# Massive-STEPS: Massive Semantic Trajectories for Understanding POI Check-ins – Dataset and Benchmarks

**Wilson Wongso, Hao Xue, Flora Salim**
School of Computer Science and Engineering
University of New South Wales
{w.wongso, hao.xue1, flora.salim}@unsw.edu.au

## Abstract

Understanding human mobility through Point-of-Interest (POI) recommendation is increasingly important for applications such as urban planning, personalized services, and generative agent simulation. However, progress in this field is hindered by two key challenges: the over-reliance on older datasets from 2012-2013 and the lack of reproducible, city-level check-in datasets that reflect diverse global regions. To address these gaps, we present Massive-STEPS (Massive Semantic Trajectories for Understanding POI Check-ins), a large-scale, publicly available benchmark dataset built upon the Semantic Trails dataset and enriched with semantic POI metadata. Massive-STEPS spans 12 geographically and culturally diverse cities and features more recent (2017-2018) and longer-duration (24 months) check-in data than prior datasets. We benchmarked a wide range of POI recommendation models on Massive-STEPS using both supervised and zero-shot approaches, and evaluated their performance across multiple urban contexts. By releasing Massive-STEPS, we aim to facilitate reproducible and equitable research in human mobility and POI recommendation. The dataset and benchmarking code are available at: https://github.com/cruiseresearchgroup/Massive-STEPS.

## 1 Introduction

**Importance of Human Mobility Data and Modeling** Human mobility data and modeling are essential for understanding how individuals interact with and move through physical spaces. This understanding enables a wide range of applications, including urban planning [63], travel service recommendations [7], improved commercial advertising strategies [59], and Point-of-Interest (POI) recommendation [6, 22, 68]. Recently, human mobility data has become even more crucial with the increasing use of large language models (LLM) agents to simulate human-like behavior and routines [71, 20]. However, while simulated human mobility data are starting to gain popularity [31, 9], they may not accurately reflect real-world human behavior [33], highlighting the value of evaluating on real-world data. These advancements are enabled by and large with the rise of Location-based Social Networks (LBSNs), which generate vast amounts of spatio-temporal data through user check-ins [68, 22]. This rich data source has allowed the development of POI recommendation systems that leverage users' historical visiting behaviors to suggest relevant locations. Such systems enhance user engagement through personalization and provide commercial value to both users and businesses by aligning recommendations with individual preferences and available services [68, 6].

**Literature Gaps** Our paper addresses three critical gaps in POI recommendation research and datasets, as detailed in Section 2.1. First, as shown in Fig. 1, the field remains dominated by studies focused on just two cities, New York and Tokyo, based on the Foursquare dataset curated in [57].

Table 1: **Comparison of check-in datasets commonly used for POI recommendation tasks**. GSCD [57, 56] and Semantic Trails [29] are global datasets not grouped into individual cities, whereas others perform city-level grouping. [†]Replicable indicates whether city boundaries are clearly defined or can be reliably reconstructed.

| Dataset | Scale | | | Completeness | Usability | |
|---|---|---|---|---|---|---|
| | #cities | Years | #months | POI Attributes | Replicable[†] | Open-source |
| GSCD [57, 56] | Varies | 2012-2013 | 17 | Coordinates, Category | N/A | ✓ |
| Semantic Trails [29] | Varies | **2012-2013, 2017-2018** | **24** | Category | N/A | ✓ |
| NYC and Tokyo [57] | 2 | 2012-2013 | 11 | Coordinates, Category | ✓ | ✓ |
| Gowalla-CA [4, 62] | 1 | 2009-2010 | 21 | Coordinates, Category | ✓ | ✓ |
| AgentMove [7] | 12 | 2012-2013 | 17 | Coordinates, Category | ✗ | ✗ |
| **Massive-STEPS** | 12 | **2012-2013, 2017-2018** | **24** | **Coordinates, Category, Name, Address** | ✓ | ✓ |

This dataset, collected in 2012-2013, raises concerns about its temporal quality, as many POIs may no longer exist and user behavior may have changed [61]. While some recent studies have expanded to other cities [65, 28, 7], they often rely on the Global-scale Check-in Dataset (GSCD) [55, 56], which, despite its large coverage, is also from 2012-2013 and contains nearly 50% erroneous entries [29].

Second, most existing studies are difficult to reproduce, either due to the lack of clearly defined geographic boundaries or the unavailability of the datasets themselves, hindering fair comparison and replication. Finally, we join recent efforts [63] in advocating for the inclusion of low-resource and underrepresented cities in evaluation. Expanding beyond well-studied urban centers is essential for building more generalizable and universally applicable POI recommendation models. Table 1 summarizes these limitations in terms of geographic coverage, temporal span, and reproducibility.

**Massive-STEPS Dataset** In this paper, we introduced the Massive Semantic Trajectories for Understanding POI Check-ins (Massive-STEPS) Dataset, derived from the Semantic Trails dataset (STD) [29]. Massive-STEPS includes high-quality check-ins from 2012-2013 and 2017-2018, providing more modern and updated POI check-in data. This supports longitudinal POI recommendation studies and addresses the limitations of older datasets commonly used in prior studies. The dataset covers 12 diverse cities across multiple regions, including East, West, and Southeast Asia, North and South America, Australia, the Middle East, and Europe. Notably, we placed a deliberate emphasis on under-explored regions by including cities such as Jakarta, Kuwait City, and Petaling Jaya, filling a key gap in POI recommendation research that has largely focused on major urban centers. We further enriched STD by aligning it with Foursquare's Open Source Places dataset, incorporating metadata such as POI coordinates, POI names, and addresses—details unavailable in the original STD.

**Benchmark Tasks** To demonstrate the utility of this dataset, we conducted an extensive benchmark on two tasks: (1) supervised POI recommendation and (2) zero-shot POI recommendation. Our benchmark covers a wide range of models, including traditional approaches, deep learning-based models, and more recent LLM-based methods. The goal of POI recommendation task is to predict a set of POIs that a user is likely to visit based on their current check-in trajectory and historical behavior. This reflects real-world applications such as personalized POI recommendations in location-based services. In addition, the scale of our dataset allows us to examine how urban features influence POI recommendation accuracy. Building on prior hypotheses, we propose a new insight: cities with more evenly distributed POI categories tend to be harder to model, as the absence of a dominant POI category makes user behavior less predictable.

**Contribution** This paper introduces the Massive Semantic Trajectories for Understanding POI Check-ins (Massive-STEPS) dataset, addressing gaps in existing POI recommendation research. Current POI check-in datasets are often only from 2012-2013, skewed to a few cities, and lack semantic metadata, hindering the development of robust and globally applicable models. While datasets like GSCD and STD offer broad geographic coverage, they either suffer from an older timespan, contain erroneous data, or have missing information. Massive-STEPS overcomes these issues by providing high-quality check-ins from 2012-2013 and 2017-2018, improving temporal quality for longitudinal POI recommendation studies. The dataset spans 12 diverse cities across multiple regions, with a focus on low-resource cities overlooked in previous research. Additionally,
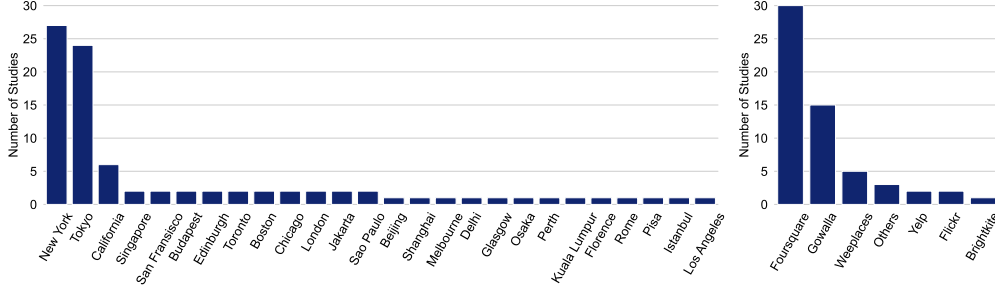
Figure 1: **Distribution of POI recommendation studies** modeled on specific cities, filtered from Table IV in [68]. We identified and counted studies that explicitly mentioned city names, revealing the skewness of existing research, which is saturated around New York and Tokyo. In addition, we include the distribution of studies by LBSN platform, showing that Foursquare is by far the most commonly used source of check-in data. The list of identified studies is shown in Table 5.

Massive-STEPS is enriched with metadata through alignment with Foursquare's Open Source Places, providing crucial details such as POI geographical coordinates, POI names, and addresses. We also conducted an extensive benchmark on both supervised and zero-shot POI recommendation tasks, evaluating a wide range of models, including traditional methods, deep learning approaches, and recent LLM-based techniques. We further analyzed which urban features affect POI recommendation accuracy and found that cities with no dominant POI category tend to be harder to predict. By releasing this dataset and benchmark code publicly, we facilitate open and reproducible research, enabling future advancements in POI recommendation studies.

## 2 Related Works

### 2.1 Existing Datasets

A survey conducted by [68] outlines the broad landscape of POI recommendation research, covering a wide range of models and architectures used in prior studies. While it offers a high-level overview of the datasets used, it lacks a dedicated discussion or evaluation of POI datasets. We address this gap by analyzing commonly used datasets in the field and positioning our dataset within this context.

**LBSN Check-in Data Sources**    Building on the tabular summary provided in [68], which offers a representative overview of the broader POI recommendation literature, we investigated which datasets are most commonly used in prior studies. From their original table (Table IV), we filtered entries pertaining specifically to POI and next POI recommendation tasks and identified (1) the most frequently used LBSN check-in data sources and (2) the most commonly studied cities. As shown in Fig. 1, Foursquare remains the dominant source of LBSN data in existing studies, appearing in almost 50% of the surveyed works. While several variants of Foursquare datasets have been employed, the most widely used are the NYC and Tokyo Dataset [57] (often abbreviated as FSQ-NYC and FSQ-TKY) and the Global-scale Check-in Dataset (GSCD) [55, 56], curated by the same authors. Other LBSN sources occasionally used include Gowalla [4], Brightkite [4], and Weeplaces [26].

**Saturated to Two Cities and Old Timespan**    Due to the widespread use of FSQ-NYC and FSQ-TKY [57], the majority of POI recommendation studies are disproportionately focused on these two cities, as illustrated in Fig. 1. While there is nothing inherently problematic about studying NYC and Tokyo, there has been growing interest in expanding research to a broader range of cities, particularly those that are underexplored or considered low-resource [63], as cultural and regional differences influence collective mobility behaviors. For instance, in some urban contexts, residents tend to commute to central business districts in the morning, whereas in others, nightlife activities such as visiting bars after work are more common [55]. Ensuring diverse geographic coverage is increasingly important, especially as LLMs are adopted for POI recommendation tasks. LLMs are known to exhibit geographical biases, often underperforming in regions with lower socioeconomic conditions [27]. It remains an open question whether LLM-based POI recommendation models can effectively generalize across diverse urban environments.

In addition, because many studies rely on the FSQ-NYC and FSQ-TKY, they are often constrained to the timespan it covers: check-in data from 2012 to 2013. However, POI data is inherently dynamic: venues may have closed, relocated, or changed in category over time. [61] underscores the importance of validating the temporal quality of POI datasets by recording whether and when a venue's information has been updated to reflect real-world changes. This is particularly critical, as recommender systems should avoid suggesting POIs that no longer exist or have undergone substantial changes (e.g., a former bookstore converted into a coworking space) and behave dynamically over longitudinal periods [53]. Moreover, behavioral patterns captured over a decade ago may no longer align with modern user preferences and routines. For example, the opening of a new train station may significantly shift commuting patterns and the popularity of surrounding POIs.

**Low Data Quality: Erroneous Entries**   More recently, researchers have begun leveraging the broader Global-scale Check-in Dataset (GSCD) [55, 56], which spans 415 cities across 77 countries. Despite its wider geographic coverage, GSCD is temporally limited to the same 2012-2013 period as FSQ-NYC and FSQ-TKY, and thus suffers from similar issues of temporal quality. More critically, [29] demonstrated that GSCD suffers from significant data quality issues, with over 14 million check-ins (about 44%) of the dataset flagged as erroneous due to anomalous user behavior. These include (1) repeated check-ins at the same venue, (2) check-ins occurring within implausibly short time intervals (less than one minute), and (3) transitions between venues that would require travel speeds exceeding Mach 1, which are physically unreasonable.

