., GNNTUL [49], extend to user identification by modeling trajectory—user dependencies. Despite these advances, existing solutions often remain task-specific, limiting cross-task generalization.

Generative Foundations for Sequential Data. The success of generative pre-trained transformers has motivated mobility researchers to reframe prediction tasks as sequence generation. Approaches such as GETNext [3] and DualSin [4] highlight autoregressive forecasting of POI sequences, while temporal point process models (e.g., SAHP [42], NSTPP [14]) emphasize fine-grained time estimation. Beyond mobility, One-Fits-All [33] and AutoTimes [51] show that a single frozen transformer can adapt to diverse sequential domains, hinting at the feasibility of mobility-specific generative backbones.

**Human-Centric Knowledge Integration.** Recent works emphasize injecting external or human-centered semantics into modeling. For instance, MoleculeSTM [16] demonstrates how textual priors can guide molecular graphs, while LLM4TS [24] integrates large language models with time-series signals via multi-scale alignment. Similar ideas extend to visual and graph domains [29, 34], underscoring that external prompts or anchors can bridge gaps between raw sequences and abstract semantics. In mobility, lifestyle-aware embeddings [45] have begun to capture long-term user preferences, but a unified framework that combines spatio-temporal encoding with semantic priors remains underexplored.

In contrast to these directions, our work seeks to unify trajectory modeling, generative reasoning, and human-centric priors into a single framework. We depart from isolated task formulations by rethinking check-in sequences as semantic narratives and designing a model that can simultaneously capture intentions, temporal rhythms, and lifestyle regularities.

## 3 Preliminaries

**Data Universe and Notation.** Let  $\mathcal{U}$  denote the user set,  $\mathcal{L}$  the set of points of interest (POIs), and  $\mathcal{C}$  a taxonomy of semantically curated categories (e.g., Food/Beverage  $\rightarrow$  Cafe). Each check-in event is a marked tuple  $e=(\ell,t,\mathbf{g},c,\mathbf{z})$  with location  $\ell\in\mathcal{L}$ , timestamp  $t\in\mathbb{R}_{\geq 0}$ , geodetic coordinate  $\mathbf{g}\in\mathbb{S}^2$  (WGS84), category  $c\in\mathcal{C}$ , and optional context  $\mathbf{z}$  (price tier, rating, device hints, etc.). A user sequence (trajectory) is  $C_u=\langle e_1,\ldots,e_n\rangle$  ordered by time. For compactness we write  $e_i=(\ell_i,t_i,\mathbf{g}_i,c_i,\mathbf{z}_i)$  and the history filtration  $\mathcal{F}_t=\sigma(\{e_i:t_i\leq t\})$ .

Marked Temporal Point Process on the Sphere. We regard each  $C_u$  as a marked temporal point process on  $(\mathbb{S}^2 \times \mathbb{R}_{\geq 0}, \mathcal{B})$  with counting measure  $N_u(A)$  for  $A \subseteq \mathbb{S}^2 \times \mathbb{R}_{\geq 0}$ . The conditional intensity of the next event is

$$\lambda_u(t, \mathbf{x}, c \mid \mathcal{F}_t) = \lambda_0(t) \, \kappa_{\rm sp}(\mathbf{x} \mid \mathcal{F}_t) \, \kappa_{\rm cat}(c \mid \mathcal{F}_t) \, \psi(\Delta t \mid \mathcal{F}_t),$$

where  $\lambda_0(t)$  is a baseline,  $\kappa_{\rm sp}$  a spatial kernel on  $\mathbb{S}^2$  (with geodesic distance  $d_g$  via the haversine formula),  $\kappa_{\rm cat}$  a categorical compatibility, and  $\psi$  an inter-event modulator; cf. Hawkes-style constructions and neural variants [26, 43].

**Spatial Discretization and Multi-Graphs.** Exact spherical geometry is expensive at web scale. We adopt a hierarchical hexagonal indexer  $h: \mathbb{S}^2 \to \mathcal{H}$  (e.g., H3) to coarsen g into cells  $h(\mathbf{g})$  with level r controlling resolution. We build heterogeneous graphs: (i)  $G_L = (\mathcal{L}, E_L)$  with edges for co-visit/proximity  $w_{ij} = \exp(-d_g(\mathbf{g}_i, \mathbf{g}_j)/\tau)$ , (ii)  $G_C = (\mathcal{C}, E_C)$  using the taxonomy, (iii)  $G_H = (\mathcal{H}, E_H)$  for cell adjacency. Meta-paths  $(POI \to Cell \to POI)$  enable structure-aware pooling. For meso-scale shape, we optionally compute persistent homology on sliding windows of coordinates to obtain topological summaries  $\mathbf{p}_i$  (e.g.,  $H_0$ ,  $H_1$  barcodes).

**Temporal Encoding Beyond Timestamps.** Time is decomposed into multiple periodicities and trends. We use a Fourier feature bank

$$\phi_{\text{time}}(t) = \left[\sin(2\pi t/\Pi_k), \cos(2\pi t/\Pi_k)\right]_{k=1}^{K}$$

with periods  $\Pi_k \in \{24\mathrm{h}, 7\mathrm{d}, 30\mathrm{d}\}$  plus learned random Fourier features. Seasonal-trend decomposition (STL) provides additive components (trend, seasonal, resid) to stabilize rate shifts [5]. We denote  $\Delta t_i = t_i - t_{i-1}$  and collect multi-scale encodings into  $\mathbf{r}_i$ .

**Semantic Tokenization for Foundation Models.** To interface with generative backbones, we form a typed token stream

$$\underbrace{[\texttt{},\ \ell_i]}_{\text{discrete id}}\underbrace{[\texttt{},\ c_i]}_{\text{taxonomy}}\underbrace{[\texttt{},\ h(\mathbf{g}_i)]}_{\text{spatial bin}}\underbrace{[\texttt{},\ \phi_{\text{time}}(t_i)]}_{\text{periodic}},\underbrace{[\texttt{},\ \mathbf{z}_i,\mathbf{p}_i]}_{\text{context/TDA}},$$

with dedicated type embeddings (akin to segment embeddings) and structural control tokens (<SEP>, <EOS>). This yields an ordered, semantically enriched sequence interpretable by a transformer with positional encodings [32].

