

Learning Fine-Grained Urban Mobility Dynamics through Large Model-Enhanced Multimodal Representations

Tianhong Zhao^a, Jianbin Li^a, Jinzhou Cao^{a,*}, Wei Tu^b, Filip Biljecki^c,
Shengao Yi^d, Zhilu Yuan^b

^a*Spatio-temporal Intelligence Center & School of Artificial Intelligence, Shenzhen
Technology University, Shenzhen 518118, China*

^b*School of Architecture and Urban Planning, Shenzhen University, Shenzhen 518060,
China*

^c*Department of Architecture, National University of Singapore, Singapore 117356*

^d*Department of City and Regional Planning, University of Pennsylvania, Philadelphia,
PA 19104 USA*

Abstract

Accurately predicting fine-grained urban mobility is essential for optimizing transportation, accessibility, and urban management. However, existing approaches often depend on dynamic data such as trajectories or signaling records, which are sparsely available across cities, thereby limiting their applicability and generalizability to new urban contexts. To address these limitations, this study proposes a Large Model Enhanced Multimodal Representations (LMEMR) framework to learn hourly grid-level mobility dynamics solely from static geospatial data—including remote sensing imagery, building data, street view imagery, and points of interest—which are widely accessible. Large vision–language models are employed to generate natural-language descriptions of each modality, enriching the data with human-understandable semantics. A dual-level contrastive learning strategy aligns raw and textual features both within and across modalities, mitigating semantic gaps and enhancing multimodal consistency. Spatial dependencies are modeled through a graph attention network, and temporal dynamics are captured via a transformer encoder to produce 24-hour mobility sequences.

*Corresponding author

Email address: caojinzhou@sztu.edu.cn (Jinzhou Cao)

Results from Shenzhen demonstrate that LMEMR outperforms the baseline CLIP model, achieving an R^2 of 0.856 and an 18.07% reduction in MAE. Ablation experiments confirm the effectiveness of semantic enhancement, spatial graph reasoning, and cross-modal fusion. Overall, this research reveals the potential of static multimodal data for dynamic mobility inference, offering a scalable, interpretable, and privacy-friendly solution for smart city planning and management.

Keywords: Urban mobility prediction, Multimodal learning, Vision–language models, Contrastive learning, Graph attention networks

1. Introduction

Cities are inherently dynamic systems shaped by continuous human movement, interaction, and adaptation. This dynamic nature is most directly reflected in urban mobility patterns, which describe how people travel through the city throughout the day. In this study, urban mobility patterns are defined as the temporal sequence of outbound flows within grids over a 24-hour period [1, 2]. Understanding and predicting these patterns is fundamental for effective transportation management, equitable allocation of public services, and resilient urban planning [3]. Accurate modeling of urban mobility provides essential information for optimizing transit operations, mitigating congestion, and designing cities that can adapt to daily fluctuations and long-term changes in human activity [4].

Despite extensive research on urban activity modeling, most existing approaches remain limited to activity prediction or functional classification, without systematically characterizing the fine-grained temporal evolution of urban mobility [5, 6, 7]. Large-scale geographic datasets commonly available for urban analysis—such as remote sensing imagery (RSI), street view imagery (SVI), building data, and points of interest (POI)—are essentially static. In contrast to trajectory or signaling data that track real-time urban mobility dynamics, static spatial data provide distinct benefits in terms of wide availability, spatial coverage, and cross-temporal transferability [8]. These datasets record the physical form and functional attributes of the built environment, but do not directly capture the temporal dynamics of residents’ travel demand [9]. To bridge this gap, this study aims to develop a framework that integrates multi-source static geographic data—which is widely accessible across cities—to accurately infer hourly, grid-level outbound flows.

27 Previous research has shown that static geographic data have the poten-
28 tial to reveal certain dynamic patterns of human activity. For example, POI
29 data, which reflect the spatial distribution of urban functional areas, have
30 been widely used to predict human mobility and regional activity levels [10].
31 Jiang et al. proposed a transfer learning method based on POI embeddings
32 and deep learning, which achieved significant improvements in the prediction
33 of human mobility in cities [11]. Meanwhile, RSI provides objective infor-
34 mation on the physical environment, such as building density, green space,
35 and road networks. These features have the potential to infer human ac-
36 tivity patterns from urban morphology [12, 13]. Complementarily, SVI pro-
37 vide finer-grained representations of the city from a human perspective [14].
38 SVI captures details such as amenities, building façades, and cultural styles,
39 revealing urban functions and reflecting mobility patterns on streets with
40 similar functions [15, 16]. Zhang et al. showed that combining high-level
41 semantic features from SVI via deep convolutional neural networks with taxi
42 mobility data can effectively forecast hourly taxi demand [17].

43 Nevertheless, three major challenges remain to be addressed. *First, se-*
44 *mantic representation is insufficient.* Previous methods mainly depend on
45 low-level features, such as pixel values, lacking high-level semantic under-
46 standing of regional functions and behavioral implications. For example,
47 high-density built-up areas can correspond to residential neighborhoods or in-
48 dustrial parks, exhibiting drastically different mobility patterns—distinctions
49 that raw visual or structural features cannot capture [18]. *Second, fusing het-*
50 *erogeneous modalities remains challenging.* RSI, SVI, building, and POI data
51 differ fundamentally in form and semantics. A naive concatenation method
52 often introduces noise and semantic shifts, affecting the robustness and gener-
53 alization of the model [19, 20]. *Third, spatial neighborhood dependencies are*
54 *often underexplored.* Urban mobility patterns are influenced not only by the
55 intrinsic attributes of a grid but also by the functions and transport facilities
56 of its surrounding context. For example, residential grids adjacent to metro
57 stations experience substantial demand surges during morning peaks—yet
58 such non-local dependencies are difficult to capture without explicit spatial
59 graph structures [21].

60 To address these challenges, we propose the *Large Model Enhanced Mul-*
61 *timodal Representations (LMEMR)* framework, which transforms static mul-
62 timodal geospatial data into dynamic representations of human mobility.
63 LMEMR consists of three key modules that correspond to the identified
64 challenges: (i) To enhance semantic representation, a *multimodal semantic*

65 *enhancement* module leverages large vision–language models (VLMs/LLMs)
66 to enrich the high-level semantics of RSI, building, SVI, and POI data; (*ii*)
67 To improve multimodal fusion, a *cross-modal fusion* module aligns modality-
68 specific embeddings through contrastive learning and attention-based inte-
69 gration for hourly mobility prediction; and (*iii*) To capture spatial dependen-
70 cies, a *spatial-context representation learning* module employs a graph atten-
71 tion network (GAT) to model local and non-local neighborhood interactions.
72 By jointly encoding semantic, spatial and multimodal features, LMEMR en-
73 ables an interpretable and accurate inference of fine-grained urban mobility
74 patterns.

75 2. Related Work

76 2.1. Urban Mobility Pattern Prediction

77 Driven by urban big data, urban mobility prediction has evolved from ag-
78 gregate estimation to fine-grained spatiotemporal forecasting [22, 23]. Exist-
79 ing research generally falls into three categories: short-term demand, origin-
80 destination (OD) flow, and mobility pattern prediction. Short-term demand
81 prediction focuses on immediate variations, where CNN–LSTM architectures
82 have proven effective in capturing nonlinear dependencies for applications
83 like metro and ride-hailing forecasting [24, 25]. OD flow prediction models
84 dynamic interzonal interactions. While traditional gravity models laid the
85 foundation, recent graph-based and attention-driven networks have signifi-
86 cantly improved large-scale estimation [26, 27].

87 Mobility pattern prediction aims to reconstruct fine-grained temporal
88 rhythms. Arenas et al. identified distinct patterns for functional zones us-
89 ing clustering [28], while recent works employ Transformer and GCNs to
90 model spatiotemporal dependencies [29]. To explain behavioral differences,
91 semantic data (POIs, SVI) have been integrated via knowledge graphs [30] or
92 geographic-semantic GCNs [31]. Recent advances also explore LLM-powered
93 semantic synthesis for dynamic mobility prediction [32] and deep learning
94 approaches for predicting mobility flows directly from satellite imagery [33].
95 This research line shifts focus from aggregate demand to understanding struc-
96 tural mobility regularities linking the built environment with human dynam-
97 ics; however, most existing methods still rely on dynamic data, leaving static-
98 only mobility inference largely unexplored.

99 *2.2. Multimodal Fusion for Geospatial Data*

100 Multimodal fusion of geospatial data (RSI, Building, SVI, POI) has be-
101 come central to urban computing [34, 35]. Research has shifted from single-
102 modal classification [36] to integrating multiple sources for tasks like func-
103 tional zone identification [37]. Recent advances in GNNs and contrastive
104 learning facilitate joint modeling of heterogeneous data, enhancing regional
105 representations [38] and mitigating label scarcity [39]. Current progress fo-
106 cuses on cross-modal semantic alignment [40], multi-level spatial graph mod-
107 eling [41], and contrastive strategies for consistency [42].

108 Existing fusion methods are generally categorized as feature-based, alignment-
109 based, or contrast-based. Feature-based concatenation often lacks semantic
110 alignment and generalization [43]. Alignment-based methods map modalities
111 to shared spaces but can be sensitive to noise [44]. Contrast-based frame-
112 works explicitly enforce cross-modal consistency, offering robust generaliza-
113 tion across heterogeneous data [45, 46]. Recent stochastic multimodal fusion
114 approaches [47] and multi-scale contrastive learning frameworks [48] further
115 enhance robust spatial representations across diverse urban tasks. These
116 efforts mark a shift toward semantically aligned and structurally informed
117 integration, yet few adopt a dual-level (intra- and inter-modal) contrastive
118 strategy.

119 *2.3. Semantic Representation of Geospatial Data*

120 Despite data growth, urban models struggle with the semantic gap—
121 the disconnect between low-level features and high-level human concepts [49,
122 50]. Traditional representations (e.g., spectral indices) lack behavioral mean-
123 ing [51]. While POI embeddings [52] and land-use classification [53] offer
124 partial solutions, they are limited by fixed taxonomies and lack contextual
125 flexibility. Consequently, models often describe where activities occur with-
126 out explaining why.