To address these limitations, [29] introduced the Semantic Trails Dataset (STD), which applies systematic filtering procedures to enhance data quality. STD comprises two subsets: a cleaned version of GSCD covering 2012-2013 (STD 2013), and a newer collection of check-ins from 2017-2018 (STD 2018), sourced from Foursquare Swarm. STD 2018 also spans a wider range of cities, making it valuable for capturing globally distributed user behavior, in contrast to GSCD's focus on densely populated urban centers. Both subsets follow the same rigorous filtering criteria, resulting in a higher-quality check-in dataset for POI recommendation tasks. Given these improvements, we adopted STD as the source for our check-in dataset.

**Poor Reproducibility**   Another persistent challenge in POI recommendation research is the lack of reproducibility in dataset preprocessing. While some recent studies utilize datasets like GSCD to cover a wide range of cities, they often omit important details needed for replicating their data filtering processes. For example, AgentMove [7] conducted city-level filtering based on a minimum distance between each trajectory and the city center, but they did not specify how the city center was defined or what distance threshold was used. Generally, city filtering can be done in two ways: (1) by selecting points within a fixed radius from a central coordinate, or (2) by applying precise administrative boundaries, such as those from OpenStreetMap. The radius-based approach is less reliable, as it assumes a circular city shape and requires a subjective definition of the city's center. In contrast, the boundary-based method is more robust and reproducible, and is the approach we adopted in this work. Without standardized, transparent filtering criteria and publicly available datasets, it remains difficult to ensure fair comparisons across POI recommendation methods.

## 2.2   Understanding Urban Features and POI Recommendation

POI recommendation studies that evaluate models across multiple city-level datasets often include analyses to assess how well their methods generalize across different urban contexts. It is well understood that POI recommendation accuracy metrics (e.g., Acc@k, NDCG@k) can vary substantially between cities and can be interpreted as a proxy for how easy or difficult a city is to model. The assumption is that higher performance reflects more predictable or structured mobility patterns. This viewpoint is consistent with prior work highlighting the role of cultural and urban-specific factors in shaping mobility behaviors [55, 34].

Several studies have proposed hypotheses connecting specific urban features to modeling difficulty. For example, GETNext [60] hypothesized that cities with fewer check-ins and higher spatial sparsity of POIs are harder to model. STHGCN [54] suggested that a larger number of user trajectories improves predictive accuracy by providing richer collaborative signals, whose architecture is designed to leverage. LLM4POI [22] proposed that cities with a greater variety of POI categories are easier to model due to LLMs' contextual reasoning capabilities, whereas cities covering a broader geographic area tend to be more difficult to model. In the zero-shot POI recommendation setting, AgentMove [7]

4

Table 2: **Summary statistics** of the 12 Massive-STEPS subsets, including the number of users, trajectories, POI locations, total check-ins, and train, validation, and test sample counts. For comparison, we also include statistics from existing Foursquare-based [57] and Gowalla-based [4] datasets. [†]Due to variations in dataset preprocessing across studies, we report the version used in [54].

| City | Users | Trajectories | POIs | Check-ins | #train | #val | #test |
|------|-------|--------------|------|-----------|--------|------|-------|
| **NYC and Tokyo Check-in Dataset[†] [57]** | | | | | | | |
| New York | 1,048 | 14,130 | 4,981 | 103,941 | 72,206 | 1,400 | 1,347 |
| Tokyo | 2,282 | 65,499 | 7,833 | 405,000 | 274,597 | 6,868 | 7,038 |
| **Gowalla[†] [4, 62]** | | | | | | | |
| California | 3,957 | 45,123 | 9,690 | 238,369 | 154,253 | 3,529 | 2,780 |
| **Massive-STEPS** | | | | | | | |
| Beijing | 56 | 573 | 1,127 | 1,470 | 400 | 58 | 115 |
| Istanbul | 23,700 | 216,411 | 53,812 | 544,471 | 151,487 | 21,641 | 43,283 |
| Jakarta | 8,336 | 137,396 | 76,116 | 412,100 | 96,176 | 13,740 | 27,480 |
| Kuwait City | 9,628 | 91,658 | 17,180 | 232,706 | 64,160 | 9,166 | 18,332 |
| Melbourne | 646 | 7,864 | 7,699 | 22,050 | 5,504 | 787 | 1,573 |
| Moscow | 3,993 | 39,485 | 17,822 | 105,620 | 27,639 | 3,949 | 7,897 |
| New York | 6,929 | 92,041 | 49,218 | 272,368 | 64,428 | 9,204 | 18,409 |
| Petaling Jaya | 14,308 | 180,410 | 60,158 | 506,430 | 126,287 | 18,041 | 36,082 |
| São Paulo | 5,822 | 89,689 | 38,377 | 256,824 | 62,782 | 8,969 | 17,938 |
| Shanghai | 296 | 3,636 | 4,462 | 10,491 | 2,544 | 364 | 728 |
| Sydney | 740 | 10,148 | 8,986 | 29,900 | 7,103 | 1,015 | 2,030 |
| Tokyo | 764 | 5,482 | 4,725 | 13,839 | 3,836 | 549 | 1,097 |

reported two key findings: (1) geospatial biases inherent in LLMs can hinder prediction quality across cities, and (2) LLMs are influenced by city-specific mobility patterns.

Building on these insights, we used Massive-STEPS to explore how urban features affect POI recommendations. Its diverse set of 12 cities allows for a comprehensive analysis across different cultural and urban contexts. We analyzed the correlation between urban features and model accuracy, and based on the results, proposed a new hypothesis that contrasts previous findings in the literature.

# 3 Dataset

## 3.1 Creation Process

Massive-STEPS is derived from the two subsets of STD [29], incorporating check-ins from both the 2013 and 2018 subsets. We utilize two additional components from STD: (1) the **cities** metadata file, which provides the latitude and longitude of administrative regions (e.g., towns, suburbs) along with their corresponding country codes obtained from GeoNames; and (2) the POI **category** mapping, which links each Foursquare Category ID to its descriptive name (e.g., "Restaurant"). Based on this metadata, each POI is thus associated with several attributes: Foursquare Place ID, Foursquare Category ID, category name, latitude/longitude of the administrative region, the administrative region name, and the country code. For anonymization purposes and model training compatibility, we apply ordinal encoding to the Place IDs and Category IDs, assigning each a unique integer index.

### 3.1.1 Preprocessing

**Trajectory Grouping**   Most POI recommendation models operate on sequences of check-ins, commonly referred to as trajectories. The model is tasked with predicting the next POIs a user is likely to visit, given the current trajectory. STD conveniently provides pre-grouped trajectories (trails) by applying a time interval-based grouping: for each user, check-ins that occur within a time interval of $\delta_\tau = 8$ hours are grouped into the same trajectory.

**Matching Trajectories to Target Cities**   To obtain city-specific datasets, we matched trajectories to the target cities. For each city, we obtain geographic boundaries from OpenStreetMap and retrieve its GeoJSON file via the Overpass API. The GeoJSON file contains a polygon defining the city's boundary in latitude and longitude. Using this boundary, we filter check-ins by comparing the
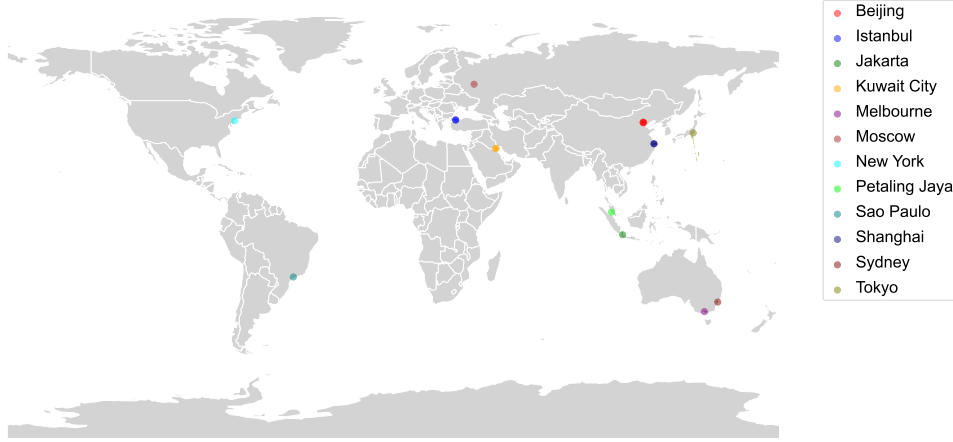
Figure 2: **World map highlighting the cities included in the Massive-STEPS dataset**.

latitude/longitude of each POI's administrative region and retain only those that are within the city's polygon. This ensures that all retained trajectories are spatially grounded within the designated city.

**Filtering Short Trajectories and Inactive Users**　To ensure data quality, we apply an additional filtering step by removing trajectories with fewer than two check-ins and excluding users with fewer than three trajectories. This prevents the model from learning from overly sparse or irrelevant data.

**Train, Validation, and Test Splits**　We split trajectories into training, validation, and test sets in a ratio of 7:1:2, following [7]. We ensure that all users in the test set appear at least once in the training or validation set, a common practice in prior studies [60, 54].

### 3.1.2　POI Enrichment via Foursquare OS Places

Since the POIs in STD include their corresponding Foursquare Place IDs, we matched them directly with entries in the Foursquare OS Places dataset using these IDs as the key. This one-to-one ID correspondence allows for a straightforward join operation, enriching each POI with additional metadata such as its precise latitude and longitude, name (e.g., of a restaurant or subway station), and address. However, not all POIs in the Foursquare OS Places dataset include the full metadata, particularly those categorized as private residences, which are excluded due to privacy restrictions.

### 3.2　Description and Addressing Literature Gaps

**Massive Semantic Trajectories for Understanding POI Check-ins (Massive-STEPS) dataset** is a city-level POI check-in dataset derived from Semantic Trails [29], comprising user check-in trajectories from 12 cities: Beijing, Istanbul, Jakarta, Kuwait City, Melbourne, Moscow, New York, Petaling Jaya, São Paulo, Shanghai, Sydney, and Tokyo. It features anonymized POI check-ins enriched with geographical metadata to support spatiotemporal and sequential modeling tasks. City-level statistics, along with comparisons to existing datasets, are presented in Table 2. Fig. 2 shows a world map highlighting the locations of all cities included in the Massive-STEPS dataset. Table 6 shows the available fields in the dataset and provides an example for each field.

**Massive-STEPS offers a more comprehensive and diverse representation of urban mobility** compared to typical POI check-in datasets. As shown in Table 2, datasets like FSQ-NYC and FSQ-TKY [57] contain fewer than 10,000 candidate POI locations. In contrast, cities in Massive-STEPS cover significantly more POIs: Massive-STEPS New York has over 49,000 POIs, while Massive-STEPS Jakarta exceeds 76,000. Massive-STEPS Istanbul, one of the largest subsets, features a large user base of 23,700, offering a broad range of user behaviors. Although some Massive-STEPS subsets are smaller than their FSQ counterparts (e.g., Tokyo), we attribute this to the strict filtering procedures applied by STD to remove erroneous entries, as explained in Section 2.1. This scale introduces

6

Table 3: **Benchmark results on POI recommendation task**. The metric reported is Acc@1. Full results, including other metrics, are available in Section C.3. **Bold** indicates the best performance for each city, while underline indicates the second-best.

| Model | Beijing | Istanbul | Jakarta | Kuwait City | Melbourne | Moscow | New York | Petaling Jaya | São Paulo | Shanghai | Sydney | Tokyo |
|-------|---------|----------|---------|-------------|-----------|--------|----------|---------------|-----------|----------|--------|-------|
| **FPMC** | 0.000 | 0.026 | 0.029 | 0.021 | 0.062 | 0.059 | 0.032 | 0.026 | 0.030 | 0.084 | 0.075 | 0.176 |
| **RNN** | 0.085 | 0.077 | 0.049 | 0.087 | 0.059 | 0.075 | 0.061 | 0.064 | 0.097 | 0.055 | 0.080 | 0.133 |
| **LSTPM** | 0.127 | 0.142 | 0.099 | 0.180 | 0.091 | 0.151 | 0.099 | 0.099 | 0.158 | 0.099 | 0.141 | 0.225 |
| **DeepMove** | 0.106 | 0.150 | 0.103 | 0.179 | 0.083 | 0.143 | 0.097 | 0.112 | 0.160 | 0.085 | 0.129 | 0.201 |
| **GETNext** | 0.433 | 0.146 | 0.155 | 0.175 | 0.100 | 0.175 | 0.134 | 0.139 | 0.202 | 0.115 | 0.181 | 0.180 |
| **STHGCN** | 0.453 | 0.241 | 0.197 | 0.225 | 0.168 | 0.223 | 0.146 | 0.174 | 0.250 | 0.193 | 0.227 | 0.250 |

additional computational challenges. Models that rely on dense POI-to-POI adjacency matrices, for instance, may require more efficient implementations to avoid excessive memory consumption.