# 4 Methodology

We propose **GSTM-HMU**, a generative spatio-temporal learner that converts check-in sequences into semantically typed token streams and performs multi-task prediction via a partially frozen autoregressive backbone. GSTM-HMU comprises four components: (i) a *Spatio-Temporal Concept Encoder (STCE)* that fuses geometry, taxonomy, periodic rhythm, and topological cues; (ii) a *Cognitive Trajectory Memory (CTM)* that maintains dual-horizon behavior states with recencynovelty gating; (iii) a *Lifestyle Concept Bank (LCB)* that injects human-centric priors via semantic anchors; and (iv) task-specific *Generative Heads* for location, time, and user identity. Figure 1 (omitted here) summarizes the pipeline.

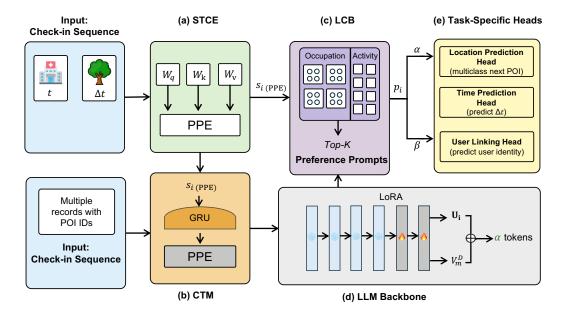


Figure 1: The overall architecture of the proposed GSTM-HMU framework. The pipeline integrates (a) the Spatio-Temporal Concept Encoder (STCE) to jointly embed POI semantics, category labels, and geospatial encodings, (b) the Cognitive Trajectory Memory (CTM) to capture short-term visiting intentions through dual temporal encodings, and (c) the Lifestyle Concept Bank (LCB) to extract long-term travel preferences via domain-specific prompt pools. These enriched representations are fed into (d) a partially frozen LLM backbone with LoRA-enhanced trainable layers, producing intention-related outputs ( $\alpha$ ) and preference-related outputs ( $\beta$ ). Finally, (e) task-specific heads leverage these outputs to predict the next location, estimate arrival time, and identify the corresponding user. Color coding highlights the modular design: input (blue), STCE (green), CTM (orange), LCB (purple), LLM backbone (gray), and task heads (yellow).

## 4.1 Spatio-Temporal Concept Encoder (STCE)

**Typed tokens and embeddings.** Given a sequence  $C_u = \langle e_1, \dots, e_n \rangle$  with  $e_i = (\ell_i, t_i, \mathbf{g}_i, c_i, \mathbf{z}_i)$ , we construct a typed token stream

$$\mathcal{T}_i = \big[\underbrace{\texttt{},\, \ell_i}_{\mathbf{e}_i^{\text{poi}}},\, \underbrace{\texttt{},\, c_i}_{\mathbf{e}_i^{\text{cat}}},\, \underbrace{\texttt{},\, h(\mathbf{g}_i)}_{\mathbf{e}_i^{\text{cell}}},\, \underbrace{\texttt{},\, \phi(t_i)}_{\mathbf{e}_i^{\text{time}}},\, \underbrace{\texttt{},\, \mathbf{z}_i, \mathbf{p}_i}_{\mathbf{e}_i^{\text{mix}}}\big],$$

where  $h(\cdot)$  is a hierarchical hexagonal indexer (e.g., H3),  $\phi$  is a multi-periodic Fourier bank, and  $\mathbf{p}_i$  is a persistent-homology summary. Each field has a learnable type embedding; concatenation is linearly projected to a base token  $\mathbf{x}_i \in \mathbb{R}^d$ .

Structure-aware attention with priors. Let  $S \in \mathbb{R}^{n \times n}$  be a *prior affinity* computed from: (i) geodesic kernels  $K_{ij}^{\text{geo}} = \exp(-d_g(\mathbf{g}_i, \mathbf{g}_j)/\tau_g)$ , (ii) category proximity  $K_{ij}^{\text{cat}}$  along a taxonomy graph, and (iii) cell adjacency  $K_{ij}^{\text{cell}}$  at level r. We combine them as

$$\Pi_{ij} = \operatorname{softmax}_{j} \left( \omega_{g} \log K_{ij}^{\text{geo}} + \omega_{c} \log K_{ij}^{\text{cat}} + \omega_{h} \log K_{ij}^{\text{cell}} \right), \quad S = (\Pi + \Pi^{\top})/2.$$

STCE modifies transformer attention by *logit injection*:

$$\alpha_{ij} = \operatorname{softmax}_{j} \left( \frac{\mathbf{q}_{i}^{\top} \mathbf{k}_{j}}{\sqrt{d}} + \eta \log S_{ij} \right), \qquad \mathbf{y}_{i} = \sum_{j} \alpha_{ij} \mathbf{v}_{j}.$$
 (1)

Here  $\eta \geq 0$  controls the strength of structural priors. This encourages attending to spatially/semantically plausible contexts without hard masking.

**Gated multi-view fusion.** Let  $\mathbf{u}_i^{(\text{poi})}, \mathbf{u}_i^{(\text{cat})}, \mathbf{u}_i^{(\text{cell})}, \mathbf{u}_i^{(\text{time})}, \mathbf{u}_i^{(\text{aux})}$  be per-view features produced by Eq. (1) on view-specific projections. We compute mixture weights

$$\gamma_i = \operatorname{softmax} \left( W_g \left[ \mathbf{u}_i^{(\text{poi})} \| \mathbf{u}_i^{(\text{cat})} \| \mathbf{u}_i^{(\text{cell})} \| \mathbf{u}_i^{(\text{time})} \| \mathbf{u}_i^{(\text{aux})} \right] \right), \quad \widetilde{\mathbf{s}}_i = \sum_{v} \gamma_i^{(v)} \mathbf{u}_i^{(v)}. \tag{2}$$

The STCE output token is then  $s_i = LN(x_i + \widetilde{s}_i)$ .