127 To bridge this gap, recent studies integrate natural-language priors. Large
128 multimodal models (e.g., CLIP, Qwen-VL) enable vision–language align-
129 ment, reshaping spatial representation [54, 55, 56]. This “GeoAI + LLM”
130 trend facilitates tasks like image annotation and knowledge extraction [57].
131 Benchmarks for evaluating VLMs on urban scene perception [58] reveal that
132 models achieve stronger alignment on objective properties than subjective
133 appraisals. Contrastive learning further aligns embeddings in a unified se-
134 mantic space [59]. Building on these developments, our work employs VLMs

135 to generate textual descriptions from geospatial inputs and aligns them with
 136 spatial features via contrastive learning for mobility modeling.

137 **3. Methodology**

138 *3.1. Problem Formulation and Notations*

139 This study addresses the task of inferring fine-grained urban mobility
 140 dynamics from static geospatial data. Let $\mathcal{G} = \{g_1, g_2, \dots, g_N\}$ denote a set
 141 of N non-overlapping urban grids partitioning the study area, where each
 142 grid corresponds to a spatial unit of size $500 \text{ m} \times 500 \text{ m}$. For each grid g_i , the
 143 input comprises multimodal static geospatial data from four modalities:

$$\mathcal{X}_i = \left\{ \mathbf{x}_i^{\text{rsi}}, \mathbf{x}_i^{\text{bld}}, \mathbf{x}_i^{\text{svi}}, \mathbf{x}_i^{\text{poi}} \right\}, \quad (1)$$

144 where $\mathbf{x}_i^{\text{rsi}}$, $\mathbf{x}_i^{\text{bld}}$, $\mathbf{x}_i^{\text{svi}}$, and $\mathbf{x}_i^{\text{poi}}$ denote remote sensing imagery, building height
 145 maps, street view images, and POI feature vectors, respectively. The model
 146 predicts a 24-hour mobility sequence $\hat{\mathbf{y}}_i = [\hat{y}_{i,1}, \hat{y}_{i,2}, \dots, \hat{y}_{i,24}] \in \mathbb{R}^{24}$ for each
 147 grid, where $\hat{y}_{i,t}$ denotes the predicted outbound flow volume at hour t . Given
 148 a training set $\{(\mathcal{X}_i, \mathbf{y}_i)\}_{i=1}^N$, the model learns a mapping function $f_\theta : \mathcal{X} \rightarrow \hat{\mathbf{y}}$
 149 by minimizing the prediction error:

$$\mathcal{L} = \sum_{i=1}^N \sum_{t=1}^{24} \ell(\hat{y}_{i,t}, y_{i,t}), \quad (2)$$

150 where $\ell(\cdot, \cdot)$ denotes a suitable loss function.

151 Table 1 summarizes the key notations used throughout the methodology.

152 *3.2. Overview of LMEMR*

153 The proposed LMEMR establishes a framework that transforms static
 154 multimodal geospatial data into dynamic representations of human mobil-
 155 ity, as shown in Fig. 1. The LMEMR comprises three key components.
 156 (i) The *Multimodal Semantic Enhancement* module uses VLMs/LLMs to
 157 enhance the semantic representation of RSI, building, SVI, and POI data.
 158 (ii) The *Graph-Enhanced Spatial-Context Learning* module models spatial
 159 relationships through GATs to capture both local and non-local neighbor-
 160 hood dependencies. (iii) The *Cross-Modal Fusion and Temporal Modeling*
 161 module integrates multimodal features via attention-based fusion and pre-
 162 dicted 24-hour grid-level mobility sequences using a Transformer encoder. By

Table 1: Summary of main notations.

Symbol	Definition
$\mathcal{M} = \{\text{rsi, bld, svi, poi}\}$	Set of data modalities.
\mathbf{x}_i^m, s_i^m	Raw input and textual description for grid i and modality m .
$E_v^m(\cdot), E_t^m(\cdot)$	Spatial and semantic encoders for modality m .
$\mathbf{H}^m, \tilde{\mathbf{H}}^m$	Intra-modally aligned and graph-enhanced feature matrices for modality m .
\mathbf{z}_i	Cross-modally fused embedding for grid i .
$\hat{y}_{i,t}$	Predicted outbound flow for grid i at hour t .
$\mathcal{L}_{\text{cla}}, \mathcal{L}_{\text{reg}}, \mathcal{L}_{\text{CL}}$	Classification, regression, and contrastive losses.

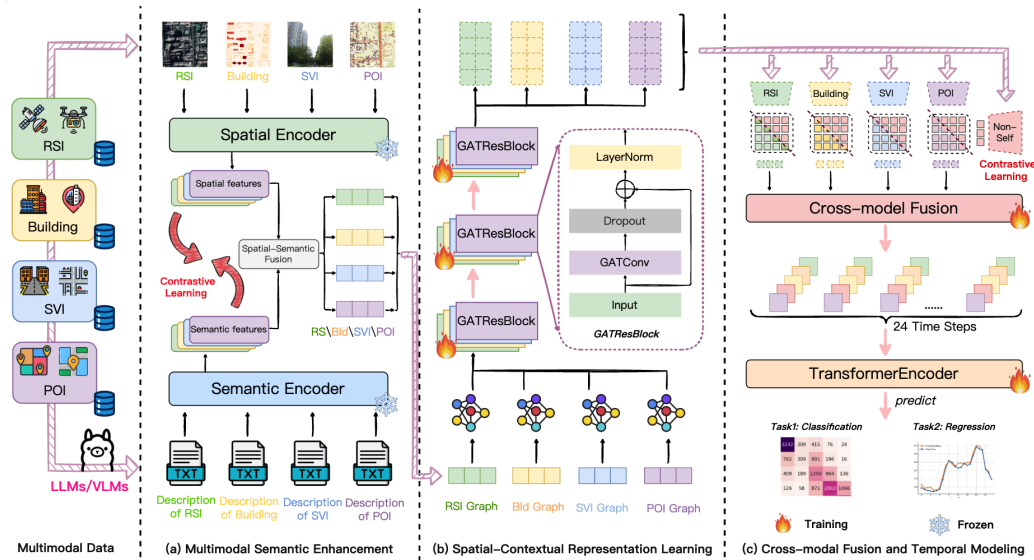


Figure 1: Overview of the proposed LMEMR. The framework consists of three modules: (a) multimodal semantic enhancement, (b) spatial-contextual representation learning, and (c) cross-modal fusion with temporal modeling.

163 combining semantic alignment, graph-based spatial reasoning, and temporal
 164 sequence modeling, LMEMR enables interpretable high-fidelity inference of
 165 fine-grained urban mobility dynamics from static multimodal data.

166 *3.3. Multimodal Semantic Enhancement*

167 This module enriches the semantic expressiveness of static geospatial data
168 by aligning raw spatial features with text-derived semantics. The process
169 consists of two stages: (i) *Semantic enrichment*, utilizing VLMs to generate
170 descriptive text for each modality; and (ii) *Intra-modal alignment*, bridging
171 the gap between spatial and semantic representations to capture behavior-
172 related urban patterns.

173 *3.3.1. Semantic Enrichment of Multimodal Geospatial Data*

174 As shown in Fig. 2, each 500 m \times 500 m grid is characterized by four be-
175 haviorally informative modalities. *RSI* captures large-scale morphological
176 cues (e.g., layouts, vegetation), while *Building* height maps reflect vertical
177 urban intensity and trip-generation potential. *SVI* offers human-scale visual
178 attributes, sampled via 20 random images per grid. For grids with sparse
179 road networks containing fewer than 20 available images, oversampling with
180 replacement is applied to maintain consistent input dimensions while preserv-
181 ing realistic visual features. *POI* vectors describe functional composition and
182 accessibility. Although rich in structural data, these modalities lack explicit
183 semantic abstraction, limiting the model’s ability to interpret high-level con-
184 cepts like “commercial corridors.”

185 To address this, we leverage VLMs to inject interpretable semantic pri-
186 ors into the spatial representation [60]. Specifically, Qwen-VL processes vi-
187 sual inputs (*RSI*, *Building*, *SVI*) for its strong visual understanding, while
188 DeepSeek-R1 handles structured *POI* data via chain-of-thought reasoning
189 that first analyzes category distributions and then synthesizes functional de-
190 scriptions. These models serve as plug-and-play components that can be
191 replaced by alternatives without architectural changes (see Table 8 for a sys-
192 tematic comparison). A unified three-level instruction prompt guides the
193 generation of multi-dimensional urban semantics: (i) *land features and func-*
194 *tional zoning*; (ii) *estimated activity intensity (0–100)*; and (iii) *predicted*
195 *24-hour departure trends (0–1)*. The resulting descriptions are encoded via
196 Qwen3-Embedding for subsequent alignment.

197 *3.3.2. Intra-modal Representation Alignment*

198 This phase aligns VLM-derived semantics with original spatial features
199 to ensure consistency. It proceeds through dual-branch encoding, contrastive
200 alignment, and feature fusion.

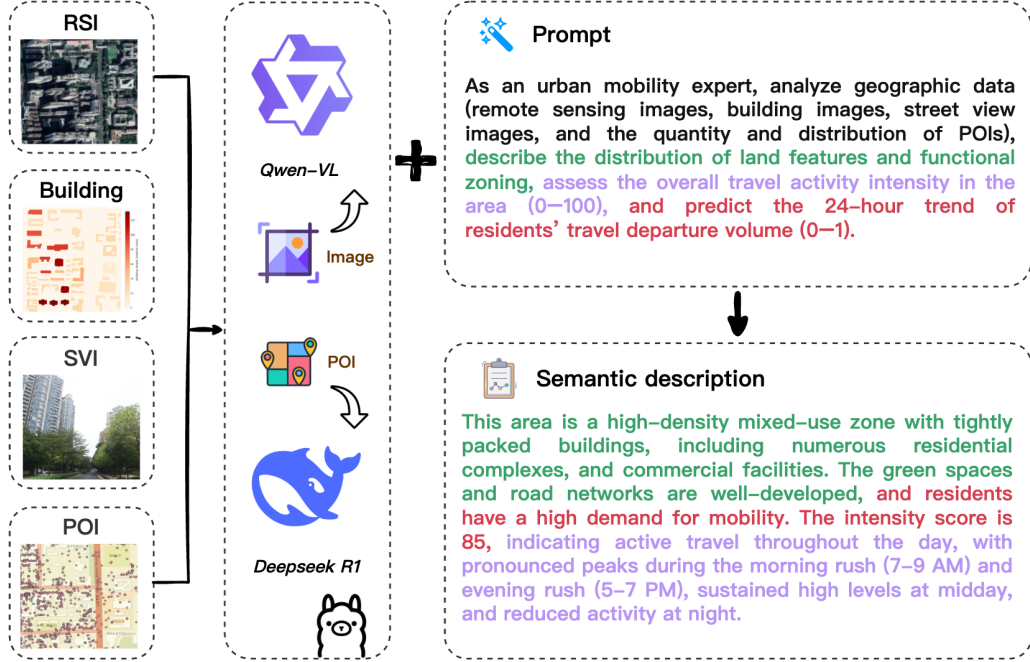


Figure 2: Illustration of text generation using LLMs/VLMs. Multiple models are employed to generate natural-language descriptions from multimodal geospatial data inputs.