Beyond scale, Massive-STEPS addresses the oversaturation of FSQ-NYC and FSQ-TKY in POI recommendation research. Notably, Massive-STEPS includes low-resource and previously under-explored cities in human mobility studies, such as Petaling Jaya and Kuwait City, both of which are among the cities with the highest number of check-ins from STD. This broader coverage opens new research opportunities for studying location-based behaviors across diverse cultural and geographic contexts. Furthermore, since Massive-STEPS is based on STD, it benefits from the carefully filtered, high-quality check-ins and a longer, more recent timespan. These characteristics make Massive-STEPS a more relevant and reliable resource for modeling human mobility patterns.

**Massive-STEPS is designed to be easily extended to other geographical regions**. Since the data processing code is open-source and fully reproducible, adding a new city only requires its geographic boundaries from OpenStreetMap. Moreover, Massive-STEPS is scalable to higher levels of geographic granularity, enabling the creation of provincial, state, and country-level POI check-in datasets, which support collective mobility studies at broader geographic scales.

## 4 Benchmark

### 4.1 POI Recommendation

This benchmark covers the common task of POI recommendation, where the goal is to predict where a user will go next based on their previous check-ins. The input to the model is a trajectory, which is a sequence of places the user has visited. The model is expected to suggest a set of $K$ POIs that the user might visit next. This task can also be framed as re-ranking a set of $N$ candidate POIs in the city based on the user's recent behavior. This task reflects a typical real-world use case in POI recommendation systems and human mobility modeling. It is a **supervised** task as the model is trained using all available historical trajectories of the user, allowing it to learn personalized patterns of movement over time. Appendix C provides the details on problem formulation, hyperparameters, experimental setups, and the full results with all evaluation metrics.

**Dataset Preparation** By design, our dataset is structured to support consistency and reproducibility, making it straightforward to use for supervised POI recommendation tasks. We adopted the predefined trajectories from the original STD, which has grouped sequences of check-ins into trajectories based on a fixed time interval (see Section 3.1.1). Since all input features in our dataset have been numerically encoded to facilitate model training, we simply used them across all experiments.

For this task, models typically use four kinds of features: (1) social features: user ID; (2) spatial features: POI ID, geographic coordinates; (3) temporal features: check-in timestamp; and (4) categorical features: POI category. Since not all POIs have exact geographic coordinates (see Section 3.1.2), we use the geographic coordinates of their administrative region as a proxy for all POIs.

**Models** We evaluated three kinds of models: (1) Markov-based methods, (2) classical deep learning models, and (3) Transformer-based graph neural networks:

- **FPMC** [32]: A classical baseline that combines first-order Markov chains with matrix factorization to model personalized next-location predictions.

- **RNN** [41], **LSTPM** [35], and **DeepMove** [8]: Recurrent neural networks designed to capture sequential dependencies, with varying mechanisms to incorporate spatio-temporal context.
- **GETNext** [60] and **STHGCN** [54]: Transformer-based graph neural networks to model social, spatial, and temporal dependencies.

**Evaluation Metrics**   We used two popular metrics in POI recommender systems: Acc@k, which checks if the true POI appears in the top-k predicted results, and NDCG@k, which measures the ranking quality of the suggested results.

**Results**   As shown in Table 3, STHGCN consistently achieves the highest accuracy across all cities, highlighting the effectiveness of Transformer-based GNNs. GETNext often ranks second, performing well in most cities, and older RNN-based models such as LSTPM and DeepMove also remain competitive in a few cities. The best-performing model achieves an average Acc@1 of 22.9%, which is comparable to results from previous studies using datasets of a similar size [7].

We investigated which urban features affect POI recommendation accuracy by computing Spearman correlations between various city features and model performance. We found that **category entropy**, based on Shannon entropy, strongly correlates with recommendation accuracy ($r = -0.736$), as shown in Fig. 6. Namely, cities with more evenly distributed POI categories and no dominant POI category tend to have lower accuracy and are thus harder to predict. This finding is consistent with the results of prior work on other datasets. Further details are provided in Appendix D.

### 4.2   Zero-shot POI Recommendation

An increasingly popular approach in POI recommendation leverages the zero-shot capabilities of LLMs, allowing models to perform zero-shot POI recommendation without any additional training. This setting challenges LLMs to generalize to any target city and personalize recommendations for any user. The task mirrors the supervised counterpart: given a user's trajectory, the model predicts the user's next destination by ranking a set of $K$ POIs. Appendix E provides the details on problem formulation, prompts, experimental setups, and the full results with all evaluation metrics.

**Dataset Preparation**   Since LLMs are central for this task, trajectories must be transformed into textual prompts [52, 51]. We adapted the prompts used in [7], which implemented the three LLM methods evaluated in our study. Notably, because LLMs can effectively leverage contextual information, not all input features need to be numerically encoded. The evaluated methods utilize only a subset of features: the timestamp, POI category name, and POI ID.

**Methods and Models**   We evaluated three LLM-based prompting methods:

- **LLM-Mob** [42]: One of the earliest methods to use LLMs for next POI prediction, prompting LLMs with both historical and current (contextual) trajectories.
- **LLM-ZS** [1]: A simplified version of LLM-Mob that retains the use of historical and contextual trajectories but simplifies its prompt design.
- **LLM-Move** [10]: Extends previous prompting methods by introducing a RAG-like approach, retrieving nearby POIs as candidates, and ranking them by geographic distance to the user's most recent visit.

To ensure a more robust evaluation, we tested each method using four different LLMs: one closed-source API (Gemini 2.0 Flash [36]) and three open-source models (Qwen 2.5 7B [38], Llama 3.1 8B [12], and Gemma 2 9B [37]). All open-source models are instruction-tuned variants, quantized to INT4 using AWQ [25], and served through vLLM [21].

**Evaluation Metrics**   We used the same metrics as in the supervised setting, Acc@k and NDCG@k, for the same reasons: to assess predictive accuracy and to evaluate the quality of the ranking.

**Results**   As shown in Table 4, LLM-Move [10] performed the best among the three methods, likely due to its well-crafted prompts that supply potential candidate POIs, rather than relying solely on historical or contextual trajectories like LLM-Mob and LLM-ZS. This also supports the use of

Table 4: **Benchmark results on zero-shot POI recommendation task**. The metric reported is Acc@1. Full results, including other metrics, are available in Section E.3. **Bold** indicates the best performance for each city, while <u>underline</u> indicates the second-best.

| Method | LLM | Beijing | Istanbul | Jakarta | Kuwait City | Melbourne | Moscow | New York | Petaling Jaya | São Paulo | Shanghai | Sydney | Tokyo |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **LLM-Mob** | Gemini 2 Flash | <u>0.115</u> | 0.080 | 0.100 | 0.095 | 0.060 | 0.130 | 0.095 | 0.090 | 0.130 | 0.055 | 0.060 | 0.140 |
| | Qwen 2.5 7B | 0.058 | 0.035 | 0.105 | 0.080 | 0.030 | 0.090 | 0.070 | 0.030 | 0.090 | 0.040 | 0.035 | 0.110 |
| | Llama 3.1 8B | 0.000 | 0.020 | 0.055 | 0.030 | 0.010 | 0.030 | 0.025 | 0.010 | 0.030 | 0.005 | 0.020 | 0.005 |
| | Gemma 2 9B | <u>0.115</u> | 0.075 | 0.105 | 0.080 | 0.055 | 0.100 | 0.070 | 0.055 | 0.085 | 0.050 | 0.030 | 0.145 |
| **LLM-ZS** | Gemini 2 Flash | 0.058 | 0.090 | 0.110 | 0.080 | 0.065 | 0.125 | 0.080 | 0.110 | 0.150 | 0.065 | 0.060 | <u>0.160</u> |
| | Qwen 2.5 7B | 0.038 | 0.040 | 0.065 | 0.050 | 0.040 | 0.080 | 0.050 | 0.045 | 0.095 | 0.045 | 0.045 | 0.120 |
| | Llama 3.1 8B | 0.077 | 0.040 | 0.045 | 0.060 | 0.040 | 0.080 | 0.055 | 0.030 | 0.030 | 0.060 | 0.040 | 0.110 |
| | Gemma 2 9B | 0.096 | 0.045 | 0.105 | 0.070 | 0.050 | 0.080 | 0.075 | 0.065 | 0.075 | 0.050 | 0.045 | 0.110 |
| **LLM-Move** | Gemini 2 Flash | 0.096 | **0.205** | **0.295** | **0.220** | **0.225** | <u>0.220</u> | **0.235** | **0.210** | **0.285** | **0.170** | **0.230** | **0.250** |
| | Qwen 2.5 7B | **0.192** | <u>0.175</u> | 0.115 | <u>0.160</u> | 0.110 | **0.230** | <u>0.120</u> | 0.135 | 0.155 | 0.095 | <u>0.125</u> | **0.250** |
| | Llama 3.1 8B | 0.058 | <u>0.015</u> | 0.015 | <u>0.010</u> | 0.040 | 0.005 | <u>0.035</u> | 0.040 | 0.045 | 0.020 | <u>0.055</u> | 0.030 |
| | Gemma 2 9B | 0.096 | 0.100 | <u>0.235</u> | 0.120 | <u>0.115</u> | 0.110 | 0.115 | <u>0.175</u> | <u>0.195</u> | <u>0.105</u> | 0.125 | 0.130 |

administrative region coordinates as a proxy when exact POI coordinates are unavailable (e.g., due to privacy issues). Across different LLMs, Gemini 2.0 Flash achieved the highest accuracy across all three prompting strategies, with Qwen 2.5 7B and Gemma 2 9B following closely as strong open-source alternatives. In contrast, we found that Llama 3.1 8B often struggled to follow prompt instructions and frequently produced irrelevant predictions, especially when using the prompts of LLM-Mob and LLM-Move.

Notably, these zero-shot methods outperformed, matched, or came close to the performance of supervised baselines in several cities (e.g., Jakarta, Kuwait City, Moscow), demonstrating their effectiveness even without additional fine-tuning. These results indicate that LLMs can effectively leverage the contextual information included in trajectory prompts to guide their predictions. Since LLMs are pre-trained, they do not require additional training, which reduces both computation time and cost. While inference may require more powerful hardware, it can still be faster overall than training supervised models from scratch.

## 5 Limitations and Future Work

Firstly, Massive-STEPS is derived from the Semantic Trails dataset and thus inherits its biases and potential errors, which may propagate through downstream tasks. Additionally, the dataset is sparse in several cities, which can impact model training quality and limit cross-city generalization. Secondly, Massive-STEPS focuses solely on trajectories and POI metadata, without including user demographic or social information (e.g., age, social connections). This restricts its applicability for personalized or socially-aware POI recommendation tasks. Thirdly, while our benchmarking covers a wide range of models and cities to emphasize replicability and geographic breadth, we did not perform extensive hyperparameter tuning, which may affect the peak performance of the evaluated models.

While the current version of Massive-STEPS already surpasses existing POI check-in datasets in scale and diversity, we envision expanding it further to include even more cities worldwide. Our dataset creation pipeline is designed to be easily extensible, allowing researchers to integrate additional cities with minimal effort. This allows Massive-STEPS to evolve into a truly global POI check-in dataset.

## 6 Conclusion

In this paper, we presented the Massive-STEPS dataset to address longstanding limitations in POI recommendation research, particularly the reliance on older, geographically saturated, and non-reproducible check-in datasets. Massive-STEPS offers a large-scale, semantically enriched resource spanning 12 cities across diverse global regions and two time periods, supporting both longitudinal and cross-city analyses. The dataset includes rich semantic information such as venue name, address, category, and coordinates. We also provide benchmark results for both supervised and zero-shot POI recommendation methods, illustrating the dataset's utility across model types and urban contexts. By releasing Massive-STEPS and our evaluation pipeline publicly, we aim to advance open, reproducible, and globally inclusive research in human mobility and POI recommendation systems. We released our dataset under the Apache 2.0 license.

## Acknowledgments and Disclosure of Funding

## References

[1] C. Beneduce, B. Lepri, and M. Luca. Large language models are zero-shot next location predictors. *arXiv preprint arXiv:2405.20962*, 2024.

[2] G. Cao, S. Cui, and I. Joe. Improving the spatial–temporal aware attention network with dynamic trajectory graph learning for next point-of-interest recommendation. *Information Processing & Management*, 60(3):103335, 2023.