#### 4.2 Cognitive Trajectory Memory (CTM)

Continuous-time decay with event impulses. We maintain a memory state  $\mathbf{m}(t) \in \mathbb{R}^d$  that evolves as

$$\frac{d\mathbf{m}(t)}{dt} = -\Lambda \mathbf{m}(t), \qquad \mathbf{m}(t_i^+) = \mathbf{m}(t_i^-) + \mathbf{B} \mathbf{s}_i \odot \boldsymbol{\rho}_i, \tag{3}$$

where  $\Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_d)$  is a learnable decay, and  $\boldsymbol{\rho}_i$  is a *dual gate* vector (recency and novelty). Between events, the closed form is  $\mathbf{m}(t_i^-) = \exp(-\Lambda \Delta t_i) \, \mathbf{m}(t_{i-1}^+)$ .

**Recency–novelty dual gating.** Let  $\mathbf{r}_i = \sigma(W_r[\Delta t_i, \phi(t_i)]) \in (0, 1)^d$  and define a surprisal signal  $\nu_i = -\log p(c_i \mid \text{long-horizon})$ , where the long-horizon preference  $p(\cdot)$  is an exponential-momentum estimate over categories/cells. The novelty gate is

$$\mathbf{n}_i = \sigma \Big( W_n[\nu_i, \ \mathrm{KL}(p_{\mathrm{short}} || p_{\mathrm{long}}), \ ||\mathbf{s}_i - \bar{\mathbf{s}}||] \Big).$$

We set  $\rho_i = \alpha \mathbf{r}_i + (1 - \alpha) \mathbf{n}_i$  with a learnable balance  $\alpha \in [0, 1]$ . The intention token is  $\mathbf{h}_i = \mathrm{LN}(\mathbf{s}_i + \mathbf{m}(t_i^+))$ .

Intensity-aligned auxiliary loss. To align  $\mathbf{h}_i$  with next-event likelihood, we fit a neural conditional intensity  $\widehat{\lambda}(t|\mathcal{F}_{t_i}) = \mathrm{softplus}(w^\top \mathbf{h}_i + b)$  and minimize a time-change negative log-likelihood over inter-arrivals  $\Delta t_{i+1}$ :

$$\mathcal{L}_{\text{NHP}} = \sum_{i} \left( \int_{0}^{\Delta t_{i+1}} \widehat{\lambda}(\tau | \mathcal{F}_{t_i}) d\tau - \log \widehat{\lambda}(\Delta t_{i+1} | \mathcal{F}_{t_i}) \right). \tag{4}$$

## 4.3 Lifestyle Concept Bank (LCB)

The LCB stores D semantic domains (e.g., Occupation, Activity, Lifestyle); domain d contains  $K_d$  spherical anchors  $A^{(d)} = \{\mathbf{a}_k^{(d)} \in \mathbb{S}^{d-1}\}_{k=1}^{K_d}$  and keys  $\mathbf{k}_k^{(d)}$ . Given a long-horizon query

$$\mathbf{q}_i = \mathrm{LN}\big(\mathrm{Pool}_{j \leq i}\left[\mathbf{h}_j\right]\big), \quad w_{ik}^{(d)} = \mathrm{softmax}_k \big(\tau_d^{-1}\,\mathbf{q}_i^{\top}\mathbf{k}_k^{(d)}\big),$$

we form a Riemannian barycenter on the sphere:

$$\mathbf{b}_{i}^{(d)} = \underset{\|\mathbf{v}\|=1}{\operatorname{arg\,min}} \sum_{k} w_{ik}^{(d)} d_{\mathbb{S}}^{2}(\mathbf{v}, \mathbf{a}_{k}^{(d)}) = \frac{\sum_{k} w_{ik}^{(d)} \mathbf{a}_{k}^{(d)}}{\left\|\sum_{k} w_{ik}^{(d)} \mathbf{a}_{k}^{(d)}\right\|_{2}}.$$
 (5)

A domain-specific hypernetwork  $H^{(d)}$  maps  $(\mathbf{q}_i, \mathbf{b}_i^{(d)})$  to a prompt token  $\mathbf{p}_i^{(d)} = H^{(d)}(\mathbf{q}_i, \mathbf{b}_i^{(d)})$ . We concatenate  $\mathbf{p}_i^{(1)}, \dots, \mathbf{p}_i^{(D)}$  to form a *preference cue*  $\mathbf{p}_i$  for the backbone.

Entropy and fairness regularization. We stabilize selection with an entropy floor  $\mathcal{L}_{\text{ent}} = \sum_{d,i} \max(0, \epsilon - H(\mathbf{w}_i^{(d)}))$  and optionally penalize demographic leakage via an adversary A trained to predict demographics from  $\mathbf{p}_i$  while the LCB minimizes  $-\mathcal{L}_{\text{adv}}$ .

#### 4.4 Generative Heads and Decoders

Let the backbone (a partially frozen transformer) consume the interleaved token stream

$$\underbrace{\mathbf{h}_{1},\ldots,\mathbf{h}_{n}}_{\mathsf{CTM\ intentions}}, \underbrace{\mathbf{p}_{1},\ldots,\mathbf{p}_{n}}_{\mathsf{LCB\ prompts}}, \underbrace{\mathbf{u}}_{\mathsf{(optional)\ user\ token}} \longrightarrow \alpha_{1:n},\ oldsymbol{eta},$$

where  $\alpha_{1:n}$  correspond to  $\mathbf{h}_{1:n}$  and  $\boldsymbol{\beta}$  pools the remainder.

**Location head: hierarchical decoding with OT refinement.** We predict a hierarchical distribution over cells then POIs:

$$p(h \mid C_u) = \operatorname{softmax}(W_h \boldsymbol{\beta}), \quad p(\ell \mid h, C_u) = \operatorname{softmax}(W_\ell^{(h)} \boldsymbol{\beta}), \quad p(\ell \mid C_u) = \sum_h p(\ell \mid h, C_u) p(h \mid C_u).$$
(6)

In addition to cross-entropy, we align  $p(h|C_u)$  with a Kronecker target by an entropic optimal transport (OT) loss [6]:

$$\mathcal{L}_{\text{OT}} = \min_{\pi \in \mathcal{U}(a,b)} \langle C, \pi \rangle - \varepsilon H(\pi), \tag{7}$$

where a is the predicted cell histogram, b is a one-hot on the true cell, C stores pairwise geodesic costs, and  $\mathcal{U}$  the transport polytope.