201 *Dual-branch Encoding.* For each modality $m \in \{\text{rsi}, \text{bld}, \text{svi}, \text{poi}\}$, we extract
 202 parallel representations from the spatial input \mathbf{x}_i^m and textual description s_i^m .
 203 The spatial encoder $E_v^m(\cdot)$ maps \mathbf{x}_i^m to latent vector $\mathbf{f}_i^m \in \mathbb{R}^{d_f}$. Specifically,
 204 we employ a pre-trained ResNet-152 [61] for image-based modalities and an
 205 MLP for POI vectors. Simultaneously, the semantic encoder $E_t^m(\cdot)$ (Qwen3-
 206 Embedding [62]) processes s_i^m to produce $\mathbf{g}_i^m \in \mathbb{R}^{d_f}$. The pair $(\mathbf{f}_i^m, \mathbf{g}_i^m)$ thus
 207 encodes the grid from both physical and semantic perspectives.

208 *Intra-modal Contrastive Alignment.* To unify these branches, we employ a
 209 bidirectional contrastive framework [63] (Fig. 3). For all N grids, the pair
 210 $(\mathbf{f}_i^m, \mathbf{g}_i^m)$ constitutes a positive sample, while mismatched pairs $(\mathbf{f}_i^m, \mathbf{g}_j^m)_{i \neq j}$
 211 serve as negatives. We define directional similarity distributions as:

$$\begin{aligned}
 p_{i \rightarrow j}^m &= \frac{\exp(\cos(\mathbf{f}_i^m, \mathbf{g}_j^m)/\tau)}{\sum_{k=1}^N \exp(\cos(\mathbf{f}_i^m, \mathbf{g}_k^m)/\tau)}, \\
 q_{i \rightarrow j}^m &= \frac{\exp(\cos(\mathbf{g}_i^m, \mathbf{f}_j^m)/\tau)}{\sum_{k=1}^N \exp(\cos(\mathbf{g}_i^m, \mathbf{f}_k^m)/\tau)},
 \end{aligned} \tag{3}$$

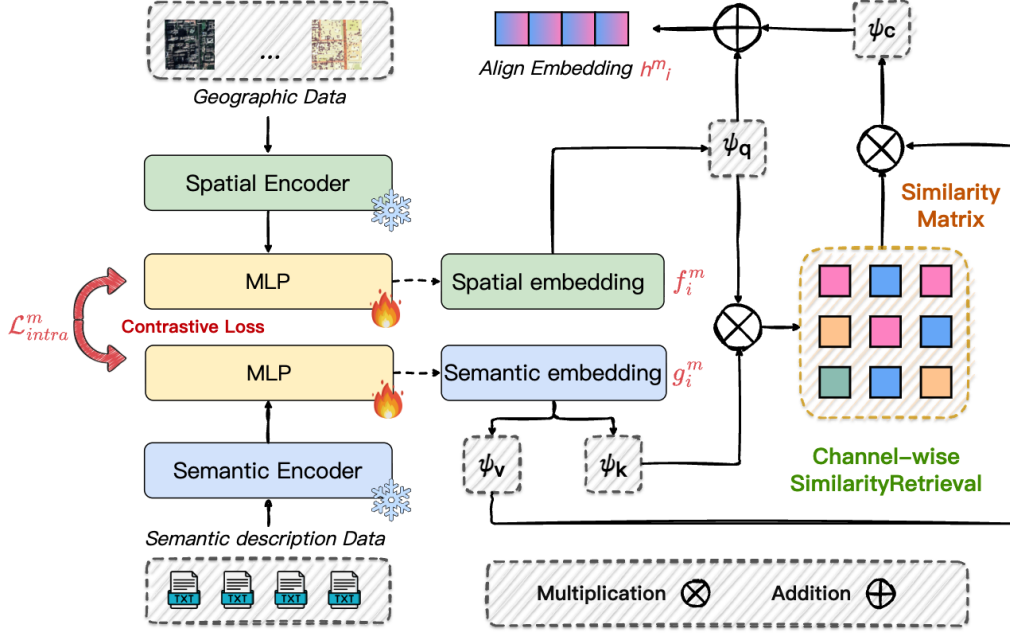


Figure 3: Architecture of the intra-modal contrastive alignment and feature fusion module.

212 where τ is the temperature coefficient. The bidirectional InfoNCE loss en-
 213 forces mutual alignment:

$$\mathcal{L}_{intra}^m = -\frac{1}{2N} \sum_{i=1}^N \left[\log p_{i \rightarrow i}^m + \log q_{i \rightarrow i}^m \right]. \quad (4)$$

214 This optimization creates a coherent latent space where spatial and textual
 215 features mutually reinforce each other.

216 *Intra-modal Feature Fusion.* Post-alignment, we integrate the branches using
 217 channel-wise similarity retrieval [60]. Let \mathbf{F}^m and \mathbf{G}^m denote the stacked spa-
 218 tial and semantic embeddings. We project these into subspaces via ψ_q, ψ_k, ψ_v
 219 to compute a similarity matrix:

$$\mathbf{S}^m = \text{rowsoftmax}((\psi_q(\mathbf{F}^m))^\top \psi_k(\mathbf{G}^m)). \quad (5)$$

220 The aligned representation \mathbf{H}^m is obtained via residual fusion:

$$\mathbf{H}^m = \omega^c(\psi_v(\mathbf{G}^m) \mathbf{S}^m) \oplus \mathbf{F}^m, \quad (6)$$

221 where ω^c is a learnable transformation. The resulting \mathbf{H}^m preserves geometric
 222 structure while incorporating high-level semantics, preparing robust inputs
 223 for cross-modal fusion.

224 3.4. Spatial-Contextual Representation Learning

225 Following intra-modal alignment, we address the spatial dependencies in-
 226 herent in urban environments. Nearby grids often exhibit strong correlations
 227 in morphology and function due to the continuity of the built environment
 228 and socioeconomic interactions [64]. To capture these contextual relation-
 229 ships, we employ a graph-based neural network to model spatial dependencies
 230 among grids, as illustrated in Fig. 1. The process consists of two key stages:
 231 (i) *Graph construction and initialization*, which establishes a shared spatial
 232 graph topology with modality-specific node features; and (ii) *Attention-based*
 233 *Spatial Aggregation*, which leverages GAT to adaptively propagate and ag-
 234 gregate contextual features from local neighborhoods.

235 3.4.1. Graph construction and initialization

236 To explicitly model such spatial relations, we construct a spatial graph
 237 $\mathcal{G}^m = (\mathcal{V}^m, \mathcal{E}^m)$ for each data modality $m \in \{\text{rsi, bld, svi, poi}\}$, sharing the
 238 same adjacency topology across modalities while initializing nodes with modality-
 239 specific features. Each node $v_i^m \in \mathcal{V}^m$ corresponds to an urban grid, and each
 240 edge $(v_i^m, v_j^m) \in \mathcal{E}^m$ encodes the spatial adjacency between two grids. The
 241 node v_i^m is initialized with the intra-modally aligned embedding \mathbf{h}_i^m obtained
 242 from the previous stage. The node feature matrix for the modality m is
 243 expressed as

$$\mathbf{H}^{(0,m)} = [\mathbf{h}_1^m, \mathbf{h}_2^m, \dots, \mathbf{h}_N^m]^\top \in \mathbb{R}^{N \times d_f}, \quad (7)$$

244 where N is the number of urban grids and d_f is the embedding dimension.
 245 Spatial relationships are encoded in an adjacency matrix $\mathbf{A}^m \in \mathbb{R}^{N \times N}$, where
 246 $A_{ij}^m = 1$ if the grids i and j are spatially adjacent (e.g., share a boundary)
 247 and $A_{ij}^m = 0$ otherwise. This shared graph topology, combined with modality-
 248 specific node features, provides the foundation for capturing spatial context
 249 and non-local dependencies in the subsequent graph-based representation
 250 learning. We denote the updated graph-enhanced representation as $\tilde{\mathbf{H}}^m =$
 251 $[\tilde{\mathbf{h}}_1^m, \tilde{\mathbf{h}}_2^m, \dots, \tilde{\mathbf{h}}_N^m]^\top$.

252 3.4.2. Attention-based Spatial Aggregation

253 For each modality, we adopt a GAT [65] to adaptively propagate con-
 254 textual information among neighboring grids. As shown in Fig. 1, three

255 stacked *GATResBlocks* are employed to deepen the receptive field and cap-
 256 ture multi-scale spatial dependencies. At the l -th layer, the hidden state is
 257 first normalized and then passed through a graph attention operator:

$$\mathbf{h}_i^{(l+1,m)} = \text{GELU} \left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij}^m \mathbf{W}^m \mathbf{h}_j^{(l,m)} \right), \quad (8)$$

258 where $\mathcal{N}(i) = \{j \mid A_{ij}^m = 1\}$, \mathbf{W}^m is a learnable weight matrix, and the
 259 attention coefficients α_{ij}^m are normalized using a softmax function:

$$\alpha_{ij}^m = \text{softmax}_{j \in \mathcal{N}(i)} \left(\text{LeakyReLU} \left((\mathbf{a}^m)^\top [\mathbf{W}^m \mathbf{h}_i^{(l,m)} \parallel \mathbf{W}^m \mathbf{h}_j^{(l,m)}] \right) \right), \quad (9)$$

260 with \mathbf{a}^m being a learnable attention vector and \parallel denoting the concatenation
 261 of features. Dropout and residual connections are applied to stabilize training
 262 and enhance convergence.

$$\tilde{\mathbf{h}}_i^{(m)} = \mathbf{h}_i^{(l,m)} + \text{Dropout}(\mathbf{h}_i^{(l+1,m)}). \quad (10)$$

263 This attention-based residual design (denoted as *GATResBlock*) allows the
 264 network to dynamically assign importance weights to neighboring nodes, ef-
 265 fectively capturing multi-scale spatial dependencies across different modali-
 266 ties.