[3] L. Chen and G. Zhu. Self-supervised contrastive learning for itinerary recommendation. *Expert Systems with Applications*, 268:126246, 2025.

[4] E. Cho, S. A. Myers, and J. Leskovec. Friendship and mobility: user movement in location-based social networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1082–1090, 2011.

[5] M. Davtalab and A. A. Alesheikh. A poi recommendation approach integrating social spatio-temporal information into probabilistic matrix factorization. *Knowledge and Information Systems*, 63:65–85, 2021.

[6] J. Ding, G. Yu, Y. Li, D. Jin, and H. Gao. Learning from hometown and current city: Cross-city poi recommendation via interest drift and transfer learning. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 3(4), Sept. 2020.

[7] J. Feng, Y. Du, J. Zhao, and Y. Li. AgentMove: A large language model based agentic framework for zero-shot next location prediction. In L. Chiruzzo, A. Ritter, and L. Wang, editors, *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1322–1338, Albuquerque, New Mexico, Apr. 2025. Association for Computational Linguistics.

[8] J. Feng, Y. Li, C. Zhang, F. Sun, F. Meng, A. Guo, and D. Jin. Deepmove: Predicting human mobility with attentional recurrent networks. In *Proceedings of the 2018 World Wide Web Conference*, WWW '18, page 1459–1468, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee.

[9] J. Feng, Z. Yang, F. Xu, H. Yu, M. Wang, and Y. Li. Learning to simulate human mobility. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '20, page 3426–3433, New York, NY, USA, 2020. Association for Computing Machinery.

[10] S. Feng, H. Lyu, F. Li, Z. Sun, and C. Chen. Where to move next: Zero-shot generalization of llms for next poi recommendation. In *2024 IEEE Conference on Artificial Intelligence (CAI)*, pages 1530–1535. IEEE, 2024.

[11] S. Feng, F. Meng, L. Chen, S. Shang, and Y. S. Ong. Rotan: A rotation-based temporal attention network for time-specific next poi recommendation. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 759–770, 2024.

[12] A. Grattafiori, A. Dubey, A. Jauhri, and et al. The llama 3 herd of models, 2024.

[13] S. Halder, K. H. Lim, J. Chan, and X. Zhang. Transformer-based multi-task learning for queuing time aware next poi recommendation. In *Pacific-Asia conference on knowledge discovery and data mining*, pages 510–523. Springer, 2021.

[14] S. Halder, K. H. Lim, J. Chan, and X. Zhang. Capacity-aware fair poi recommendation combining transformer neural networks and resource allocation policy. *Applied Soft Computing*, 147:110720, 2023.

[15] H. Han, M. Zhang, M. Hou, F. Zhang, Z. Wang, E. Chen, H. Wang, J. Ma, and Q. Liu. Stgcn: a spatial-temporal aware graph learning method for poi recommendation. In *2020 IEEE International Conference on Data Mining (ICDM)*, pages 1052–1057. IEEE, 2020.

[16] N. L. Ho and K. H. Lim. Poibert: A transformer-based model for the tour recommendation problem. In *2022 IEEE International Conference on Big Data (Big Data)*, pages 5925–5933. IEEE, 2022.

[17] M. B. Hossain, M. S. Arefin, I. H. Sarker, M. Kowsher, P. K. Dhar, and T. Koshiba. Caran: A context-aware recency-based attention network for point-of-interest recommendation. *IEEE Access*, 10:36299–36310, 2022.

[18] S. Jiang, W. He, L. Cui, Y. Xu, and L. Liu. Modeling long-and short-term user preferences via self-supervised learning for next poi recommendation. *ACM Transactions on Knowledge Discovery from Data*, 17(9):1–20, 2023.

[19] S. Jiang and J. Wu. Temporal-geographical attention-based transformer for point-of-interest recommendation. *Journal of Intelligent & Fuzzy Systems*, 45(6):12243–12253, 2023.

[20] W. JIAWEI, R. Jiang, C. Yang, Z. Wu, R. Shibasaki, N. Koshizuka, C. Xiao, et al. Large language models as urban residents: An llm agent framework for personal mobility generation. *Advances in Neural Information Processing Systems*, 37:124547–124574, 2024.

[21] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. E. Gonzalez, H. Zhang, and I. Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

[22] P. Li, M. de Rijke, H. Xue, S. Ao, Y. Song, and F. D. Salim. Large language models for next point-of-interest recommendation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, page 1463–1472, New York, NY, USA, 2024. Association for Computing Machinery.

[23] R. Li, J. Guo, C. Liu, Z. Li, and S. Zhang. Using attributes explicitly reflecting user preference in a self-attention network for next poi recommendation. *ISPRS International Journal of Geo-Information*, 11(8):440, 2022.

[24] Y. Li, T. Chen, P.-F. Zhang, Z. Huang, and H. Yin. Self-supervised graph-based point-of-interest recommendation. *arXiv preprint arXiv:2210.12506*, 2022.

[25] J. Lin, J. Tang, H. Tang, S. Yang, W.-M. Chen, W.-C. Wang, G. Xiao, X. Dang, C. Gan, and S. Han. Awq: Activation-aware weight quantization for llm compression and acceleration, 2024.

[26] Y. Liu, T.-A. N. Pham, G. Cong, and Q. Yuan. An experimental evaluation of point-of-interest recommendation in location-based social networks. *Proc. VLDB Endow.*, 10(10):1010–1021, June 2017.

[27] R. Manvi, S. Khanna, M. Burke, D. Lobell, and S. Ermon. Large language models are geographically biased. *arXiv preprint arXiv:2402.02680*, 2024.

[28] P. Merinov and F. Ricci. Positive-sum impact of multistakeholder recommender systems for urban tourism promotion and user utility. In *Proceedings of the 18th ACM Conference on Recommender Systems*, pages 939–944, 2024.

[29] D. Monti, E. Palumbo, G. Rizzo, R. Troncy, T. Ehrhart, and M. Morisio. Semantic trails of city explorations: How do we live a city. *arXiv preprint arXiv:1812.04367*, 2018.

[30] Y. Qin, Y. Fang, H. Luo, F. Zhao, and C. Wang. Next point-of-interest recommendation with auto-correlation enhanced multi-modal transformer network. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2612–2616, 2022.

[31] Y. Qin, H. Wu, W. Ju, X. Luo, and M. Zhang. A diffusion model for poi recommendation. *ACM Trans. Inf. Syst.*, 42(2), Nov. 2023.

[32] S. Rendle, C. Freudenthaler, and L. Schmidt-Thieme. Factorizing personalized markov chains for next-basket recommendation. In *Proceedings of the 19th International Conference on World Wide Web*, WWW '10, page 811–820, New York, NY, USA, 2010. Association for Computing Machinery.

[33] F. D. Salim, B. Dong, M. Ouf, Q. Wang, I. Pigliautile, X. Kang, T. Hong, W. Wu, Y. Liu, S. K. Rumi, M. S. Rahaman, J. An, H. Deng, W. Shao, J. Dziedzic, F. C. Sangogboye, M. B. Kjærgaard, M. Kong, C. Fabiani, A. L. Pisello, and D. Yan. Modelling urban-scale occupant behaviour, mobility, and energy in buildings: A survey. *Building and Environment*, 183:106964, 2020.

[34] K. Sun, C. Li, and T. Qian. City matters! a dual-target cross-city sequential poi recommendation model. *ACM Transactions on Information Systems*, 42(6):1–27, 2024.

[35] K. Sun, T. Qian, T. Chen, Y. Liang, Q. V. H. Nguyen, and H. Yin. Where to go next: Modeling long- and short-term user preferences for point-of-interest recommendation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(01):214–221, Apr. 2020.

[36] G. Team and et al. Gemini: A family of highly capable multimodal models, 2024.

[37] G. Team and et al. Gemma 2: Improving open language models at a practical size, 2024.

[38] Q. Team. Qwen2.5: A party of foundation models, September 2024.

[39] D. Wang, C. Chen, C. Di, and M. Shu. Exploring behavior patterns for next-poi recommendation via graph self-supervised learning. *Electronics*, 12(8):1939, 2023.

[40] D. Wang, F. Wan, D. Yu, Y. Shen, Z. Xiang, and Y. Xu. Context-and category-aware double self-attention model for next poi recommendation. *Applied Intelligence*, 53(15):18355–18380, 2023.

[41] J. Wang, J. Jiang, W. Jiang, C. Li, and W. X. Zhao. Libcity: An open library for traffic prediction. In *Proceedings of the 29th International Conference on Advances in Geographic Information Systems*, SIGSPATIAL '21, page 145–148, New York, NY, USA, 2021. Association for Computing Machinery.

[42] X. Wang, M. Fang, Z. Zeng, and T. Cheng. Where would i go next? large language models as human mobility predictors. *arXiv preprint arXiv:2308.15197*, 2023.

[43] X. Wang, Y. Liu, X. Zhou, Z. Leng, and X. Wang. Long-and short-term preference modeling based on multi-level attention for next poi recommendation. *ISPRS International Journal of Geo-Information*, 11(6):323, 2022.

[44] Y. Wang, A. Liu, J. Fang, J. Qu, and L. Zhao. Adq-gnn: Next poi recommendation by fusing gnn and area division with quadtree. In *Web Information Systems Engineering–WISE 2021: 22nd International Conference on Web Information Systems Engineering, WISE 2021, Melbourne, VIC, Australia, October 26–29, 2021, Proceedings, Part II 22*, pages 177–192. Springer, 2021.

[45] Z. Wang, J. Zeng, L. Zhong, L. Liu, M. Gao, and J. Wen. Dsdrec: Next poi recommendation using deep semantic extraction and diffusion model. *Information Sciences*, 678:121004, 2024.

[46] Y. Wu, X. Jiao, Q. Hao, Y. Xiao, and W. Zheng. Dlan: Modeling user long-and short-term preferences based on double-layer attention network for next point-of-interest recommendation. *Journal of Intelligent & Fuzzy Systems*, 46(2):3307–3321, 2024.

[47] Y. Wu, G. Zhao, M. Li, Z. Zhang, and X. Qian. Reason generation for point of interest recommendation via a hierarchical attention-based transformer model. *IEEE Transactions on Multimedia*, 26:5511–5522, 2023.

[48] J. Xia, Y. Yang, S. Wang, H. Yin, J. Cao, and P. S. Yu. Bayes-enhanced multi-view attention networks for robust poi recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 36(7):2895–2909, 2023.

[49] J. Xie and Z. Chen. Hierarchical transformer with spatio-temporal context aggregation for next point-of-interest recommendation. *ACM Transactions on Information Systems*, 42(2):1–30, 2023.

[50] X. Xu, T. Suzumura, J. Yong, M. Hanai, C. Yang, H. Kanezashi, R. Jiang, and S. Fukushima. Revisiting mobility modeling with graph: A graph transformer model for next point-of-interest recommendation. In *Proceedings of the 31st ACM international conference on advances in geographic information systems*, pages 1–10, 2023.

[51] H. Xue and F. D. Salim. Promptcast: A new prompt-based learning paradigm for time series forecasting. *IEEE Trans. on Knowl. and Data Eng.*, 36(11):6851–6864, Nov. 2024.

[52] H. Xue, B. P. Voutharoja, and F. D. Salim. Leveraging language foundation models for human mobility forecasting. In *Proceedings of the 30th International Conference on Advances in Geographic Information Systems*, SIGSPATIAL '22, New York, NY, USA, 2022. Association for Computing Machinery.

[53] T. Yabe, K. Tsubouchi, T. Shimizu, Y. Sekimoto, K. Sezaki, E. Moro, and A. Pentland. Yj-mob100k: City-scale and longitudinal dataset of anonymized human mobility trajectories. *Scientific Data*, 11(1):397, Apr 2024.

[54] X. Yan, T. Song, Y. Jiao, J. He, J. Wang, R. Li, and W. Chu. Spatio-temporal hypergraph learning for next poi recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '23, page 403–412, New York, NY, USA, 2023. Association for Computing Machinery.

[55] D. Yang, D. Zhang, L. Chen, and B. Qu. Nationtelescope: Monitoring and visualizing large-scale collective behavior in lbsns. *Journal of Network and Computer Applications*, 55:170–180, 2015.

[56] D. Yang, D. Zhang, and B. Qu. Participatory cultural mapping based on collective behavior data in location-based social networks. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 7(3):1–23, 2016.

[57] D. Yang, D. Zhang, V. W. Zheng, and Z. Yu. Modeling user activity preference by leveraging user spatial temporal characteristics in lbsns. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 45(1):129–142, 2014.