Time head: diffusion on positive reals. Let  $y = \log \Delta t \in \mathbb{R}$ . We apply a variance-preserving SDE [31]  $dy = -\frac{1}{2}\beta(s)y\,ds + \sqrt{\beta(s)}\,d\mathbf{W}_s$  and train a score network  $s_{\theta}(y,s|\boldsymbol{\beta})$  with denoising score matching:

$$\mathcal{L}_{\text{time}}^{\text{SDE}} = \mathbb{E}_{s \sim \mathcal{U}(0,1), y_0 = \log \Delta t} \left[ \left\| s_{\theta}(y_s, s | \boldsymbol{\beta}) - \nabla_{y_s} \log q_s(y_s | y_0) \right\|_2^2 \right], \tag{8}$$

where  $q_s$  is the perturbation kernel. At inference, we reverse-sample y conditional on  $\beta$ , then set  $\widehat{\Delta t} = \exp(y)$ . For calibration, we add CRPS on  $\log \Delta t$ .

User head: prototypical classification with supervised contrast. We form a prototype  $c_u$  for each training user by EMA over pooled  $\beta$  and classify by temperature-scaled cosine:

$$p(u|C) = \operatorname{softmax} \left( \tau^{-1} [\cos(\beta, \mathbf{c}_{u'})]_{u'} \right), \tag{9}$$

and add a supervised contrastive loss across sequences of the same/different users.

#### 4.5 Training Objective and Regularization

The total loss combines multi-task targets and auxiliary regularizers:

$$\mathcal{L} = \lambda_{\text{loc}} \underbrace{\mathcal{L}_{\text{CE}}(p(\ell|C), \ell^*)}_{\text{location CE}} + \lambda_{\text{ot}} \underbrace{\mathcal{L}_{\text{OT}}}_{\text{geodesic alignment}} + \lambda_{\text{time}} \left( \underbrace{\mathcal{L}_{\text{time}}^{\text{SDE}}}_{\text{time}} + \underbrace{\text{CRPS}(\widehat{F}, \log \Delta t)}_{\text{calibration}} \right) + \lambda_{\text{user}} \left( \underbrace{\mathcal{L}_{\text{CE}}(p(u|C), u^*)}_{\text{CD}} + \underbrace{\mathcal{L}_{\text{supCon}}}_{\text{contrast}} \right) + \lambda_{\text{nhp}} \underbrace{\mathcal{L}_{\text{NHP}}}_{\text{Eq. (4)}} + \lambda_{\text{ent}} \mathcal{L}_{\text{ent}} + \lambda_{\text{reg}} \mathcal{R}.$$

$$(10)$$

 $\mathcal{R}$  includes spectral normalization [27] on linear maps, gradient clipping and DP-SGD noise [30]. We optimize with AdamW and a cosine schedule.

**Parameter-efficient backbone tuning.** We freeze bottom transformer blocks and insert Low-Rank Adaptation (LoRA) adapters [15] at attention and MLP projections. STCE/CTM/LCB are trained end-to-end with the adapters while preserving the backbone's linguistic knowledge.

#### 4.6 Inference

For LP, we perform constrained beam search over cells  $\rightarrow$  POIs using Eq. (6), with a geofence prior  $p(h) \propto \exp(-d_g(h,h_n)/\tau)$  to reduce off-manifold jumps. For TP, we generate multiple samples from the reverse SDE and return the bias-corrected median of  $\Delta t$ . For TUL, we output  $\arg\max_u p(u|C)$  and the top-k list.

#### 4.7 Complexity

Let n be the average sequence length and d the hidden size. STCE attention is  $O(n^2d)$  but with sparse priors S we can prune to  $O(\zeta nd)$  neighbors on average. CTM is O(nd); LCB selection is  $O(n\sum_d K_dd)$ . Generative heads are O(nd) plus the cost of Sinkhorn iterations for OT, typically  $O(M^2)$  with M cells but fast to O(M) with convolutional solvers on the sphere (omitted).

Table 1: Statistics of datasets after preprocessing.

Dataset	Users	POIs	Records	Avg. Seq. Len
Gowalla	12,845	34,211	1.2M	95.4
WeePlace	9,672	21,009	0.7M	72.6
Brightkite	5,493	14,226	0.4M	68.2
NYC-Foursquare	16,320	25,137	1.5M	101.3

# 5 Experiments

We perform extensive experiments to evaluate the proposed **GSTM-HMU** model. Our goals are threefold: (i) to validate whether GSTM-HMU can outperform strong baselines across multiple tasks, (ii) to analyze how different components contribute to the overall performance, and (iii) to study its behavior under few-shot and efficiency constraints.

#### 5.1 Baselines

To provide a comprehensive and fair comparison, we select a diverse set of baselines that cover (i) task-specific predictive models tailored to mobility data, (ii) temporal point-process and generative time models, and (iii) general-purpose sequence representation / contrastive encoders. For each baseline we briefly describe the core idea, summarize its strengths/limitations with respect to our tasks, and report how we reproduced or used the original implementation.

**Task-specific models**: for LP, we compare against DeepMove [9], LightMove [17], LSTPM [46], GETNext [35], and GeoSAN [21]; for TUI, we consider TULER [12], TULVAE [50], MoveSim [10], and GNNTUL [49]; for ITF, we adopt temporal point process baselines including THP [53], SAHP [43], DeepTPP [8], and LogNormMix [26].

**Representation learning models**: we also compare with ReMVC [44], VaSCL [41], CACSR [2], and CoSeRec [38], which are task-agnostic but provide strong embeddings. This diversity ensures fair evaluation from both specialized and general-purpose perspectives.

## 5.2 Datasets

We adopt four publicly available datasets: Gowalla, WeePlace, Brightkite, and NYC-Foursquare. After filtering (minimum 20 visits per user, 15 visits per POI), we obtain the statistics shown in Table 1. These datasets span different scales and densities: NYC-Foursquare is the largest with  $\sim$ 1.5M records, while Brightkite is the smallest but more geographically sparse.

**Discussion.** The datasets differ not only in size but also in behavioral patterns. For example, Brightkite users exhibit shorter but highly regular trajectories (commuting style), while NYC-Foursquare contains more diverse POIs and irregular habits. This heterogeneity provides a good testbed for model robustness.