267 3.5. Cross-modal Fusion and Temporal Modeling

268 After obtaining semantically aligned and spatially contextualized embed-
 269 dings for each modality, the next step is to perform cross-modal fusion and
 270 learn the temporal dynamics of human mobility. This stage aims to (i)
 271 achieve inter-modal semantic consistency through contrastive learning, (ii)
 272 integrate complementary multimodal features via attention-based fusion, and
 273 (iii) model 24-hour temporal variations using a Transformer-based sequence
 274 encoder.

275 3.5.1. Inter-modal Contrastive Learning

276 Although intra-modal alignment ensures semantic consistency within each
 277 individual modality, the embeddings of different modalities (e.g., RSI, build-
 278 ing, SVI, and POI) may still reside in heterogeneous latent spaces. To unify

279 these representations, we design an *inter-modal contrastive learning* frame-
 280 work that encourages cross-modal correspondence among semantically re-
 281 lated regions.

282 Specifically, for any pair of modalities (m, n) , we regard the graph-enhanced
 283 embeddings $\tilde{\mathbf{h}}_i^m$ and $\tilde{\mathbf{h}}_i^n$ corresponding to the same spatial location i as a *pos-*
 284 *itive pair*, while embeddings from different locations $(\tilde{\mathbf{h}}_i^m, \tilde{\mathbf{h}}_j^n)$ ($i \neq j$) serve
 285 as *negative pairs*. Following the contrastive paradigm introduced in Sec-
 286 tion 3.3.2, we compute pairwise similarities between modality-specific em-
 287 beddings using cosine similarity and apply softmax normalization across all
 288 N grids. The bidirectional InfoNCE objective is then formulated as:

$$\mathcal{L}_{\text{inter}}^{(m,n)} = -\frac{1}{2N} \sum_{i=1}^N \left[\log p_{i \rightarrow i}^{(m,n)} + \log q_{i \rightarrow i}^{(n,m)} \right], \quad (11)$$

289 where $p_{i \rightarrow i}^{(m,n)}$ and $q_{i \rightarrow i}^{(n,m)}$ denote the normalized similarity probabilities from
 290 the modality m to n and n to m , respectively, and τ is the shared temperature
 291 coefficient that controls the sharpness of the distribution.

292 3.5.2. Multimodal Fusion

293 Let $\{\tilde{\mathbf{H}}^{\text{rsi}}, \tilde{\mathbf{H}}^{\text{bld}}, \tilde{\mathbf{H}}^{\text{svi}}, \tilde{\mathbf{H}}^{\text{poi}}\}$ denote the modality-specific embedding matri-
 294 ces from Section 3.4, where each $\tilde{\mathbf{H}}^m = [\tilde{\mathbf{h}}_1^m, \tilde{\mathbf{h}}_2^m, \dots, \tilde{\mathbf{h}}_N^m]^\top \in \mathbb{R}^{N \times d_f}$ represents
 295 the features of all N urban grids under modality m , and $\tilde{\mathbf{h}}_i^m \in \mathbb{R}^{d_f}$ denotes the
 296 feature vector of grid i . For each grid i , we construct a modality token matrix
 297 $\mathbf{T}_i = [\tilde{\mathbf{h}}_i^{\text{rsi}}, \tilde{\mathbf{h}}_i^{\text{bld}}, \tilde{\mathbf{h}}_i^{\text{svi}}, \tilde{\mathbf{h}}_i^{\text{poi}}]^\top \in \mathbb{R}^{4 \times d_f}$ and apply a lightweight Transformer to
 298 learn inter-modal dependencies:

$$\mathbf{z}_i = \text{AttnPool}(\text{Transformer}_{\text{modal}}(\mathbf{T}_i)), \quad (12)$$

299 where $\text{AttnPool}(\cdot)$ aggregates modality-aware features into a unified represen-
 300 tation. The fused feature matrix is then expressed as $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N]^\top \in$
 301 $\mathbb{R}^{N \times d_f}$, where each \mathbf{z}_i encodes the integrated semantic and structural char-
 302 acteristics of grid i .

303 3.5.3. Temporal Modeling

304 To capture the diurnal variation of human mobility, the fused embeddings
 305 \mathbf{Z} are temporally expanded into a 24-step sequence for each spatial unit. Let
 306 $\{\mathbf{e}_t\}_{t=1}^{24}$ denote learnable hour embeddings. We form per-hour contexts by
 307 conditioning the fused feature on \mathbf{e}_t :

$$\mathbf{z}_i^{(t)} = \phi([\mathbf{z}_i \parallel \mathbf{e}_t]) \in \mathbb{R}^{d_f}. \quad (13)$$

308 This yields $\mathbf{Z}_{\text{seq}} \in \mathbb{R}^{N \times 24 \times d_f}$, where each spatial unit is represented by a 24-
 309 step temporal feature sequence. A Transformer encoder [66] is then applied
 310 to model long-range dependencies across hours, producing updated temporal
 311 representations $\mathbf{U} = \text{TransformerEncoder}(\mathbf{Z}_{\text{seq}}) \in \mathbb{R}^{N \times 24 \times d_u}$.

312 For each grid i and hour t , the model outputs a probability distribution
 313 over C magnitude bins through a classification head:

$$\hat{\mathbf{p}}_{i,t} = \text{softmax}(\text{MLP}_{\text{cla}}(\mathbf{U}_{i,t})), \quad \hat{\mathbf{p}}_{i,t} \in \mathbb{R}^C. \quad (14)$$

314 Following the classification-then-regression strategy, the continuous outbound
 315 flow $\hat{y}_{i,t}$ is calculated as the expectation of the predicted bin distribution.
 316 Let the n -th bin correspond to the interval $[b_n, b_{n+1}]$ with midpoint $m_n =$
 317 $(b_n + b_{n+1})/2$. The final predicted mobility intensity is obtained by:

$$\hat{y}_{i,t} = \sum_{n=1}^C \hat{p}_{i,t}^{(n)} \cdot m_n, \quad (15)$$

318 where $\hat{p}_{i,t}^{(n)}$ is the predicted probability of bin n . This expectation-based
 319 formulation allows the model to produce smooth and physically meaningful
 320 continuous predictions while benefiting from the stability of a classification
 321 objective.

322 3.6. Loss Functions and Optimization

323 To jointly optimize multimodal, spatial, and temporal representations,
 324 the proposed framework employs two task-specific objectives—classification
 325 and regression—along with a contrastive alignment loss introduced in Sec-
 326 tions 3.3 and 3.5. The overall training objective ensures that the model si-
 327 multaneously learns to predict mobility intensity levels, capture fine-grained
 328 temporal variations, and maintain semantic consistency across modalities.

329 (i) *Classification Loss*: For each urban grid i and hour t , the model
 330 predicts a categorical travel intensity level from C discrete bins. If the
 331 ground-truth label $y_{i,t}$ is represented as a single class index (*hard labels*),
 332 the classification loss adopts a standard cross-entropy form:

$$\mathcal{L}_{\text{cla}} = -\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \log \hat{p}_{i,t}[y_{i,t}], \quad (16)$$

333 where $\hat{p}_{i,t}[y_{i,t}]$ denotes the predicted probability for the true class of grid i
 334 at hour t , and N and T represent the total number of grids and time steps
 335 ($T = 24$).

336 (ii) *Regression Loss*: To capture fine-grained variations in hourly travel
 337 intensity, a regression objective is introduced to directly estimate the contin-
 338 uous outbound flow volume. It is formulated as a mean squared error (MSE)
 339 between the predicted and ground-truth sequences:

$$\mathcal{L}_{\text{reg}} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (\hat{y}_{i,t} - y_{i,t})^2, \quad (17)$$

340 where $\hat{y}_{i,t}$ and $y_{i,t}$ denote the predicted and observed hourly flow values for
 341 grid i at hour t .

342 (iii) *Contrastive Alignment Loss*: To ensure consistent multimodal se-
 343 mantics, the total contrastive loss \mathcal{L}_{CL} combines both the intra- and inter-
 344 modal contrastive objectives defined in previous sections:

$$\mathcal{L}_{\text{CL}} = \sum_{m \in \mathcal{M}} \mathcal{L}_{\text{intra}}^m + \frac{1}{\binom{|\mathcal{M}|}{2}} \sum_{m < n} \mathcal{L}_{\text{inter}}^{(m,n)}, \quad (18)$$

345 where $\mathcal{L}_{\text{intra}}^m$ aligns spatial and semantic features within modality m , and
 346 $\mathcal{L}_{\text{inter}}^{(m,n)}$ enforces consistency between different modalities.

347 (iv) *Overall Optimization*: The final loss function jointly optimizes the
 348 classification, regression, and contrastive objectives:

$$\mathcal{L}_{\text{total}} = \alpha \mathcal{L}_{\text{cla}} + \beta \mathcal{L}_{\text{reg}} + \lambda \mathcal{L}_{\text{CL}}, \quad (19)$$

349 where α , β , and λ are weighting coefficients controlling the relative impor-
 350 tance of the three components. This multi-objective formulation enables the
 351 model to simultaneously enhance prediction accuracy, temporal smoothness,
 352 and multimodal semantic coherence across the entire training process.