[58] Q. Yang, S. Hu, W. Zhang, and J. Zhang. Attention mechanism and adaptive convolution actuated fusion network for next poi recommendation. *International Journal of Intelligent Systems*, 37(10):7888–7908, 2022.

[59] S. Yang, J. Liu, and K. Zhao. Getnext: Trajectory flow map enhanced transformer for next poi recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '22, page 1144–1153, New York, NY, USA, 2022. Association for Computing Machinery.

[60] S. Yang, J. Liu, and K. Zhao. Getnext: trajectory flow map enhanced transformer for next poi recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on research and development in information retrieval*, pages 1144–1153, 2022.

[61] L. W. Yeow, R. Low, Y. X. Tan, and L. Cheah. Point-of-interest (poi) data validation methods: An urban case study. *ISPRS International Journal of Geo-Information*, 10(11):735, 2021.

[62] Q. Yuan, G. Cong, Z. Ma, A. Sun, and N. M. Thalmann. Time-aware point-of-interest recommendation. In *Proceedings of the 36th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '13, page 363–372, New York, NY, USA, 2013. Association for Computing Machinery.

[63] Y. Yuan, Y. Zhang, J. Ding, and Y. Li. Worldmove, a global open data for human mobility. *arXiv preprint arXiv:2504.10506*, 2025.

[64] H. Zang, D. Han, X. Li, Z. Wan, and M. Wang. Cha: Categorical hierarchy-based attention for next poi recommendation. *ACM Transactions on Information Systems (TOIS)*, 40(1):1–22, 2021.

[65] J. Zhang, Y. Li, R. Zou, J. Zhang, R. Jiang, Z. Fan, and X. Song. Hyper-relational knowledge graph neural network for next poi recommendation. *World Wide Web*, 27(4), July 2024.

[66] J. Zhang, Y. Li, R. Zou, J. Zhang, R. Jiang, Z. Fan, and X. Song. Hyper-relational knowledge graph neural network for next poi recommendation. *World Wide Web*, 27(4):46, 2024.

[67] J. Zhang and W. Ma. Hybrid structural graph attention network for poi recommendation. *Expert Systems with Applications*, 248:123436, 2024.

[68] Q. Zhang, P. Yang, J. Yu, H. Wang, X. He, S.-M. Yiu, and H. Yin. A survey on point-of-interest recommendation: Models, architectures, and security. *IEEE Transactions on Knowledge and Data Engineering*, 2025.

[69] P. Zhao, H. Zhu, Y. Liu, Z. Li, J. Xu, and V. S. Sheng. Where to go next: A spatio-temporal lstm model for next poi recommendation. *arXiv preprint arXiv:1806.06671*, 2018.

[70] H. Zhong, W. He, L. Cui, L. Liu, Z. Yan, and K. Zhao. Joint attention networks with inherent and contextual preference-awareness for successive poi recommendation. *Data Science and Engineering*, 7(4):370–382, 2022.

[71] Z. Zhou, Y. Lin, D. Jin, and Y. Li. Large language model for participatory urban planning. *arXiv preprint arXiv:2402.17161*, 2024.

[72] J. Zuo and Y. Zhang. Diff-dgmn: A diffusion-based dual graph multi-attention network for poi recommendation. *IEEE Internet of Things Journal*, 2024.

## A   Existing POI Recommendation Datasets

To examine the trend of the usage of POI recommendation datasets, we filtered the comprehensive survey by [68] to extract studies that explicitly mention the cities used in their experiments. The resulting distribution is summarized in Table 5, which shows a strong concentration of studies focused on New York and Tokyo. Additionally, Fig. 1 visualizes the same data, highlighting the uneven distribution of city choices across studies. We also include information on the LBSN platforms used, revealing that Foursquare remains the predominant data source in the field. These findings underscore the need for broader, more inclusive datasets that support evaluation across a wider range of global cities.

## B   Data Visualization

We present several visualizations highlighting Massive-STEPS' scale and diversity to complement our dataset description.

In Fig. 3, we show the top 10 most frequent POI categories for each city. The distribution reflects the local culture and lifestyle across different urban areas. For example, Beijing and Shanghai have a high number of Chinese restaurants, while Melbourne and Sydney show a strong presence of cafes. In Tokyo, convenience stores and ramen shops dominate. These patterns illustrate the diversity of local culture and user interests.

Fig. 4 plots the distribution of trajectory lengths (i.e., number of check-ins per trajectory). The distribution is long-tailed, with most trajectories being relatively short, similar to the original Semantic Trails dataset. This indicates that users often make only a few check-ins per outing.

Finally, we show the distribution of user activity levels, measured by the number of trajectories per user in Fig. 5. Most users exhibit cold-start behavior, contributing only a small number of trajectories. This highlights the importance of models that are robust to sparse and short user histories.

Table 5: **Overview of POI Recommendation Studies by City and LBSN Platform**. This table is adapted from Table IV in the survey by [68] and presents a filtered list of POI recommendation studies that explicitly mention city names and their associated LBSN platforms.

| Study | Cities | LBSN |
|---|---|---|
| SSTPMF [5] | New York, Tokyo | Foursquare, Gowalla |
| ST-LSTM [69] | California, Singapore | Brightkite, Foursquare, Gowalla |
| LSMA [43] | New York, San Fransisco, Tokyo | Foursquare, Weeplaces |
| DLAN [46] | New York, Tokyo | Foursquare |
| TLR-M [13] | New York, Tokyo | Foursquare |
| GETNext [60] | New York, Tokyo, California | Foursquare, Gowalla |
| CARAN [17] | New York, Tokyo | Foursquare, Gowalla |
| JANICP [70] | New York, Tokyo | Foursquare, Weeplaces |
| Li et al. [23] | New York, Tokyo | Foursquare |
| AMACF [58] | New York, Tokyo | Foursquare, Weeplaces |
| CHA [64] | New York, Tokyo | Foursquare |
| HAT [47] | Beijing, Shanghai | Yelp, others |
| STAR-HiT [49] | New York | Foursquare, Gowalla |
| CAFPR [14] | Tokyo, California, Budapest, Melbourne | Foursquare |
| TGAT [19] | New York, Tokyo | Foursquare |
| MobGT [50] | New York | Foursquare, Gowalla |
| POIBERT [16] | Budapest, Delhi, Edinburgh, Glasgow, Osaka, Perth, Toronto | Flickr |
| AutoMTN [30] | New York, Tokyo | Foursquare |
| CCDSA [40] | New York, Tokyo, San Fransisco | Foursquare, Weeplaces |
| TDGCN [2] | Tokyo, California | Foursquare, Gowalla, Weeplaces |
| BayMAN [48] | New York | Foursquare, Gowalla |
| ROTAN [11] | New York, Tokyo, California | Foursquare, Gowalla |
| STGCN [15] | Boston, Chicago, London | Gowalla, others |
| ADQ-GNN [44] | New York, Tokyo | Foursquare, Gowalla |
| HS-GAT [67] | Boston, Chicago, London | Yelp, others |
| HKGNN [66] | New York, Jakarta, Kuala Lumpur, Sao Paulo | Foursquare |
| S2GRec [24] | New York, Tokyo | Foursquare, Gowalla |
| GSBPL [39] | New York, Tokyo | Foursquare, Gowalla |
| LSPSL [18] | New York, Tokyo | Foursquare |
| SCL [3] | Florence, Rome, Pisa, Edinburgh, Toronto | Flickr |
| LLM-Move [10] | New York, Tokyo | Foursquare |
| LLM4POI [22] | New York, Tokyo, California | Foursquare, Gowalla |
| DiffPOI [31] | Singapore, New York, Tokyo | Foursquare, Gowalla |
| DSDRec [45] | New York, Tokyo | Foursquare |
| Diff-DGMN [72] | Istanbul, Jakarta, Sao Paulo, New York, Los Angeles | Foursquare |

Table 6: **Fields available in the Massive-STEPS dataset**, including user, POI, geographic/spatial, and temporal details, along with example data for each field.

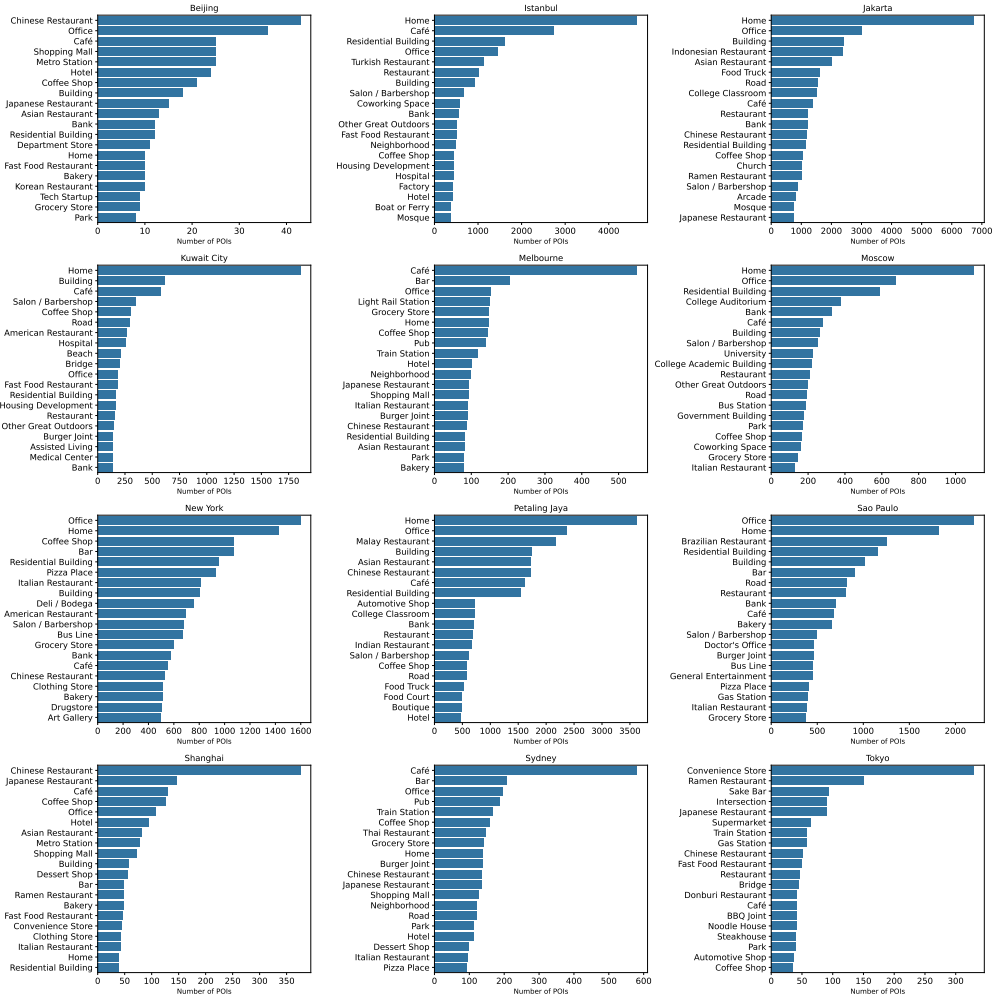| Field | Description | Example |
|---|---|---|
| trail_id | Numeric identifier of trajectory | 2013_2866 |
| user_id | Numeric identifier of user | 90 |
| venue_id | Numeric identifier of POI venue | 185 |
| latitude | Latitude of POI venue | -33.87301862604473 |
| longitude | Longitude of POI venue | 151.20668402700997 |
| name | POI name | Sydney Town Hall |
| address | Street address of POI venue | 483 George St |
| venue_category | POI category name | City Hall |
| venue_category_id | Foursquare Category ID | 4bf58dd8d48988d129941735 |
| venue_category_id_code | Numeric identifier of POI category | 72 |
| venue_city | Administrative region name | Sydney |
| venue_city_latitude | Latitude of administrative region | -33.86785 |
| venue_city_longitude | Longitude of administrative region | 151.20732 |
| venue_country | Country code | AU |
| timestamp | Check-in timestamp | 2012-04-22 08:20:00 |