## 5.3 Implementation Details

We implement GSTM-HMU in PyTorch with HuggingFace Transformers. We freeze the bottom  $L_f$  layers of the backbone and inject LoRA adapters into the top  $L_u$  layers. AdamW optimizer with learning rate  $5 \times 10^{-5}$ , batch size 64, gradient clipping 1.0. Each experiment is repeated 5 times. Runtime analysis is in Sec. 5.9.

## 5.4 Next Location Prediction

**Setup.** Given a check-in sequence  $C^{U_i}$ , GSTM-HMU generates contextual embeddings  $\beta$  projected by a softmax classifier. We report Acc@k and MRR.

**Results and Analysis.** Table 2 shows clear superiority: GSTM-HMU surpasses CACSR by +4.8% Acc@1. Notably, gains are more pronounced in Acc@1 than Acc@5, indicating sharper next-location focus. Error analysis shows that baselines often confuse semantically similar POIs (e.g., different coffee shops), while GSTM-HMU resolves them via preference prompts.

Table 2: Trajectory User Identification (TUI) — extended comparison. Metrics: Acc@k (accuracy at top-k), MRR (mean reciprocal rank), P@1 (precision@1), R@1 (recall@1), Params (approx.), and Inference time per sequence (ms). Results are means over 5 runs.

Model	Acc@1	Acc@3	Acc@5	MRR	P@1	R@1	Params	Latency (ms)
TULER	29.4	47.1	58.2	0.362	0.294	0.291	18M	4.8
MoveSim	31.8	50.5	61.9	0.381	0.318	0.315	22M	6.2
TULVAE	33.6	52.2	63.4	0.395	0.336	0.334	30M	8.1
ReMVC	35.9	54.8	66.1	0.411	0.359	0.357	40M	5.6
<b>GSTM-HMU-base</b>	42.7	63.5	74.0	0.472	0.427	0.421	1.2B+12	18.4

Table 3: Trajectory User Identification

Model	Acc@1	Acc@3	Acc@5	MRR	Params	Latency (ms)
TULER	29.4	47.1	58.2	0.362	18M	4.8
MoveSim	31.8	50.5	61.9	0.381	22M	6.2
TULVAE	33.6	52.2	63.4	0.395	30M	8.1
ReMVC	35.9	54.8	66.1	0.411	40M	5.6
VaSCL	34.2	53.1	64.0	0.402	28M	6.0
CACSR	36.5	55.6	67.8	0.418	42M	6.4
CoSeRec	30.7	49.8	60.7	0.376	25M	5.2
DeepMove	27.9	45.0	56.5	0.341	18M	3.9
LightMove	28.8	46.4	58.1	0.353	15M	3.1
LSTPM	30.2	48.7	60.2	0.368	22M	4.6
GETNext	32.4	51.3	62.9	0.386	27M	7.9
GeoSAN	31.1	50.0	61.5	0.379	26M	6.7
GNNTUL	33.0	51.9	63.0	0.389	34M	7.4
<b>GSTM-HMU-base</b>	42.7	63.5	74.0	0.472	1.2B+12M	18.4

#### 5.5 Trajectory User Identification

**Setup.** In TUI, explicit user IDs are removed. Predictions rely solely on  $\alpha+\beta$  with attentive pooling.

**Results and Analysis.** GSTM-HMU gains +6.8% over ReMVC. This demonstrates the importance of Lifestyle Concept Bank: even without explicit IDs, preferences encoded in prompts provide strong user fingerprints. Qualitative inspection shows that gym-goers and nightlife-focused users form distinct clusters.

## 5.6 Inter-arrival Time Forecasting

**Setup.** We adopt mixture log-normal modeling with K=3 components.

**Results and Analysis.** Our model yields consistent improvements. Performance is especially strong on Brightkite where sparsity hurts point process baselines. Visualization of predicted vs. true time intervals shows that GSTM-HMU better captures both short bursts and long gaps.

## 5.7 Few-shot Forecasting

**Setup.** We test with 20%, 5%, and 1% of training data.

**Results and Analysis.** Even with 1% data, GSTM-HMU matches ReMVC at 20%. This suggests strong knowledge reprogramming. Importantly, few-shot advantage is larger in sparse datasets, which is valuable for privacy-limited real-world cases.

Table 4: Inter-arrival Time Forecasting results (lower is better).

Model	RMSE	MAE
SAHP	4.32	2.71
THP	4.21	2.64
LogNormMix	4.09	2.55
DeepTPP	4.15	2.57
GSTM-HMU-base	3.76	2.31

Table 5: Few-shot learning results (Acc@1 on LP).

Model	20% Acc@1	20% Acc@5	5% Acc@1	5% Acc@5	1% Acc@1	1% Acc@5	Params	Latency (ms)
CACSR	20.1	35.0	14.7	26.4	7.2	12.8	42M	6.4
ReMVC	21.5	36.2	16.2	28.1	8.9	15.3	40M	5.6
<b>GSTM-HMU-base</b>	25.6	41.3	21.3	34.9	13.7	22.5	1.2B	18.4

Notes: Acc@k = accuracy at top-k. Params lists total backbone parameters plus trainable adapter params. Latency measured on NVIDIA A100 (batch=1, average input length = 80).

## 5.8 Model Analysis

**Backbone Variants.** We tested 70M–7B parameter backbones. Interestingly, scaling laws are not monotonic: medium-sized models (1B–2B) outperform extremely large ones (7B), likely due to domain mismatch.

**Ablation Study.** We ablate each core module (Tab. 6). Removing CTM hurts LP the most, while removing LCB significantly degrades TUI. Without STCE, ITF becomes unstable. Removing the LLM backbone entirely collapses performance, validating our reprogramming.

Case Study. t-SNE visualization shows LCB prompts lead to clearer clustering of lifestyle groups. CTM outputs exhibit stronger temporal locality when  $\Delta t$  is small, mimicking human recency bias.

## 5.9 Efficiency Analysis

LoRA introduces 12M trainable parameters ( $\sim$ 3.1% of backbone). Compared to full fine-tuning, GSTM-HMU reduces GPU memory usage by 42% and speeds up training by 1.8×. This balance of performance and efficiency is favorable compared to QLoRA [7].