353 4. Experiments and Results

354 4.1. Datasets and Experimental Setup

355 4.1.1. Datasets

356 This study selected Shenzhen, China as the study area, dividing the built-
 357 up region into 5,263 grid tiles of 500 m \times 500 m (Fig. 4). We integrated four
 358 types of static geospatial data: *RSI* from Google Earth (2.15 m resolution);
 359 *Building* height data from Gaode Map; *SVI* from Baidu Map (\approx 78k raw im-
 360 ages, 1024 \times 768, oversampled to 20 per grid); and *POI* data (\approx 1.5 million)

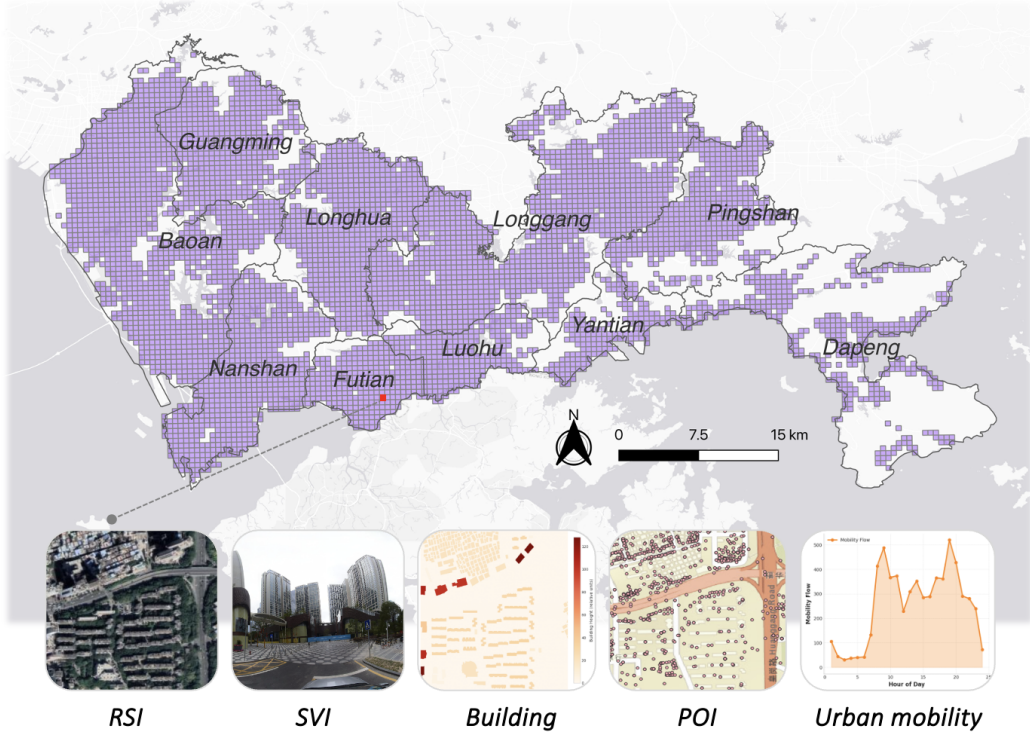


Figure 4: Study area and multimodal geospatial data overview. Shenzhen comprises 10 administrative districts, and its built-up region is partitioned into 5,263 grids of 500 m × 500 m, with four static data modalities illustrated.

361 from Gaode Map. The prediction target is the hourly outbound mobility de-
 362 rived from anonymized mobile signaling data. We utilized weekday records
 363 from June 18 to June 25, 2019. The mobility flows were log-transformed, Z-
 364 score normalized, and discretized into 10 bins. All modalities were collected
 365 within a one-year window centered on June 2019, minimizing temporal mis-
 366 alignment. The dataset was divided into training, validation, and test sets
 367 with a ratio of 6:2:2 using spatially random masks applied at the grid level.

368 4.1.2. Experimental Setup

369 We adopt *Overall Accuracy (OA)* for classification and *Coefficient of De-*
 370 *termination (R^2)*, *Mean Absolute Error (MAE)*, and *Mean Absolute Percent-*
 371 *age Error (MAPE)* for regression. The framework is implemented in PyTorch
 372 and trained on an NVIDIA RTX A6000 (48GB) for 300 epochs (Adam opti-
 373 mizer, $lr = 1 \times 10^{-4}$ with cosine annealing). Key hyperparameters: $d_f = 256$,

374 $\tau = 0.1, \alpha=0.5, \beta=0.5, \lambda=0.1$.

375 To improve efficiency and transferability, the framework separates com-
376 ponents into *frozen* and *trainable* groups. Frozen modules—VLM encoders
377 (Qwen-VL, DeepSeek-R1) for offline semantic generation, spatial encoders
378 (ResNet-152) for image features, and the text encoder (Qwen3-Embedding)—
379 leverage large-scale pretraining without fine-tuning. Trainable modules—
380 GAT layers, cross-modal Transformer fusion, temporal Transformer encoder,
381 and prediction heads—are optimized via backpropagation. Fig. 5 illustrates
382 the training process with stable convergence.

383 4.2. Results

384 This section presents the experimental results of the proposed LMEMR
385 for grid-level mobility prediction. We comprehensively evaluate the model
386 in terms of quantitative performance and interpretability. Comparisons with
387 baseline models and ablation variants are provided to demonstrate the ad-
388 vantages of the proposed modules.

389 4.2.1. Overall Performance Comparison with Baselines

390 To evaluate LMEMR, we compare it against representative machine learn-
391 ing, deep learning baselines, and specialized spatiotemporal models. Except
392 for CLIP, all baselines utilize the same 256-dimensional raw spatial encod-
393 ings \mathbf{f}_i^m (Section 3.3.2) without textual semantics or contrastive learning.
394 The baselines are:

- 395 • *Linear Regression (LR)*. Maps grid embeddings to hourly mobility in-
396 tensity via ordinary least squares [67].
- 397 • *Random Forest (RF)*. An ensemble of 500 decision trees capturing non-
398 linear feature–mobility relationships [68].
- 399 • *CNN*. A 2D convolutional model capturing local spatial correlations,
400 lacking explicit temporal reasoning [69].
- 401 • *LSTM*. A two-layer RNN modeling the 24-hour sequence of grid em-
402 beddings to predict mobility dynamics [70].
- 403 • *Transformer*. A temporal attention model learning long-range depen-
404 dencies across hourly time steps [66].

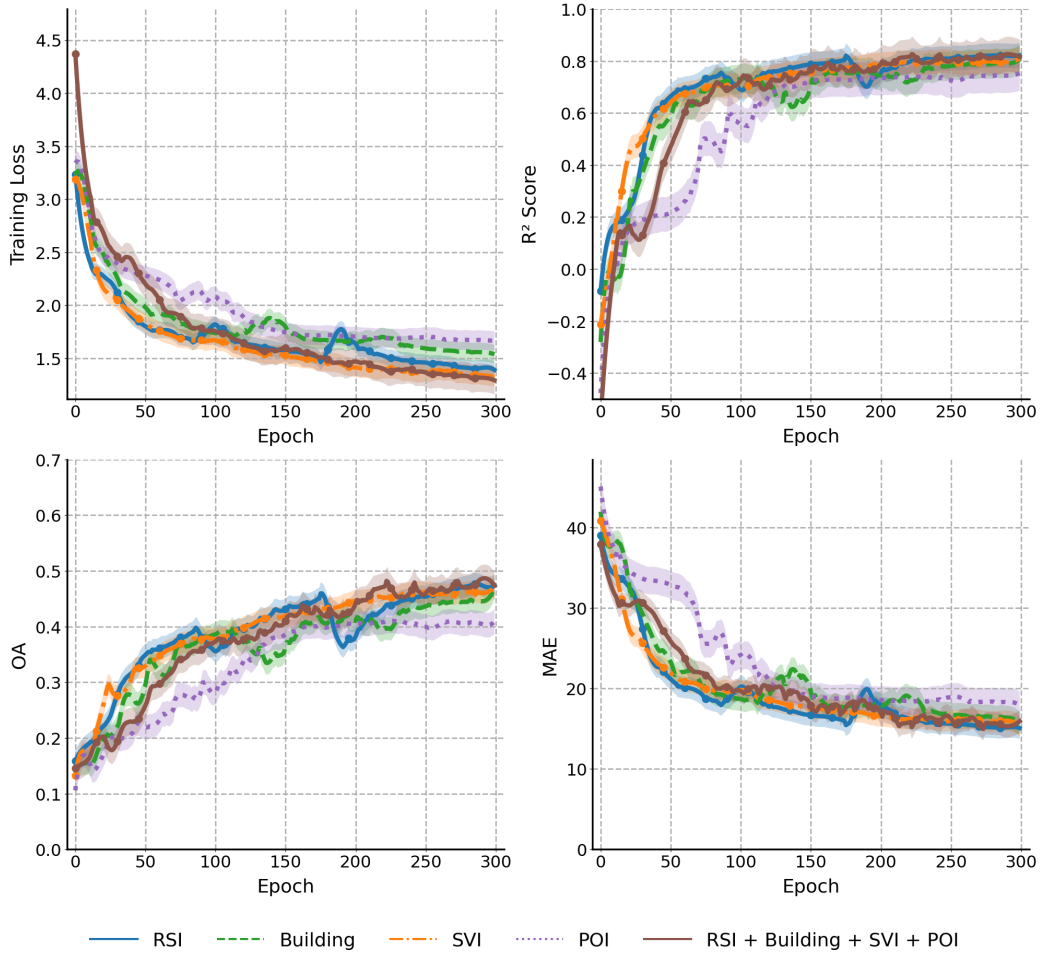


Figure 5: Training and validation curves of the proposed model. The plots illustrate the evolution of Training Loss, OA, R^2 , and MAE over 300 epochs for different modality configurations.

- 405
406
407
 • *CLIP*. A multimodal baseline using pretrained CLIP ViT-L/14 for visual encoding (RSI, Building, SVI) and CLIP Text Encoder for POI, retaining the same spatial and temporal modules as LMEMR [54].

408
409
410
411
412
 As presented in Table 2, LMEMR outperforms all baselines. Traditional models (LR, RF) show limited accuracy due to insufficient spatial-behavioral modeling. Deep learning approaches (CNN, LSTM) improve performance by capturing local spatial or temporal dependencies, with the Transformer further enhancing results via attention mechanisms. The multimodal baseline

413 (CLIP) leverages pretrained vision-language embeddings but achieves sub-
 414 optimal results ($R^2 = 0.761$), indicating that generic multimodal represen-
 415 tations lack domain-specific spatial-behavioral knowledge. LMEMR signifi-
 416 cantly surpasses all baselines with inference time of 17.2s and training time
 417 of 2,327s, confirming that jointly modeling spatial, semantic, and temporal
 418 dimensions offers a superior understanding of urban mobility with reasonable
 419 computational cost.

Table 2: Overall predictive performance comparison with baselines. Best in **bold**, second best underlined.