Figure 3: **Top 10 most frequent POI categories in each city**, highlighting local cultural and urban preferences.
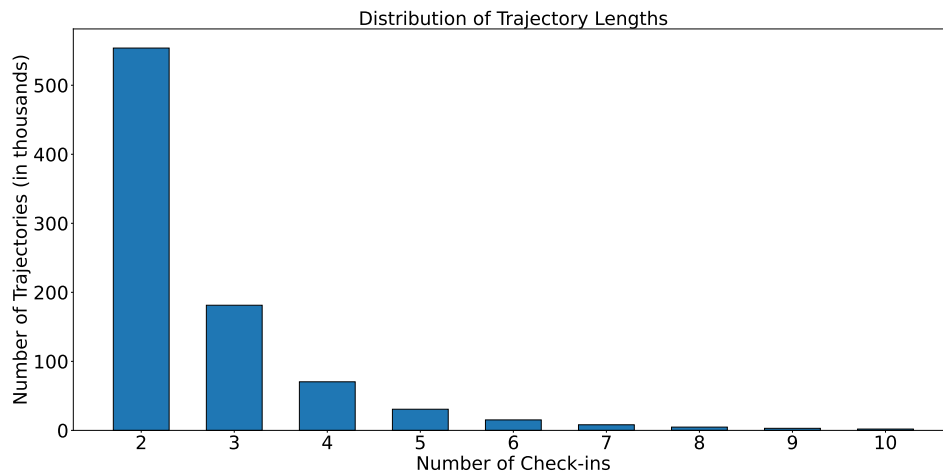


Figure 4: **Distribution of trail lengths**, showing a long-tailed pattern with most trajectories consisting of a few check-ins.
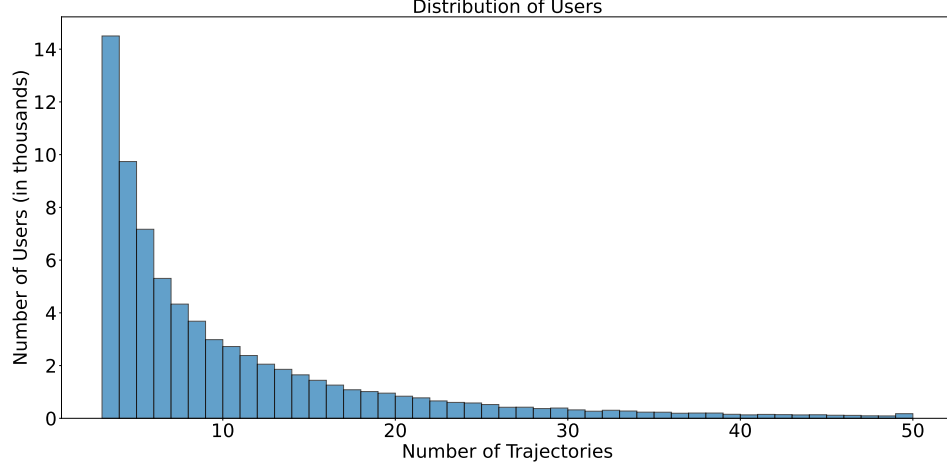
Figure 5: **Distribution of user activity** based on the number of trajectories per user, indicating a cold-start-heavy dataset.

## C  POI Recommendation: Task Details

We adopt the conventional problem formulation used in prior POI recommendation studies [68, 60, 54], which defines the task as learning user preferences and routines from historical check-ins to recommend future POIs.

### C.1  Problem Formulation

Let $\mathcal{U} = \{u_1, u_2, \ldots, u_M\}$ denote the set of users, $\mathcal{P} = \{p_1, p_2, \ldots, p_N\}$ the set of Points of Interest (POIs), and $\mathcal{T} = \{t_1, t_2, \ldots, t_K\}$ the set of timestamps, where $M, N, K \in \mathbb{N}$.

**POI Definition**    Each POI $p \in \mathcal{P}$ is represented as a tuple:

$$p = \langle \phi, \lambda, \kappa, \alpha, \beta, \gamma \rangle,$$

where:

- $\phi$ and $\lambda$ are the latitude and longitude,
- $\kappa$ is the POI category (e.g., *restaurant*, *park*),
- $\alpha$ is the unique POI identifier,
- $\beta$ is the textual address, and
- $\gamma$ is the POI name.

**Check-in Definition**    A check-in is a tuple $c = \langle u, p, t \rangle \in \mathcal{U} \times \mathcal{P} \times \mathcal{T}$, indicating that user $u$ visited POI $p$ at timestamp $t$.

**Trajectory Definition**    A trajectory for user $u$ is defined as a temporally ordered sequence of check-ins within a fixed time interval $\delta\tau = 8$ hours. Each trajectory $T_u^i(t)$ up to timestamp $t$ is defined as:

$$T_u^i(t) = \{(p_1, t_1), (p_2, t_2), \ldots, (p_k, t_k)\}$$

such that $t_1 < t_2 < \ldots < t_k = t$ and $t_k - t_{k-1} \leq \delta\tau$. Given a set of historical trajectories

$$\mathcal{T}_u = \{T_u^1, T_u^2, \ldots, T_u^L\}$$

for user $u$, where $L$ is the number of such **historical** trajectories, the goal is to recommend the POIs that $u$ is most likely to visit next after the current **contextual** trajectory $T_u'(t)$.

**POI Recommendation Task Definition**    Given a current contextual trajectory $T'_u(t)$ of user $u$ up to time $t$, along with their historical trajectories $\mathcal{T}_u$, the task of next POI recommendation is to rank all candidate POIs $p_i \in \mathcal{P}$ according to the model's predicted probability that user $u$ will visit each POI next.

Formally, the model learns a ranking function:

$$f : (T'_u(t), \mathcal{T}_u) \rightarrow \{\hat{y}_i\}_{i=1}^{|\mathcal{P}|}$$

where $\hat{y}_i$ denotes the predicted likelihood that user $u$ will visit POI $p_i$ next. Based on these scores, a ranked list of POIs is returned as recommendations.

This formulation enables POI recommendation, where the goal is to suggest a set of likely POIs that a user may visit next, based on their historical check-ins and inferred preferences. Our evaluation metrics, Acc@k and NDCG@k, assess whether the ground-truth POI appears among the top-$k$ ranked candidates, reflecting the quality of the recommended set. In particular, Acc@1 captures the stricter task of *immediate* next POI prediction, measuring whether the top-ranked POI matches the user's actual next visit.

## C.2    Experiment and Implementation Details

For training and evaluation, we used the LibCity[1] library [41], which provides implementations of classical baselines including FPMC [32], RNN [41], LSTPM [35], and DeepMove [8]. The training hyperparameters are listed in Table 7 and, unless otherwise noted, follow the default configurations provided by LibCity.

For GETNext[2] [60] and STHGCN[3] [54], we adapted the original source code released by the respective authors. Due to variations in dataset sizes and training costs across cities, we applied different hyperparameters for some cities, as detailed in Table 8.

Table 7: **Hyperparameters for Markov-based methods and recurrent networks baselines**.

| Hyperparameter | FPMC | RNN | LSTPM | DeepMove |
|---|---|---|---|---|
| Batch Size | 20 | 20 | 20 | 20 |
| Learning Rate | 5e-4 | 1e-3 | 1e-4 | 1e-3 |
| Max Epoch | 1 | 30 | 40 | 30 |
| Location Embedding Size | 64 | 500 | 500 | 500 |
| Hidden Embedding Size | N/A | 500 | 500 | 500 |
| Dropout | N/A | 0.3 | 0.8 | 0.5 |

Table 8: **Hyperparameters for Transformer-based graph neural networks**.

| Model | Cities | Batch Size | LR | Epochs |
|---|---|---|---|---|
| **GETNext** | Beijing, Melbourne, Moscow, Shanghai, Sydney, Tokyo | 16 | 1e-3 | 200 |
| | Istanbul, Kuwait City, New York, Petaling Jaya, São Paulo | 16 | 1e-4 | 20 |
| | Jakarta | 16 | 5e-5 | 20 |
| **STHGCN** | Beijing, Melbourne, Shanghai, Sydney, Tokyo | 16 | 1e-4 | 20 |
| | Istanbul, Jakarta, Kuwait City, Moscow, New York, Petaling Jaya, São Paulo | 64 | 1e-4 | 20 |

All modified code implementations are available as submodules in our main dataset repository. Experiments were conducted using GPUs provided by the School of Computer Science and Engineering at UNSW Sydney via the Wolfpack computational cluster, which includes NVIDIA L4, L40S, and H100 GPUs.

## C.3    Supplementary Results

We report the full results of our supervised POI recommendation baselines in Table 9 and 10, using three evaluation metrics: Acc@1, Acc@5, and NDCG@5.

---

[1] `https://github.com/libcity/bigscity-libcity-datasets/`

[2] `https://github.com/songyangme/GETNext`

[3] `https://github.com/alipay/Spatio-Temporal-Hypergraph-Model`

Table 9: **Performance of supervised POI recommendation baselines across 6 cities**: Beijing, Istanbul, Jakarta, Kuwait City, Melbourne, and Moscow. We report three metrics: Acc@1 (A@1), Acc@5 (A@5), and NDCG@5 (N@5).

| Model | Beijing | | | Istanbul | | | Jakarta | | | Kuwait City | | | Melbourne | | | Moscow | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 |
| FPMC | 0.000 | 0.021 | 0.009 | 0.026 | 0.074 | 0.050 | 0.029 | 0.085 | 0.058 | 0.021 | 0.089 | 0.054 | 0.062 | 0.147 | 0.107 | 0.059 | 0.129 | 0.094 |
| RNN | 0.085 | 0.183 | 0.134 | 0.077 | 0.178 | 0.130 | 0.049 | 0.115 | 0.083 | 0.087 | 0.203 | 0.146 | 0.059 | 0.105 | 0.083 | 0.075 | 0.164 | 0.122 |
| LSTPM | 0.127 | 0.211 | 0.169 | 0.142 | 0.286 | 0.217 | 0.099 | 0.210 | 0.157 | 0.180 | 0.362 | 0.275 | 0.091 | 0.204 | 0.150 | 0.151 | 0.300 | 0.229 |
| DeepMove | 0.106 | 0.261 | 0.190 | 0.150 | 0.298 | 0.228 | 0.103 | 0.212 | 0.160 | 0.179 | 0.360 | 0.274 | 0.083 | 0.179 | 0.134 | 0.143 | 0.283 | 0.217 |
| GETNext | 0.433 | 0.527 | 0.486 | 0.146 | 0.268 | 0.210 | 0.155 | 0.257 | 0.209 | 0.175 | 0.322 | 0.251 | 0.100 | 0.250 | 0.179 | 0.175 | 0.335 | 0.260 |
| STHGCN | 0.453 | 0.640 | 0.552 | 0.241 | 0.385 | 0.318 | 0.197 | 0.334 | 0.270 | 0.225 | 0.394 | 0.314 | 0.168 | 0.318 | 0.247 | 0.223 | 0.382 | 0.308 |

Table 10: **Performance of supervised POI recommendation baselines across 6 cities**: New York, Petaling Jaya, São Paulo, Shanghai, Sydney, and Tokyo. We report three metrics: Acc@1 (A@1), Acc@5 (A@5), and NDCG@5 (N@5).

| Model | New York | | | Petaling Jaya | | | São Paulo | | | Shanghai | | | Sydney | | | Tokyo | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 |
| FPMC | 0.032 | 0.090 | 0.061 | 0.026 | 0.084 | 0.057 | 0.030 | 0.079 | 0.055 | 0.084 | 0.154 | 0.120 | 0.075 | 0.180 | 0.131 | 0.176 | 0.291 | 0.239 |
| RNN | 0.061 | 0.119 | 0.092 | 0.064 | 0.148 | 0.107 | 0.097 | 0.191 | 0.147 | 0.055 | 0.120 | 0.090 | 0.080 | 0.164 | 0.125 | 0.133 | 0.254 | 0.197 |
| LSTPM | 0.099 | 0.206 | 0.155 | 0.099 | 0.222 | 0.163 | 0.158 | 0.319 | 0.243 | 0.099 | 0.195 | 0.149 | 0.141 | 0.265 | 0.206 | 0.225 | 0.394 | 0.315 |
| DeepMove | 0.097 | 0.195 | 0.149 | 0.112 | 0.234 | 0.175 | 0.160 | 0.310 | 0.240 | 0.085 | 0.168 | 0.128 | 0.129 | 0.240 | 0.188 | 0.201 | 0.362 | 0.288 |
| GETNext | 0.134 | 0.263 | 0.202 | 0.139 | 0.254 | 0.200 | 0.202 | 0.360 | 0.286 | 0.115 | 0.230 | 0.177 | 0.181 | 0.347 | 0.266 | 0.180 | 0.361 | 0.275 |
| STHGCN | 0.146 | 0.259 | 0.207 | 0.174 | 0.301 | 0.241 | 0.250 | 0.425 | 0.344 | 0.193 | 0.329 | 0.264 | 0.227 | 0.378 | 0.307 | 0.250 | 0.432 | 0.350 |

# D    Analyzing Urban Features and POI Recommendation Performance

As discussed in Section 2.2, several hypotheses have been proposed to explain why POI recommendation models perform better in certain cities than others. These hypotheses aim to uncover how various urban features affect model performance. For example, Gowalla-CA [4, 62] often yields lower accuracy compared to FSQ-NYC and FSQ-TKY [57], suggesting that some cities may be inherently harder to model. In this analysis, we focus on supervised models only.