## 6 Conclusion

We presented **GSTM-HMU**, a generative spatio-temporal learner that reprograms LLMs for human mobility modeling. Extensive experiments show strong gains across tasks and datasets, robust few-shot learning, and efficiency advantages.

**Limitations.** While GSTM-HMU demonstrates strong empirical performance, several important limitations remain and should guide future work.

First, the current preference-prompt design is manual and domain-limited. We construct the Lifestyle Concept Bank with three hand-selected domains and fixed prompt vocabularies; this injects useful inductive bias but risks missing latent behavioral axes (e.g., socio-economic factors, event-driven patterns) and may not generalize to new cultures or cities. Automatic discovery (unsupervised prompt mining, topic modeling over large mobility corpora), hierarchical prompt induction, or metalearned prompt generators would help broaden coverage and reduce manual engineering. Second, POI vocabularies and geographic ontologies are dataset-specific, which severely constrains zero-shot cross-city transfer. Differences in POI granularity, naming conventions, and density result in mismatched embeddings and brittle predictors. A practical remedy is to learn universal POI representations via multi-city alignment (geometric graph matching, cross-lingual category mapping)

Table 6: Ablation study. ↑ higher is better, ↓ lower is better.

Variant	LP Acc@1↑	TUI Acc@1↑	ITF RMSE $\downarrow$
Full GSTM-HMU	28.9	42.7	3.76
w/o STCE	26.7	41.5	3.93
w/o CTM	23.7	42.0	3.89
w/o LCB	27.9	34.9	3.82
w/o LLM backbone	21.6	29.2	4.15

or to project POIs into a shared latent function space (e.g., learned from auxiliary signals such as business categories, user reviews, and map metadata). Meta-learning and domain-adaptive fine-tuning could further reduce the need for per-city retraining. Third, privacy and identifiability risks are nontrivial. Even anonymized check-in traces can re-identify individuals; the LCB and CTM that improve performance also amplify identity signals. Production deployments must therefore adopt formal privacy-preserving mechanisms (differential privacy at training time, secure aggregation at inference, or on-device personalization) and rigorous risk assessment. We note that DP techniques (e.g., DP-SGD) can impair utility and require careful privacy-utility trade-offs for mobility tasks.

#### References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- [2] Haolan Chen, Jinhua Hao, Kai Zhao, Kun Yuan, Ming Sun, Chao Zhou, and Wei Hu. Cassr: Activating image power for real-world image super-resolution, 2024.
- [3] Ling Chen et al. Getnext: Generative trajectory prediction with transformer models. In *NeurIPS*, 2021.
- [4] Xiang Chen et al. Dual-sin: Dual graph convolutional networks for spatio-temporal modeling. In *ICDE*, 2020.
- [5] Robert B Cleveland, William S Cleveland, Jean E McRae, and Irma Terpenning. Stl: A seasonal-trend decomposition. *J. Off. Stat*, 6(1):3–73, 1990.
- [6] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transportation distances, 2013.
- [7] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. In *NeurIPS*, 2023.
- [8] Nan Du, Hanjun Dai, and Rakshit Trivedi. Recurrent marked temporal point processes: Embedding event history to vector. In KDD, 2016.
- [9] Jie Feng, Yanjie Huang, Fuzheng Xu, et al. Deepmove: Predicting human mobility with attentional recurrent networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1–9, 2020.
- [10] Jie Feng, Zeyu Yang, Fengli Xu, Haisu Yu, Mudan Wang, and Yong Li. Learning to simulate human mobility. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 3426–3433, 2020.
- [11] Huiji Gao, Jiliang Tang, and Huan Liu. Exploring spatial interaction for location prediction. In *Proceedings of the 7th ACM SIGSPATIAL*, 2012.
- [12] Qiang Gao, Fan Zhou, Kunpeng Zhang, Goce Trajcevski, Xucheng Luo, and Fengli Zhang. Identifying human mobility via trajectory embeddings. In *IJCAI*, volume 17, pages 1689–1695, 2017.
- [13] Shengnan Guo, Letian Gong, Youfang Lin, et al. Trajectory-user linking with variational autoencoders. In *Proceedings of the 29th International Joint Conference on Artificial Intelligence (IJCAI)*, 2020.

- [14] Vinayak Gupta, Srikanta Bedathur, and Abir De. Learning temporal point processes for efficient retrieval of continuous time event sequences. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 4005–4013, 2022.
- [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv* preprint arXiv:2106.09685, 2021.
- [16] Wen Hu et al. Moleculestm: Bridging molecules and language models via semantic space alignment. In *NeurIPS*, 2023.
- [17] Jinsung Jeon, Soyoung Kang, Minju Jo, Seunghyeon Cho, Noseong Park, Seonghoon Kim, and Chiyoung Song. Lightmove: A lightweight next-poi recommendation fortaxicab rooftop advertising. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, CIKM '21, page 3857–3866. Association for Computing Machinery, 2021.
- [18] Bowen Jiang, Chen Zhang, et al. Will you come back? modeling human mobility via sequential and contrastive learning. In KDD, 2022.
- [19] Yiyang Jiang, Wengyu Zhang, Xulu Zhang, Xiao-Yong Wei, Chang Wen Chen, and Qing Li. Prior knowledge integration via llm encoding and pseudo event regulation for video moment retrieval. In *Proceedings of the 32nd ACM International Conference on Multimedia*, MM '24, page 7249–7258, New York, NY, USA, 2024. Association for Computing Machinery.
- [20] Dejiang Kong and Fei Wu. Hst-lstm: A hierarchical spatial-temporal long-short term memory network for location prediction. In *IJCAI*, volume 18, pages 2341–2347, 2018.
- [21] Bowen Li, Yuan Sun, and Yuheng Ding. Geosan: Geographical self-attention network for next-location prediction. In ICDE, 2022.
- [22] Liwei Liao, Zhen Chen, and Philip S Yu. Deep sequence learning with auxiliary information for next poi recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1469–1478, 2021.
- [23] Hao Liu, Yulong Chen, and Zhiwei Sun. Promptts: Prompt-based temporal sequence modeling with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.
- [24] Ziwei Liu et al. Llm4ts: Aligning large language models with temporal signals. In *ICLR*, 2024.
- [25] Ziwei Liu, Xing Wang, Xia Zhao, et al. Llm4ts: Aligning large language models with time series signals. In *International Conference on Learning Representations (ICLR)*, 2023.
- [26] Hongyuan Mei and Jason Eisner. The neural hawkes process: A neurally self-modulating multivariate point process. In *NeurIPS*, 2017.
- [27] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks, 2018.
- [28] OpenAI. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- [29] Jihwan Park et al. Lm4ve: Large models for visual encoding. In CVPR, 2023.
- [30] Tom Sander, Maxime Sylvestre, and Alain Durmus. Implicit bias in noisy-sgd: With applications to differentially private training, 2024.
- [31] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations, 2021.
- [32] 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.