Model	OA \uparrow	R^2 \uparrow	MAE \downarrow	MAPE \downarrow	Training (s)	Inference (s) \downarrow
LR	0.237	0.452	26.794	0.405	135	10.3
RF	0.226	0.521	24.916	0.374	<u>271</u>	10.3
CNN	0.411	0.734	18.904	0.259	1283	<u>11.9</u>
LSTM	0.423	0.733	17.882	0.264	1225	12.3
Transformer	0.425	0.755	17.525	0.249	1451	14.9
CLIP	<u>0.430</u>	<u>0.761</u>	<u>16.769</u>	<u>0.245</u>	1542	14.9
LMEMR	0.521	0.856	13.738	0.187	2327	17.2

420 4.2.2. Performance in Temporal Dynamics

421 Fig. 6 demonstrates that LMEMR accurately reproduces diurnal mobil-
 422 ity dynamics across diverse urban environments. The model successfully
 423 captures the *morning-peak* in dense residential areas (Baoan urban village),
 424 the *evening-peak* in office clusters (Futian CBD), the *stable flow* in trans-
 425 port hubs (Longhua railway station), and the *dual-peak* pattern in mixed-use
 426 districts (Luohu). The predictions align closely with observed trends, show-
 427 ing high fidelity in both magnitude and temporal variation, with only minor
 428 deviations in highly dynamic zones.

429 4.2.3. Performance in Mobility Intensity Classification

430 Fig. 7 compares the classification performance across ten mobility inten-
 431 sity bins. LMEMR exhibits the strongest diagonal concentration (OA=0.521),
 432 significantly outperforming Transformer (OA=0.425) and RF (OA=0.226).
 433 Specifically, LMEMR reduces misclassification in the dominant low-to-mid
 434 intensity bins (0–5) and maintains robustness in high-activity bins (6–8), ef-
 435 fectively mitigating the over-smoothing and underestimation issues observed

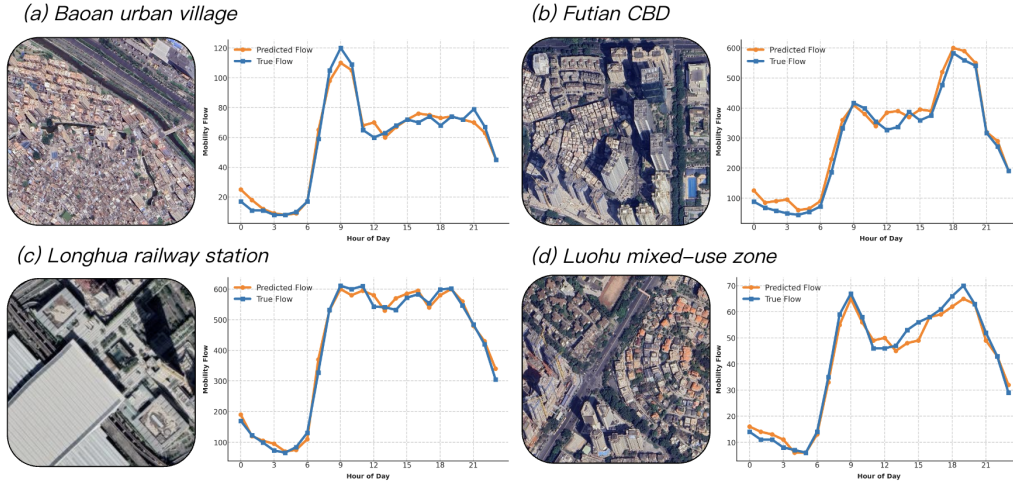


Figure 6: Predicted and observed hourly mobility in representative urban areas. Each pair shows a high-resolution remote sensing image (left) and its corresponding flow curve (right): (a) Baoan urban village (morning-peak), (b) Futian CBD (evening-peak), (c) Longhua railway station (stable flow), and (d) Luohu mixed-use zone (dual peaks).

436 in baselines. This confirms that integrating multimodal semantics and spa-
 437 tial reasoning enhances the model’s precision in distinguishing fine-grained
 438 mobility states.

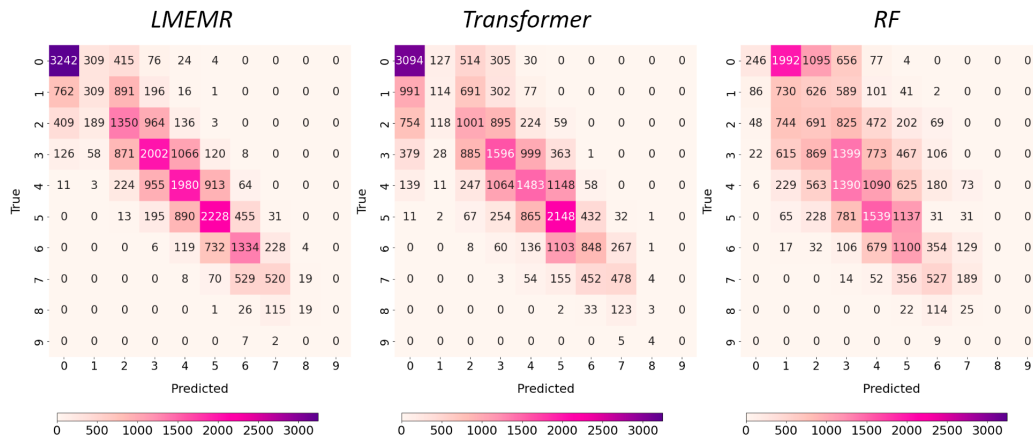


Figure 7: Confusion matrices of LMEMR, Transformer, and RF across ten mobility intensity bins (0–9). LMEMR achieves the strongest diagonal concentration and highest accuracy (OA=0.521), clearly outperforming Transformer (0.425) and RF (0.226).

439 *4.2.4. Cross-City Generalization Analysis*

440 To validate the generalizability of LMEMR across diverse urban con-
 441 texts, we conducted independent experiments on three additional cities in
 442 the Greater Bay Area: Guangzhou (mega city, 8,054 grids), Dongguan (large
 443 industrial city, 5,022 grids), and Zhuhai (medium-sized tourism city, 2,173
 444 grids), each trained and evaluated independently using identical protocols
 445 (6:2:2 split). As shown in Table 3, LMEMR consistently outperforms the
 446 CLIP baseline across all four cities, achieving high predictive accuracy with
 447 R^2 ranging from 0.734 to 0.863 across cities of varying scales and economic
 448 structures. Guangzhou achieves the highest R^2 (0.863), surpassing even
 449 Shenzhen (0.856), likely influenced by its larger sample size. Dongguan
 450 shows the largest relative improvement in R^2 over CLIP, suggesting that
 451 VLM-generated semantics offer notable benefits for cities with different ur-
 452 ban structures. Zhuhai maintains competitive performance ($R^2 = 0.810$)
 453 despite the smallest sample size. These results confirm that LMEMR gener-
 454 alizes effectively across diverse urban contexts.

Table 3: Cross-city generalization validation across four cities in the Greater Bay Area.

City	Model	OA \uparrow	R^2 \uparrow	MAE \downarrow	MAPE \downarrow
Shenzhen	LMEMR	0.521	0.856	13.738	0.187
	CLIP	0.430	0.761	16.769	0.245
Guangzhou	LMEMR	0.578	0.863	10.943	0.238
	CLIP	0.513	0.777	12.934	0.283
Dongguan	LMEMR	0.421	0.734	5.967	0.273
	CLIP	0.374	0.644	6.339	0.337
Zhuhai	LMEMR	0.596	0.810	3.909	0.244
	CLIP	0.543	0.733	4.138	0.370

455 *4.3. Ablation Studies*

456 To isolate the contribution of each proposed component, all ablation vari-
 457 ants are independently retrained from scratch using the same training con-
 458 figuration (300 epochs, Adam optimizer, lr= 1×10^{-4}) and identical data
 459 splits. For each variant, only the targeted module is modified while all other
 460 components remain unchanged, ensuring fair and controlled comparisons.

461 *4.3.1. Effectiveness of the Semantic Enhancement Module*

462 To evaluate the impact of VLM-derived semantics and contrastive align-
 463 ment, four variants are developed: (i) *w/o Semantics*, which removes VLM-
 464 generated textual descriptions and uses only raw visual features from the
 465 spatial encoder E_v^m ; (ii) *w/o Intra-CL*, which removes the intra-modal con-
 466 trastive loss $\mathcal{L}_{\text{intra}}^m$ while keeping all other components; (iii) *Prompt-Lite*,
 467 which uses simplified prompts without behavioral cues; and (iv) *Prompt-*
 468 *Shuffle*, which employs randomized captions to test robustness.

469 As summarized in Table 4, LMEMR achieves superior performance across
 470 all metrics. The removal of textual semantics leads to a significant degrada-
 471 tion, indicating that high-level semantic cues are essential for predictive accu-
 472 racy. Similarly, the performance drop in *w/o Intra-CL* confirms the necessity
 473 of explicitly aligning visual and textual representations. The intermediate
 474 results of *Prompt-Lite* and the poor performance of *Prompt-Shuffle* further
 475 validate that behavior-aware prompting and accurate text–location corre-
 476 spondence are critical for capturing travel dynamics. Qualitative analysis in
 477 Fig. 8 further illustrates this contribution. LMEMR accurately reproduces
 478 diurnal peaks, whereas the variant without semantics underestimates ampli-
 479 tudes and misaligns timing. The generated descriptions explicitly link built-
 480 environment characteristics (e.g., “mixed-use,” “office clusters”) with mobility
 481 behaviors, serving as semantic priors that guide the model to better associate
 482 spatial structures with travel rhythms.

Table 4: Ablation of the semantic enhancement module. Best in **bold**, second best underlined.