Prior studies [60, 54, 22] have suggested several features as potential explanatory variables, including:

- Number of unique check-ins,
- Number of unique trajectories,
- Number of unique POI categories,
- Geographical area (larger areas are assumed to be harder to model), and
- POI density or spatial sparsity (i.e., unique POIs per unit area).

We also propose several additional features for consideration:

- Number of unique POIs,
- Check-in density (unique check-ins per area),
- Trajectory density (unique trajectories per area), and
- Category entropy, our proposed feature capturing category diversity.

**Category entropy**, based on Shannon entropy, measures how evenly POI categories are distributed in a city. A higher entropy suggests that check-ins are spread across a wide variety of categories, while a lower entropy indicates a concentration in fewer types. The formula for Shannon entropy is:

$$H = -\sum_{i=1}^{N} p_i \log(p_i) \tag{1}$$

where $p_i$ is the proportion of venues in category $i$, and $N$ is the total number of POI categories. The proportion $p_i$ is defined as:

$$p_i = \frac{c_i}{\sum_{j=1}^{N} c_j} \tag{2}$$

where $c_i$ is the count of venues in category $i$. In other words, $p_i$ represents the fraction of all venues that belong to category $i$.

Moreover, previous studies have primarily focused on only three datasets: FSQ-NYC, FSQ-TKY, and Gowalla-CA. In contrast, Massive-STEPS provides broader coverage across 12 cities, enabling a more comprehensive and robust analysis. To examine the relationship between urban features and model performance, we averaged the three evaluation metrics across six supervised baselines for each city and computed the Spearman correlation with each candidate feature. To further support our findings, we also included the results of GETNext [60] and STHGCN [54] on FSQ-NYC, FSQ-TKY, and Gowalla-CA, calculated their corresponding urban features, and averaged the reported metrics of each city. Fig. 6 presents the correlations between all nine features and the average performance metric.

Among all features, **category entropy** shows the strongest correlation with model performance, with a Spearman correlation of $r = -0.736$ ($p = 0.002$). This suggests that cities with more evenly distributed POI categories *tend* to yield lower prediction accuracy. Intuitively, when no single category dominates (a city has roughly equal proportions of restaurants, cafes, homes, and other POIs), it becomes more difficult for models to predict a user's next destination. In these cases, user behavior is more varied and less predictable. On the other hand, cities with more skewed category distributions (e.g., mostly food places or mostly residential areas) tend to have more consistent patterns of movement, making them easier for models to learn and predict. Interestingly, our finding contradicts the hypothesis proposed by LLM4POI [22], which suggests that FSQ-NYC is easier to model than Gowalla-CA due to the former's vast number of POI categories, which were supposed to provide richer contextual signals for the model. Our results indicate that it is not the number of categories that matters, but rather how these categories are distributed.

## E    Zero-shot POI Recommendation: Task Details

### E.1    Problem Formulation

The zero-shot POI recommendation task follows the same problem formulation as its supervised counterpart (see Section C.1). The key difference is that in this setting, the model parameters remain frozen and the models are pre-trained, rather than trained from randomly initialized weights.

### E.2    Experiment and Implementation Details

**Preprocessing**    We adopted the AgentMove[4] library [7], which provides implementations of three LLM methods: LLM-Mob [42], LLM-ZS [1], and LLM-Move [10]. The preprocessing steps used by AgentMove are as follows.

First, we selected 200 random users from the test set and sampled one random trajectory for each user. This trajectory serves as the **context stays**, representing the current trajectory to be predicted. The **historical stays** are composed of the most recent 15 trajectories from the same user, drawn from the training set. Each check-in is described by four attributes: the hour (in 12-hour format), the day of the week, the POI ID, and the POI category name.

Second, the LLMs are set to return outputs in JSON format, generating the top 5 predicted POI IDs along with an explanation of their reasoning. Following the AgentMove setup and to ensure replicability, we set the generation parameters as follows: a temperature of 0.0, a maximum output length of 1000 tokens, and an input context window capped at 2000 tokens.

**Prompting**    Prompt templates for each method, LLM-Mob, LLM-ZS, and LLM-Move, are presented in Listing 1, 2, and 3, respectively.

```
1  Your task is to predict a user's next location based on his/her
       activity pattern.
```

---

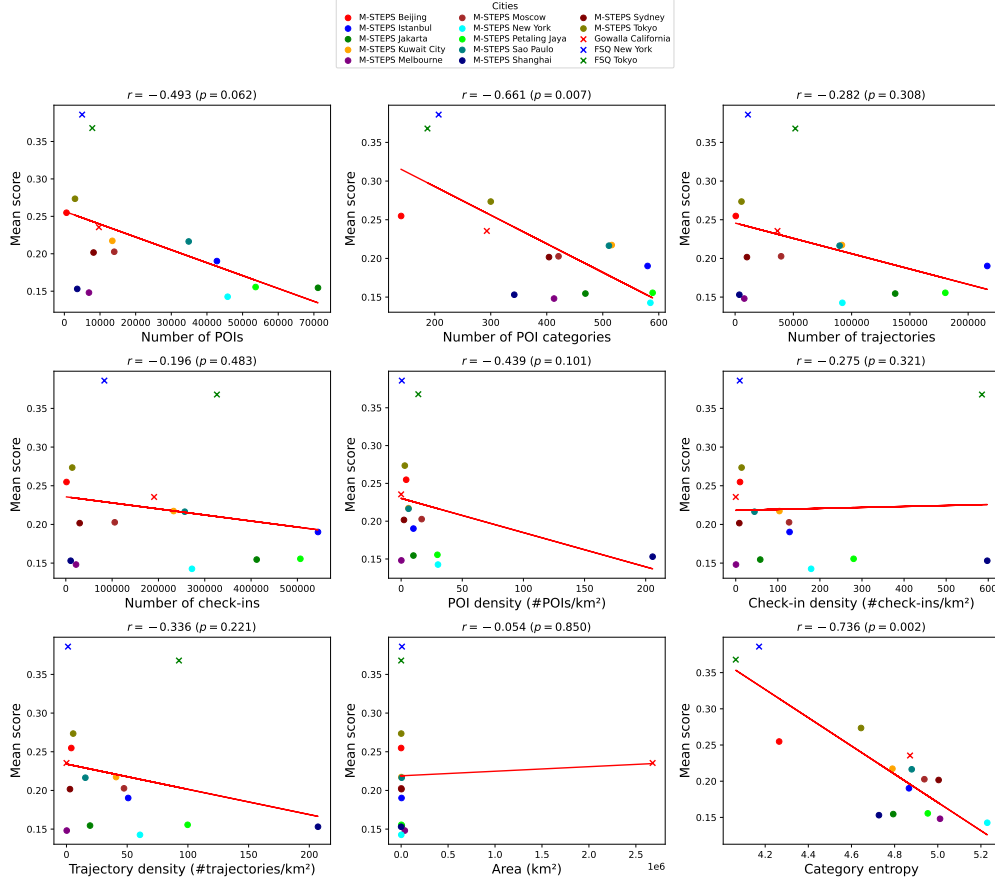[4]https://github.com/tsinghua-fib-lab/agentmove/

Figure 6: **Spearman correlation between nine candidate urban features and the mean score of POI recommendation models** across 15 cities, including Massive-STEPS (ours), FSQ [57], and Gowalla [4, 62].

```
2  You will be provided with <history> which is a list containing this
       user's historical stays, then <context> which provide contextual
       information
3  about where and when this user has been to recently. Stays in both <
       history> and <context> are in chronological order.
4  Each stay takes on such form as (start_time, day_of_week, duration,
       place_id). The detailed explanation of each element is as follows:
5  start_time: the start time of the stay in 12h clock format.
6  day_of_week: indicating the day of the week.
7  duration: an integer indicating the duration (in minute) of each stay.
        Note that this will be None in the <target_stay> introduced later
       .
8  place_id: an integer representing the unique place ID, which indicates
        where the stay is.
9
10 Then you need to do next location prediction on <target_stay> which is
        the prediction target with unknown place ID denoted as <
       next_place_id> and
11 unknown duration denoted as None, while temporal information is
       provided.
12
13 Please infer what the <next_place_id> might be (please output the 10
       most likely places which are ranked in descending order in terms
       of probability), considering the following aspects:
```

```
14 1. the activity pattern of this user that you learned from <history>,
      e.g., repeated visits to certain places during certain times;
15 2. the context stays in <context>, which provide more recent
      activities of this user;
16 3. the temporal information (i.e., start_time and day_of_week) of
      target stay, which is important because people's activity varies
      during different time (e.g., nighttime versus daytime)
17 and on different days (e.g., weekday versus weekend).
18
19 Please organize your answer in a JSON object containing following keys
      :
20 "prediction" (the ID of the five most probable places in descending
      order of probability) and "reason" (a concise explanation that
      supports your prediction). Do not include line breaks in your
      output.
21
22 The data are as follows:
23 <history>: {historical_stays}
24 <context>: {context_stays}
25 <target_stay>: {target_time, target_day_of_week}
```

Listing 1: Prompt for LLM-Mob

```
1 Your task is to predict <next_place_id> in <target_stay>, a location
      with an unknown ID, while temporal data is available.
2
3 Predict <next_place_id> by considering:
4 1. The user's activity trends gleaned from <historical_stays> and the
      current activities from  <context_stays>.
5 2. Temporal details (start_time and day_of_week) of the target stay,
      crucial for understanding activity variations.
6
7 Present your answer in a JSON object with:
8 "prediction" (IDs of the five most probable places, ranked by
      probability) and "reason" (a concise justification for your
      prediction).
9
10 The data:
11 <historical_stays>: {historical_stays}
12 <context_stays>: {context_stays}
13 <target_stay>: {target_time, target_day_of_week}
```

Listing 2: Prompt for LLM-ZS

```
1 <long-term check-ins> [Format: (POIID, Category)]: {historical_stays}
2 <recent check-ins> [Format: (POIID, Category)]: {context_stays}
3 <candidate set> [Format: (POIID, Distance, Category)]: {candidates}
4 Your task is to recommend a user's next point-of-interest (POI) from <
      candidate set> based on his/her trajectory information.
5 The trajectory information is made of a sequence of the user's <long-
      term check-ins> and a sequence of the user's <recent check-ins> in
       chronological order.
6 Now I explain the elements in the format. "POIID" refers to the unique
       id of the POI, "Distance" indicates the distance (kilometers)
      between the user and the POI, and "Category" shows the semantic
      information of the POI.
7
8 Requirements:
9 1. Consider the long-term check-ins to extract users' long-term
      preferences since people tend to revisit their frequent visits.
10 2. Consider the recent check-ins to extract users' current perferences
      .
11 3. Consider the "Distance" since people tend to visit nearby pois.
12 4. Consider which "Category" the user would go next for long-term
      check-ins indicates sequential transitions the user prefer.
```

```
13
14  Please organize your answer in a JSON object containing following keys
        :
15  "prediction" (10 distinct POIIDs of the ten most probable places in <
        candidate set> in descending order of probability), and "reason" (
        a concise explanation that supports your recommendation according
        to the requirements). Do not include line breaks in your output.
```
Listing 3: Prompt for LLM-Move

**Models and Implementations**  We use the following LLMs in our experiments:

- Gemini 2.0 Flash (`gemini-2.0-flash`),
- Qwen 2.5 7B Instruct (`Qwen2.5-7B-Instruct-AWQ`)[5],
- Llama 3.1 8B Instruct (`Meta-Llama-3.1-8B-Instruct-AWQ-INT4`)[6],
- Gemma 2 9B Instruct (`gemma-2-9b-it-AWQ-INT4`)[7].

All open-source models are quantized using AWQ [25] and served via vLLM [21]. Inference of open-source models was conducted on NVIDIA A100 GPUs provided by SaturnCloud[8] through the NVIDIA Academic Grant Program. We accessed Gemini via the official API and through the Google Cloud Research Credits program. All modified code implementations are publicly available in our main dataset repository.

### E.3  Supplementary Results

We provide the full results of our zero-shot POI recommendation results in Table 11 and 12, providing three metrics: Acc@1, Acc@5, and NDCG@5.