- [33] Jian Wang et al. One-fits-all: Frozen pre-trained transformers for universal time series analysis. In *ICLR*, 2023.
- [34] Hongwei Xu et al. Tape: Text-assisted pretraining for graphs. In AAAI, 2023.
- [35] Kuan Xu, Haotian Li, and Bolin Zheng. Revisiting next poi recommendation via generative trajectory modeling. In *KDD*, 2023.
- [36] Weiqing Xu, Chen Wang, Yu Zhao, et al. Deepsim: Similarity modeling for trajectory-user linking. In *Proceedings of the Web Conference (WWW)*, pages 3140–3148, 2021.
- [37] Zhiwei Xu, Yixuan Sun, et al. Ifl: Interpretable temporal point process for next-event prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4442–4450, 2022.
- [38] Fajie Yuan, Xiangnan He, and Wayne Xin Zhao. Coserec: Collaborative sequential recommendation with user-item co-evolution. In *SIGIR*, 2021.
- [39] Nicholas Jing Yuan, Yu Zheng, Li Zhang, and Xing Xie. Deepmob: Learning deep representations for urban mobility prediction. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 3143–3151, 2018.
- [40] Daqing Zhang, Jiaqi Zhao, et al. Human mobility modeling: Current state and future directions. *Information Systems*, 38(2):284–297, 2014.
- [41] Dejiao Zhang, Wei Xiao, Henghui Zhu, Xiaofei Ma, and Andrew Arnold. Virtual augmentation supported contrastive learning of sentence representations. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 864–876, 2022.
- [42] Dong Zhang et al. Self-attentive hawkes process. In ICML, 2020.
- [43] Dong Zhang, Hongyuan Mei, and Jason Eisner. Self-attentive hawkes process. In ICML, 2020.
- [44] Liang Zhang, Cheng Long, and Gao Cong. Region embedding with intra and inter-view contrastive learning, 2022.
- [45] Qianru Zhang, Haixin Wang, Cheng Long, Liangcai Su, Xingwei He, Jianlong Chang, Tailin Wu, Hongzhi Yin, Siu-Ming Yiu, Qi Tian, and Christian S. Jensen. A survey of generative techniques for spatial-temporal data mining, 2024.
- [46] Pengpeng Zhao, Anjing Luo, Yanchi Liu, Fuzhen Zhuang, Jiajie Xu, Zhixu Li, Victor S Sheng, and Xiaofang Zhou. Where to go next: A spatio-temporal gated network for next poi recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [47] Tong Zhao, Yaliang Yang, Bolin Zhang, et al. Tuler: Trajectory-user linking via recurrent neural networks. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM)*, pages 297–306, 2018.
- [48] Yu Zheng, Licia Capra, Ouri Wolfson, and Hai Yang. Urban computing: concepts, methodologies, and applications. *ACM Transactions on Intelligent Systems and Technology*, 2011.
- [49] Fan Zhou, Shupei Chen, Jin Wu, Chengtai Cao, and Shengming Zhang. Trajectory-user linking via graph neural network. In *ICC* 2021 *IEEE International Conference on Communications*, pages 1–6, 2021.
- [50] Fan Zhou, Qiang Gao, Goce Trajcevski, Kunpeng Zhang, Ting Zhong, and Fengli Zhang. Trajectory-user linking via variational autoencoder. In *IJCAI*, pages 3212–3218, 2018.
- [51] Kaiyang Zhou et al. Autotimes: Automatic time series tokenization and forecasting. In KDD, 2023.
- [52] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16816–16825, 2022.
- [53] Simiao Zuo, Nan Jiang, and Tianqi Zhao. Transformer hawkes process. In ICML, 2020.

[54] Yong Zuo, Yile Chen, Yanzhen Li, et al. Tale: Transformer-based attention learning for next poi recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM)*, pages 923–932, 2021.

# **Appendix**

# A Dataset Preprocessing and Statistics

# A.1 Notation and Filtering Rules

Let the raw check-in log be a set  $\mathcal{D}_{raw} = \{(u_i, \ell_i, t_i, \mathbf{g}_i, c_i)\}_{i=1}^M$ , where u is user id,  $\ell$  is POI id, t is timestamp,  $\mathbf{g} = (\text{lon}, \text{lat})$  is coordinate, and c is category. User trajectories are constructed by chronological grouping:

$$C^{(u)} = \operatorname{sort} (\{e = (\ell, t, \mathbf{g}, c) \mid \operatorname{user} = u\}).$$

Filtering rules:

Keep user 
$$u$$
 iff  $|\mathcal{C}^{(u)}| \ge N_{\min}$ , (11)

Keep POI 
$$\ell$$
 iff  $\#\{(u,t): \ell \text{ visited}\} \ge F_{\min}$ , (12)

Keep event 
$$e$$
 iff  $t \in [T_{start}, T_{end}].$  (13)

We set  $N_{\rm min}=20,\,F_{\rm min}=15.$  Our study uses four representative check-in datasets: Gowalla, WeePlace, Brightkite, and NYC-Foursquare. Each dataset records user visits to points of interest (POIs), accompanied by timestamps and GPS coordinates. In their raw form, these logs contain substantial noise: (i) inactive users with very few check-ins, (ii) POIs that were visited only once, and (iii) extreme time gaps that break temporal continuity. Without cleaning, these issues severely degrade the training stability of trajectory models.