Variant	OA \uparrow	R ² \uparrow	MAE \downarrow	MAPE \downarrow
w/o Semantics	0.491	0.825	14.833	0.213
w/o Intra-CL	0.505	0.837	14.285	0.209
Prompt-Lite	<u>0.514</u>	<u>0.850</u>	<u>13.783</u>	<u>0.199</u>
Prompt-Shuffle	0.475	0.811	15.426	0.228
LMEMR	0.521	0.856	13.738	0.187

483 *4.3.2. Contribution of Spatial-Contextual Representation Learning*

484 To quantify the role of spatial dependencies, we compare against two
 485 variants: (i) *w/o Spatial*, where grid embeddings are fed directly into the

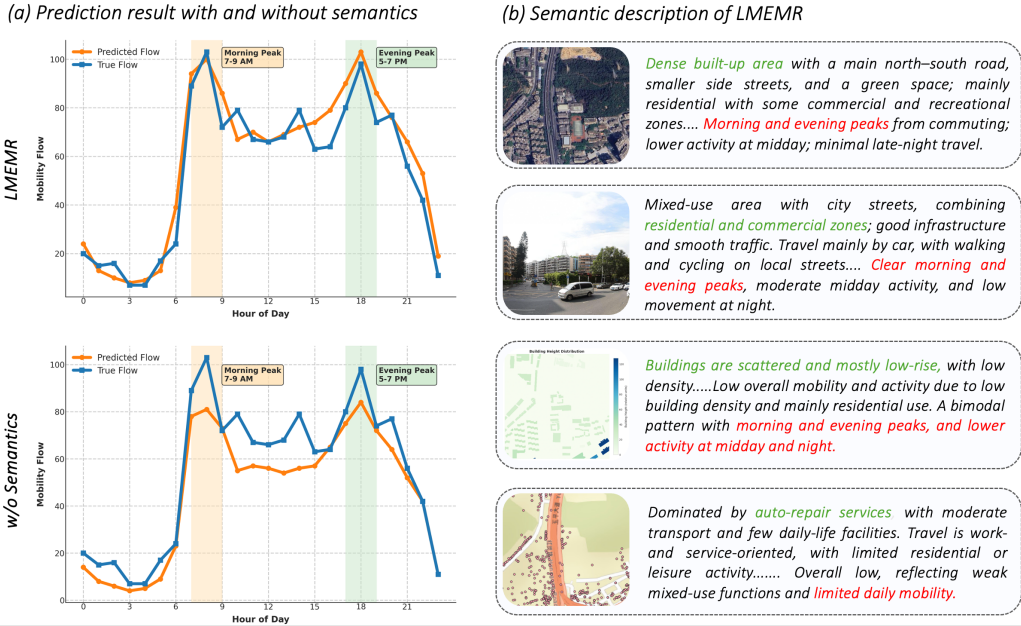


Figure 8: Comparison of mobility flow predictions and semantic interpretation. (a) LMEMR captures morning and evening peaks more accurately than the version without semantics. (b) VLM-generated descriptions provide explanatory context for the predicted patterns.

486 temporal encoder without spatial aggregation (GAT removed); and (ii) GCN,
 487 which replaces the attention-based GAT with fixed-weight mean aggregation
 488 using standard graph convolution [71].

489 Table 5 reports the results. The complete LMEMR framework yields
 490 the highest accuracy, while removing the spatial module causes a substan-
 491 tial drop in R^2 (from 0.856 to 0.764), underscoring that spatial context is
 492 indispensable for grid-level prediction. The GCN variant improves upon the
 493 non-spatial baseline but remains inferior to the GAT-based approach, demon-
 494 strating the advantage of adaptive attention in capturing non-uniform spatial
 495 interactions. Fig. 9 visualizes the learned attention weights. The spatial dis-
 496 tribution reveals a clear urban hierarchy, with high-weight areas concentrated
 497 in core commercial corridors and transit hubs (e.g., Luohu-Futian-Nanshan).
 498 These zones exert strong spatial influences on adjacent regions, justifying the
 499 higher contextual weights assigned by the model. Furthermore, modality-
 500 specific maps show distinct hotspots, reflecting the complementary roles of
 501 structural (RSI/Building) and functional (SVI/POI) features in shaping ur-

502 ban mobility.

Table 5: Ablation of spatial-contextual representation learning. Best in **bold**, second best underlined.

Variant	OA \uparrow	R ² \uparrow	MAE \downarrow	MAPE \downarrow
w/o Spatial	0.434	0.764	17.434	0.247
GCN	<u>0.488</u>	<u>0.820</u>	<u>15.061</u>	<u>0.217</u>
LMEMR	0.521	0.856	13.738	0.187

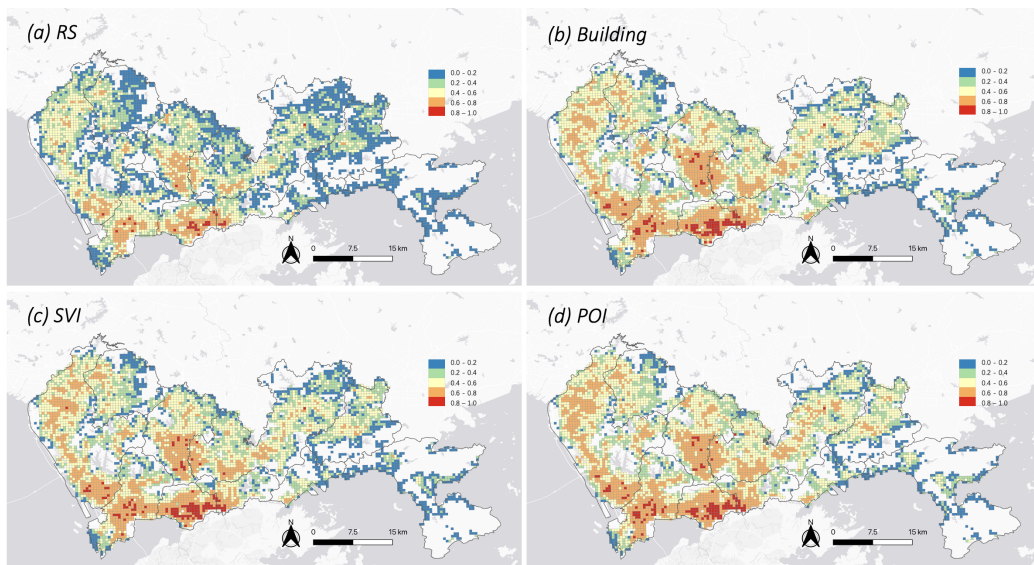


Figure 9: Spatial distribution of learned attention weights. High coefficients (red) align with core urban districts (e.g., Luohu, Futian), indicating strong spatial dependencies in dense functional zones.

503 *4.3.3. Role of Cross-modal Semantic Fusion Module*

504 This section investigates the contribution of cross-modal mechanisms
 505 through four sets of variants: (i) *w/o Inter-CL*, which disables inter-modal
 506 contrastive loss $\mathcal{L}_{\text{inter}}$ while retaining intra-modal alignment and the Trans-
 507 former fusion; (ii) *Late-Concat*, which replaces the Transformer-based fu-
 508 sion with simple feature concatenation followed by MLP; (iii) *Gated-Fusion*,
 509 which uses gated summation instead of attention-based fusion; and (iv)

510 *Single-Mod Drop*, which retains only individual modalities while keeping the
 511 GAT and temporal modules.

512 Results in Table 6 demonstrate that the full LMEMR model achieves
 513 the highest accuracy. The degradation observed in *w/o Inter-CL* confirms
 514 that explicit cross-modal alignment fosters consistent representations among
 515 heterogeneous data. Furthermore, the attention-based fusion outperforms
 516 both concatenation and gating, indicating that the multi-head mechanism
 517 better captures nonlinear cross-modal dependencies. Among single-modality
 518 variants, RSI and SVI perform relatively well, suggesting that morphology
 519 and visual perception are primary determinants of mobility, whereas POI
 520 data alone is insufficient. Overall, the findings validate the necessity of the
 521 proposed contrastive-Transformer fusion design for robust multimodal inte-
 522 gration.

Table 6: Ablation of cross-modal semantic fusion. Best in **bold**, second best underlined.

Variant	OA \uparrow	R^2 \uparrow	MAE \downarrow	MAPE \downarrow
w/o Inter-CL	0.492	0.828	14.879	0.219
Late-Concat	<u>0.507</u>	<u>0.831</u>	<u>14.377</u>	0.215
Gated-Fusion	0.495	0.827	14.787	<u>0.213</u>
RSI-only	0.482	0.820	14.924	0.222
Building-only	0.478	0.810	15.381	0.218
SVI-only	0.478	0.815	15.016	0.214
POI-only	0.435	0.764	17.571	0.243
LMEMR	0.521	0.856	13.738	0.187

523 4.3.4. Effect of Prediction Head Design

524 To evaluate the contribution of the classification-then-regression strat-
 525 egy, we compare three loss configurations: (i) *Regression only*, which uses
 526 MSE loss directly on continuous flow values; (ii) *Classification only*, which
 527 discretizes flow into intensity bins and uses cross-entropy loss without regres-
 528 sion refinement; (iii) *Classification + Regression (LMEMR)*, the full design
 529 combining both losses.

530 As shown in Table 7, the combined strategy ($R^2 = 0.856$) outperforms
 531 both regression-only ($R^2 = 0.845$) and classification-only ($R^2 = 0.828$) ap-
 532 proaches. Regression-only already achieves competitive performance, but the

Table 7: Ablation of prediction head design. Best in **bold**, second best underlined.

Variant	OA \uparrow	R^2 \uparrow	MAE \downarrow	MAPE \downarrow
Regression-Only	<u>0.497</u>	<u>0.845</u>	<u>14.102</u>	<u>0.203</u>
Classification-Only	0.492	0.828	14.696	0.215
LMEMR	0.521	0.856	13.738	0.187

533 distributional approach further improves by handling the long-tail mobility
 534 distribution through probability bins, mitigating extreme value influence dur-
 535 ing training. Classification-only performs worst, indicating that discretiza-
 536 tion alone loses important magnitude information. These results confirm
 537 that the two components are complementary.

538 5. Discussion

539 5.1. Error Analysis across Different Hours

540 As illustrated in Fig. 10, model performance exhibits a distinct diurnal
 541 pattern linked to the regularity of human mobility. Quantitative metrics in-
 542 dicate that prediction precision is highest during stable, low-mobility periods
 543 (e.g., early morning $R^2 > 0.84$) and decreases during transient phases such as
 544 morning (07:00–09:00) and evening (18:00–19:00) peaks. Specifically, MAPE
 545 peaks at 04:00 and 06:00, reflecting sensitivity to rapid shifts in travel inten-
 546 sity. This temporal variation aligns with previous findings that static features
 547 effectively capture structural determinants (e.g., land use) but struggle with
 548 volatile congestion dynamics [72]. Nevertheless, the consistently high R^2 and
 549 low average MAPE confirm that multimodal semantics and spatial-graph
 550 reasoning substantially enhance temporal robustness.