Table 11: **Performance of zero-shot POI recommendation baselines across 6 cities**: Beijing, Istanbul, Jakarta, Kuwait City, Melbourne, and Moscow. We report three metrics: Acc@1 (A@1), Acc@5 (A@5), and NDCG@5 (N@5).

| Method | Model | Beijing | | | Istanbul | | | Jakarta | | | Kuwait City | | | Melbourne | | | Moscow | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 |
| LLM-Mob | Gemini 2 Flash | 0.115 | 0.308 | 0.226 | 0.080 | 0.225 | 0.160 | 0.100 | 0.245 | 0.174 | 0.095 | 0.270 | 0.185 | 0.060 | 0.150 | 0.108 | 0.130 | 0.245 | 0.187 |
| | Qwen 2.5 7B | 0.058 | 0.385 | 0.218 | 0.035 | 0.240 | 0.148 | 0.105 | 0.245 | 0.179 | 0.080 | 0.220 | 0.155 | 0.030 | 0.130 | 0.083 | 0.090 | 0.270 | 0.185 |
| | Llama 3.1 8B | 0.000 | 0.000 | 0.000 | 0.020 | 0.110 | 0.065 | 0.055 | 0.150 | 0.104 | 0.030 | 0.100 | 0.066 | 0.010 | 0.065 | 0.040 | 0.030 | 0.100 | 0.068 |
| | Gemma 2 9B | 0.115 | 0.288 | 0.206 | 0.075 | 0.200 | 0.146 | 0.105 | 0.240 | 0.178 | 0.080 | 0.210 | 0.150 | 0.055 | 0.150 | 0.108 | 0.100 | 0.240 | 0.176 |
| LLM-ZS | Gemini 2 Flash | 0.058 | 0.385 | 0.246 | 0.090 | 0.235 | 0.166 | 0.110 | 0.250 | 0.188 | 0.080 | 0.245 | 0.167 | 0.065 | 0.160 | 0.115 | 0.125 | 0.300 | 0.217 |
| | Qwen 2.5 7B | 0.038 | 0.404 | 0.237 | 0.040 | 0.235 | 0.141 | 0.065 | 0.250 | 0.161 | 0.050 | 0.220 | 0.140 | 0.040 | 0.155 | 0.100 | 0.080 | 0.260 | 0.176 |
| | Llama 3.1 8B | 0.077 | 0.346 | 0.221 | 0.040 | 0.225 | 0.137 | 0.045 | 0.200 | 0.126 | 0.060 | 0.210 | 0.137 | 0.040 | 0.155 | 0.101 | 0.080 | 0.270 | 0.183 |
| | Gemma 2 9B | 0.096 | 0.308 | 0.217 | 0.045 | 0.225 | 0.141 | 0.105 | 0.250 | 0.180 | 0.070 | 0.230 | 0.153 | 0.050 | 0.140 | 0.100 | 0.080 | 0.290 | 0.194 |
| LLM-Move | Gemini 2 Flash | 0.096 | 0.346 | 0.218 | 0.205 | 0.385 | 0.289 | 0.295 | 0.405 | 0.350 | 0.220 | 0.380 | 0.295 | 0.225 | 0.325 | 0.275 | 0.220 | 0.400 | 0.316 |
| | Qwen 2.5 7B | 0.192 | 0.346 | 0.280 | 0.175 | 0.270 | 0.226 | 0.115 | 0.225 | 0.169 | 0.160 | 0.285 | 0.227 | 0.110 | 0.220 | 0.165 | 0.230 | 0.310 | 0.274 |
| | Llama 3.1 8B | 0.058 | 0.135 | 0.100 | 0.015 | 0.055 | 0.036 | 0.015 | 0.025 | 0.021 | 0.010 | 0.035 | 0.023 | 0.040 | 0.195 | 0.123 | 0.005 | 0.065 | 0.031 |
| | Gemma 2 9B | 0.096 | 0.365 | 0.229 | 0.100 | 0.200 | 0.155 | 0.235 | 0.290 | 0.266 | 0.120 | 0.275 | 0.202 | 0.115 | 0.275 | 0.199 | 0.110 | 0.245 | 0.185 |

## F  License and Data Usage

Our work **does not** involve the collection of new data. Instead, we derive our resulting dataset by combining and aligning two publicly available datasets, both of which are distributed under permissive licenses. We did not scrape data from the internet or use proprietary APIs to construct this dataset.

We accessed the Semantic Trails Dataset [29] via Figshare: `https://doi.org/10.6084/m9.figshare.7429076.v2`. The dataset is licensed under the Creative Commons CC0 1.0 license

---

[5] `https://huggingface.co/qwen/qwen2.5-7b-instruct-awq`
[6] `https://huggingface.co/hugging-quants/Meta-Llama-3.1-8B-Instruct-AWQ-INT4`
[7] `https://huggingface.co/hugging-quants/gemma-2-9b-it-AWQ-INT4`
[8] `https://saturncloud.io/`

Table 12: **Performance of zero-shot POI recommendation baselines across 6 cities**: New York, Petaling Jaya, São Paulo, Shanghai, Sydney, and Tokyo. We report three metrics: Acc@1 (A@1), Acc@5 (A@5), and NDCG@5 (N@5).

| Method | Model | New York | | | Petaling Jaya | | | São Paulo | | | Shanghai | | | Sydney | | | Tokyo | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 | A@1 | A@5 | N@5 |
| LLM-Mob | Gemini 2 Flash | 0.095 | 0.175 | 0.136 | 0.090 | 0.220 | 0.160 | 0.130 | 0.305 | 0.223 | 0.055 | 0.160 | 0.111 | 0.060 | 0.160 | 0.112 | 0.140 | 0.320 | 0.238 |
| | Qwen 2.5 7B | 0.070 | 0.185 | 0.131 | 0.030 | 0.195 | 0.116 | 0.090 | 0.290 | 0.188 | 0.040 | 0.170 | 0.108 | 0.035 | 0.145 | 0.091 | 0.110 | 0.350 | 0.243 |
| | Llama 3.1 8B | 0.025 | 0.090 | 0.061 | 0.010 | 0.090 | 0.050 | 0.030 | 0.165 | 0.098 | 0.005 | 0.020 | 0.013 | 0.020 | 0.085 | 0.053 | 0.005 | 0.045 | 0.025 |
| | Gemma 2 9B | 0.070 | 0.175 | 0.124 | 0.055 | 0.185 | 0.122 | 0.085 | 0.230 | 0.162 | 0.050 | 0.150 | 0.104 | 0.030 | 0.130 | 0.086 | 0.145 | 0.345 | 0.255 |
| LLM-ZS | Gemini 2 Flash | 0.080 | 0.170 | 0.129 | 0.110 | 0.210 | 0.164 | 0.150 | 0.315 | 0.235 | 0.065 | 0.160 | 0.113 | 0.060 | 0.155 | 0.111 | 0.160 | 0.380 | 0.278 |
| | Qwen 2.5 7B | 0.050 | 0.180 | 0.116 | 0.045 | 0.175 | 0.111 | 0.095 | 0.290 | 0.198 | 0.045 | 0.155 | 0.103 | 0.045 | 0.170 | 0.109 | 0.110 | 0.365 | 0.257 |
| | Llama 3.1 8B | 0.055 | 0.160 | 0.111 | 0.030 | 0.205 | 0.123 | 0.030 | 0.280 | 0.159 | 0.060 | 0.165 | 0.116 | 0.040 | 0.185 | 0.110 | 0.110 | 0.415 | 0.269 |
| | Gemma 2 9B | 0.075 | 0.185 | 0.129 | 0.065 | 0.185 | 0.126 | 0.075 | 0.300 | 0.192 | 0.050 | 0.165 | 0.112 | 0.045 | 0.155 | 0.103 | 0.110 | 0.395 | 0.263 |
| LLM-Move | Gemini 2 Flash | 0.235 | 0.415 | 0.325 | 0.210 | 0.335 | 0.273 | 0.285 | 0.415 | 0.350 | 0.170 | 0.270 | 0.221 | 0.230 | 0.420 | 0.331 | 0.250 | 0.470 | 0.368 |
| | Qwen 2.5 7B | 0.120 | 0.255 | 0.188 | 0.135 | 0.175 | 0.155 | 0.155 | 0.235 | 0.199 | 0.095 | 0.165 | 0.133 | 0.125 | 0.280 | 0.205 | 0.250 | 0.360 | 0.312 |
| | Llama 3.1 8B | 0.035 | 0.130 | 0.084 | 0.040 | 0.060 | 0.049 | 0.045 | 0.045 | 0.045 | 0.020 | 0.040 | 0.030 | 0.055 | 0.220 | 0.141 | 0.030 | 0.060 | 0.046 |
| | Gemma 2 9B | 0.115 | 0.245 | 0.183 | 0.175 | 0.235 | 0.208 | 0.195 | 0.300 | 0.252 | 0.105 | 0.150 | 0.128 | 0.125 | 0.370 | 0.254 | 0.130 | 0.305 | 0.225 |

We accessed the Foursquare Open Source Places dataset via Hugging Face: `https://huggingface.co/datasets/foursquare/fsq-os-places`. Foursquare Open Source Places is licensed under the Apache License, Version 2.0:

More details are available in Foursquare's documentation: `https://docs.foursquare.com/data-products/docs/access-fsq-os-places`.

We released our Massive-STEPS dataset under the same Apache Version 2.0 License, and have included Foursquare Open Source Places' license in our hosted dataset's README file.

## NeurIPS Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: The claims and contributions in the abstract and introduction have been clearly stated and accurately match the experimental results.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: The paper discussed the limitations of the work in the dedicated "Limitations" section (see Section 5), which discussed the limitations of our dataset and benchmarks.

   Guidelines:

   - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
   - The authors are encouraged to create a separate "Limitations" section in their paper.
   - The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
   - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
   - The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
   - The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
   - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
   - While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory assumptions and proofs**

   Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental result reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We ensured reproducibility by (1) releasing the dataset and all accompanying code to replicate our supervised and zero-shot POI recommendation experiments, and (2) providing detailed descriptions of preprocessing procedures, model configurations, training setups, and evaluation protocols throughout the paper (see Section C.2 and Section E.2).

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general. releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

    Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

    Answer: [Yes]

    Justification: The paper provides open access to the Massive-STEPS dataset and all source code, along with detailed instructions for data access, preprocessing, and experiment reproduction. This includes scripts for preparing raw and resultant aligned data, reproducing all supervised and zero-shot POI recommendation benchmarks, and evaluating both proposed methods and baselines.

    Guidelines:

    - The answer NA means that paper does not include experiments requiring code.
    - Please see the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
    - While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
    - The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
    - The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
    - The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
    - At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
    - Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental setting/details**

    Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

    Answer: [Yes]

    Justification: The paper specifies all relevant training and testing details (see Section C and E). For supervised models, we describe how hyperparameters were selected (following original papers), and for zero-shot LLM baselines, we detail the prompt formats, generation parameters, and input constraints.

    Guidelines:

    - The answer NA means that the paper does not include experiments.
    - The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
    - The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment statistical significance**

    Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

    Answer: [No]

    Justification: We did not conduct repeated experiments to calculate error bars in our benchmarking due to computational resource constraints. However, we included the p-value of our proposed hypothesis in Section D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments compute resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We stated the computational resources used and needed to reproduce our experiments in Section C.2 and Section E.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code of ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics `https://neurips.cc/public/EthicsGuidelines`?

Answer: [Yes]

Justification: We have reviewed and ensured that our research conforms to the NeurIPS Code of Ethics. The original datasets we used and on which we based our dataset have been anonymized to remove PII and have a clear license (see Section F).

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We have discussed the positive societal benefits of POI recommendation research in Section 1 and the issues of the field, such as the existence of geographical bias in LLMs (see Section 2.2), which may cause fairness issues.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: While we did not scrape a dataset from the internet and instead, based our dataset on two existing datasets, the latter have been anonymized to remove PII and do not impose a high risk for misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have cited the original paper that produced the Semantic Trails dataset and have clearly mentioned the license of Semantic Trails (CC0 1.0 Universal) and Foursquare Open Source Places (Apache 2.0), providing the URL. See Section F.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

    Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

    Answer: [Yes]

    Justification: We released the dataset and all accompanying code with documentation to replicate our experiments.

    Guidelines:

    - The answer NA means that the paper does not release new assets.
    - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
    - The paper should discuss whether and how consent was obtained from people whose asset is used.
    - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

    Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

    Answer: [NA]

    Justification: The paper does not involve crowdsourcing nor research with human subjects.

    Guidelines:

    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
    - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
    - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components, and was only used for writing, editing, or formatting purposes.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.