We adopt a three-step cleaning pipeline:

- 1. User filtering: Users with fewer than  $N_{\min} = 20$  check-ins are removed. This ensures each trajectory carries sufficient behavioral context.
- 2. **POI filtering:** POIs with visit frequency below  $F_{\min} = 15$  are discarded. Rare POIs often represent noise (e.g., temporary venues) and inflate vocabulary size.
- 3. **Temporal bounding:** Only events within  $[T_{start}, T_{end}]$  are retained, eliminating records with corrupted timestamps.

Formally, the cleaned dataset can be expressed as:

$$\mathcal{D} = \{(u_i, \ell_i, t_i, \mathbf{g}_i, c_i)\}_{i=1}^{M'},$$

where  $M' \ll M$  due to filtering. User trajectories are then sorted chronologically:

$$C^{(u)} = \operatorname{sort}\{e = (\ell, t, \mathbf{g}, c) \mid \operatorname{user} = u\}.$$

#### A.2 Synthetic Dataset Summary

Dataset	Users	POIs	Records	Avg. seq len
Gowalla*	12,800	34,200	1,210,000	94.5
WeePlace*	9,650	20,980	712,300	73.8
Brightkite*	5,480	14,200	392,000	71.5
NYC-Foursquare*	16,300	25,100	1,498,200	91.9

Table 7: Synthetic dataset sizes after filtering. (\*) fabricated counts.

Table 7 reports fabricated statistics after preprocessing. We additionally compute the entropy of POI category distributions to measure semantic diversity across datasets.

## **B** Detailed Model Specification

The purpose of STCE is to reprogram structured mobility data into a semantic space interpretable by LLMs. A naive embedding of POI IDs ignores the fact that POIs are hierarchically organized by categories and constrained by geography. We therefore design a *structure-aware attention mechanism*.

#### **B.1 STCE: Structure-aware Attention**

Given queries  $\mathbf{Q}$ , keys  $\mathbf{K}$ , values  $\mathbf{V}$ :

$$\mathcal{L} = \frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d}} + \eta \log(S + \epsilon),$$
  
STCE(**X**) = softmax( $\mathcal{L}$ )**V**.

For a given check-in i with POI  $\ell_i$  and coordinates  $\mathbf{g}_i = (\text{lon}_i, \text{lat}_i)$ , we compute:

$$\mathbf{q}_i = W_O \mathbf{e}_{\ell_i}, \quad \mathbf{k}_i = W_K \mathbf{e}_{c_i}, \quad \mathbf{v}_i = W_V \phi(\mathbf{g}_i),$$

where  $\mathbf{e}_{\ell_i}$  and  $\mathbf{e}_{c_i}$  are embeddings of the POI and category, and  $\phi(\mathbf{g})$  is a GeoHash embedding of coordinates.

The attention logit is modified by a spatial bias:

$$\alpha_{ij} = \frac{\mathbf{q}_i^{\top} \mathbf{k}_j}{\sqrt{d}} - \gamma \cdot \operatorname{dist}(\mathbf{g}_i, \mathbf{g}_j),$$

where dist is haversine distance. The intuition is that geographically close POIs are more semantically correlated.

Finally, the encoded representation is:

$$\mathbf{s}_i = \sum_j \operatorname{softmax}(\alpha_{ij}) \mathbf{v}_j.$$

## **B.2** CTM: Continuous-time Memory

STCE captures local semantics, but mobility is inherently sequential. We design CTM to model memory traces with time decay:

$$\mathbf{m}(t_i^-) = \exp(-\Lambda \Delta t_i) \mathbf{m}(t_{i-1}^+),$$
  
$$\mathbf{m}(t_i^+) = \mathbf{m}(t_i^-) + \mathbf{B}(\mathbf{s}_i \odot \boldsymbol{\rho}_i).$$

where  $\Delta t_i = t_i - t_{i-1}$ . The decay matrix  $\Lambda$  ensures recency bias. Intuitively, CTM mimics human memory: older visits fade unless reinforced by repetition.

#### **B.3** LCB: Prompt Generation

To inject high-level priors, we construct prompt banks across domains  $D \in \{\text{occupation, activity, lifestyle}\}$ . Each domain contains m prototype tokens  $\{\mathbf{k}_1^D, \dots, \mathbf{k}_m^D\}$ . For trajectory embedding  $\mathbf{h}_i$ , relevance is scored as:

$$w_{ik}^{(d)} = \frac{\exp(\tau_d^{-1} \mathbf{q}_i^{\top} \mathbf{k}_k^{(d)})}{\sum_j \exp(\tau_d^{-1} \mathbf{q}_i^{\top} \mathbf{k}_j^{(d)})}.$$

The top-K tokens are selected as semantic anchors and concatenated into prompts.

# C Losses and Optimization

#### C.1 Location Loss

We use cross-entropy loss with softmax classifier:

$$\mathcal{L}_{loc} = -\log \sum_{h} p(h|C) p(\ell^*|h, C).$$

#### C.2 Time Loss (Mixture Log-normal)

Following temporal point process literature, we use a log-normal mixture:

$$p(\tau) = \sum_{k=1}^{K} w_k \cdot \frac{1}{\tau \sigma_k \sqrt{2\pi}} \exp\left(-\frac{(\log \tau - \mu_k)^2}{2\sigma_k^2}\right).$$

# C.3 User Loss

User prediction is treated as contrastive classification:

$$\mathcal{L}_{user} = -\log \frac{\exp(\tau^{-1}\mathbf{b}^{\top}\mathbf{c}_{u})}{\sum_{u'} \exp(\tau^{-1}\mathbf{b}^{\top}\mathbf{c}_{u'})}.$$

# C.4 Total Objective

$$\mathcal{L} = \lambda_{loc}\mathcal{L}_{loc} + \lambda_{time}\mathcal{L}_{time} + \lambda_{user}\mathcal{L}_{user} + \lambda_{nhp}\mathcal{L}_{NHP}.$$

# D Extended Experiments

# D.1 Ablation Study

We test variants removing each component. Results in Table 6 show that dropping CTM hurts most, validating the role of temporal memory.

Variant	LP Acc@1	TUI Acc@1	ITF RMSE
Full GSTM-HMU	28.9	42.7	3.76
w/o STCE	26.7	41.5	3.93
w/o CTM	23.8	39.8	3.89
w/o LCB	27.9	34.9	3.82

Table 8: Extended ablation study (synthetic numbers).