551 5.2. Semantic Analysis of Mobility Patterns

552 Hierarchical clustering of 24-hour mobility vectors reveals spatially dis-
 553 tinct travel rhythms rooted in urban functionality. As shown in Fig. 11 (a),
 554 increasing clusters to $K=7$ progressively refines the urban structure from a
 555 simple built/non-built dichotomy to a polycentric hierarchy consistent with
 556 functional zones [73]. Note that $K=7$ is not claimed as an optimal value but
 557 serves as an exploratory setting that balances granularity and interpretabil-
 558 ity; the hierarchical evolution from $K=2$ to $K=7$ demonstrates that the
 559 discovered patterns remain coherent across resolutions. At this resolution,

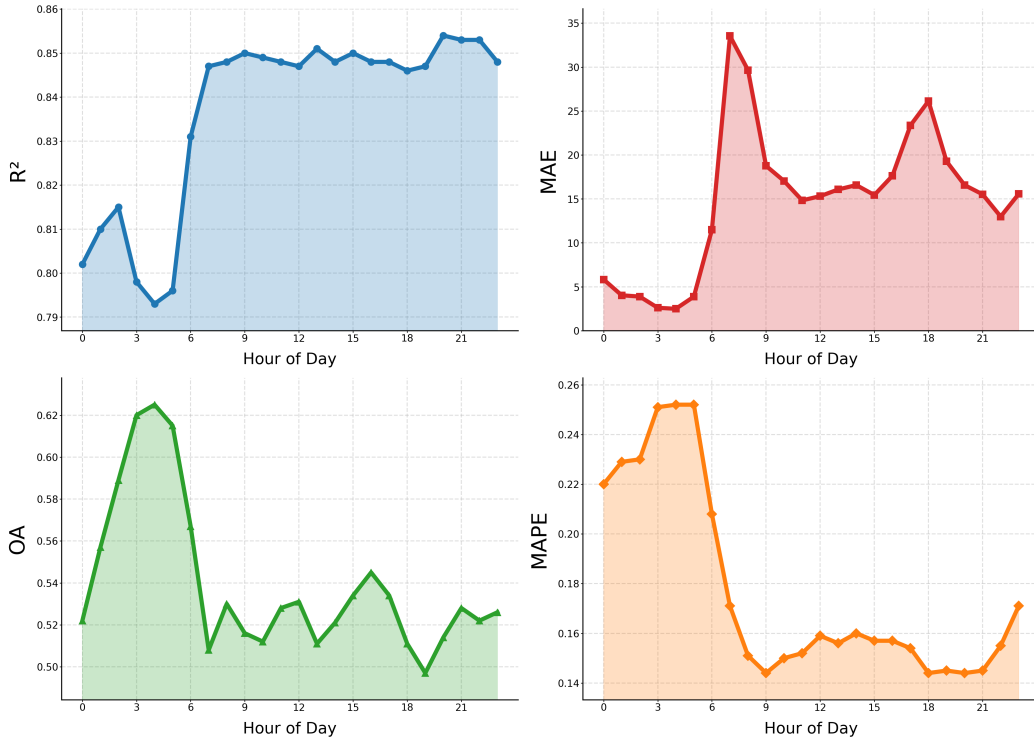


Figure 10: Hourly error analysis. Model performance across 24 hours evaluated using R^2 , OA, MAE, and MAPE.

560 four representative patterns emerge from the seven clusters, supported by
 561 VLM-generated semantics (Fig. 11 (b)): (i) *Cluster 0 (Morning-peak)*, char-
 562 acterized by 7:00–9:00 AM surges and terms like “commuting” and “residen-
 563 tial,” indicating outbound flow from housing areas; (ii) *Cluster 1 (Evening-*
 564 *peak)*, associated with “commercial district” and “night life,” reflecting post-
 565 work activity in CBDs; (iii) *Cluster 2 (Dual-peak)*, featuring bidirectional
 566 flows typical of mixed-use “transportation nodes”; and (iv) *Cluster 3 (Low-*
 567 *activity)*, corresponding to “ecological zones” with minimal human presence.
 568 These results demonstrate that integrating temporal signatures with textual
 569 semantics enables a coherent, explainable interpretation of urban dynamics.

570 5.3. Impact of Contrastive Learning

571 Fig. 12 visualizes the evolution of multimodal embedding distributions via
 572 t-SNE [74]. In the absence of contrastive learning (CL), embeddings remain
 573 scattered and modality-dependent, indicating weak semantic correspondence.

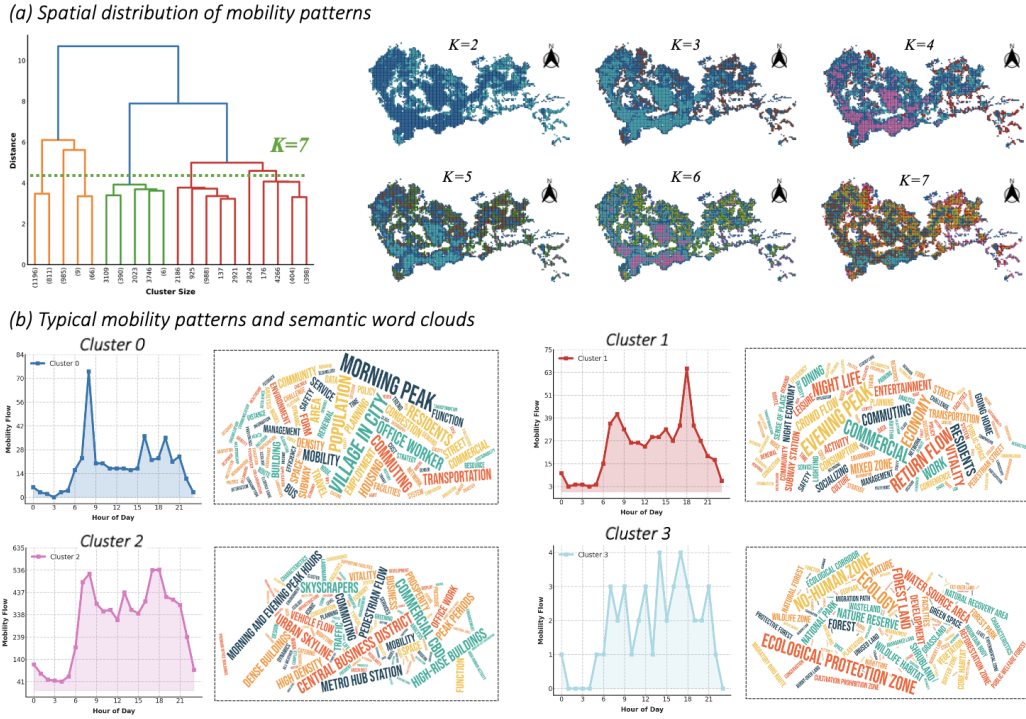


Figure 11: Spatiotemporal analysis of mobility patterns. (a) *Hierarchical evolution*: Clusters from $K=2$ to $K=7$ illustrate the structural refinement from a binary built/non-built dichotomy to a polycentric urban hierarchy. (b) *Semantic-temporal patterns*: Representative diurnal curves and VLM word clouds at $K=7$. The identified patterns include: *Cluster 0 (Morning-peak)*: residential zones driven by early commuting; *Cluster 1 (Evening-peak)*: commercial CBDs with nightlife activity; *Cluster 2 (Dual-peak)*: mixed-use areas and transport hubs; and *Cluster 3 (Low-activity)*: ecological protection zones. Word clouds link quantitative mobility rhythms to qualitative urban functions.

574 Conversely, the inclusion of CL drives the convergence of heterogeneous fea-
 575 tures into compact, semantically coherent clusters based on shared urban
 576 functions (e.g., residential vs. commercial). This explicit optimization of
 577 intra-class compactness and inter-modality similarity yields a significantly
 578 more discriminative embedding space than traditional fusion methods that
 579 rely solely on concatenation.

580 5.4. Impact of VLM Model Size

581 Table 8 presents a controlled comparison of six VLMs with varying param-
 582 eter scales. Performance generally correlates with model capacity; the largest

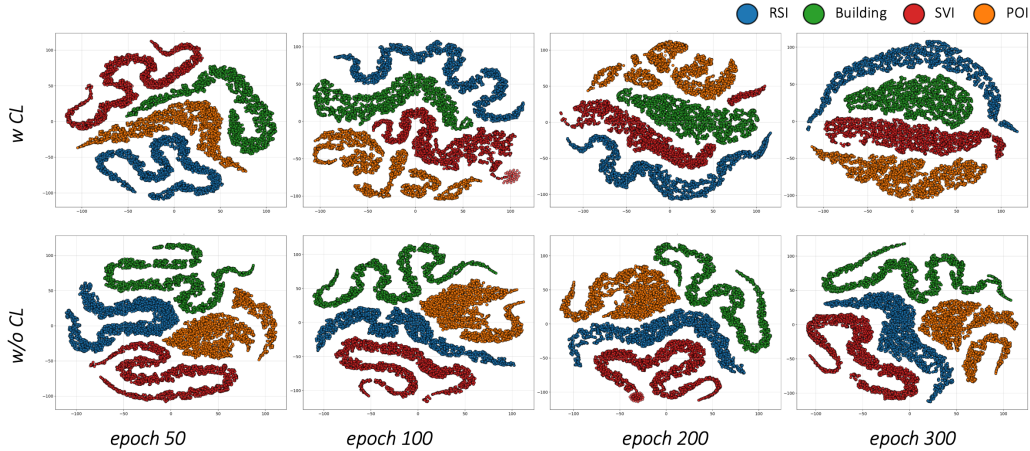


Figure 12: t-SNE visualization of multimodal embeddings with and without contrastive learning (CL). Colors denote modalities (RSI, Building, SVI, POI).

583 model, Qwen-VL-Max, achieves the best accuracy ($R^2 = 0.856$, MAE=
 584 13.738), outperforming the lightweight Gemini-Flash-Lite by a clear mar-
 585 gin. This suggests that larger models provide richer semantic reasoning for
 586 text-image alignment. However, gains diminish beyond 70B parameters,
 587 with mid-sized models (e.g., Qwen3-VL-32B, GPT-4o-mini) offering a fa-
 588 vorable balance between accuracy ($R^2 > 0.83$) and computational efficiency
 589 ($\approx 27s/image$). Consequently, models with approximately 30B parameters
 590 represent the most practical trade-off for large-scale urban semantic applica-
 591 tions.

Table 8: Performance comparison across different VLM model sizes. Time denotes the per-image semantic generation time. Best in **bold**, second best underlined.

Model	OA \uparrow	R^2 \uparrow	MAE \downarrow	MAPE \downarrow	Time (s/image)
Qwen-VL-Max	0.521	0.856	13.738	0.187	28.6
Qwen2.5-VL-72B	<u>0.501</u>	<u>0.844</u>	<u>14.103</u>	<u>0.209</u>	25.5
Qwen3-VL-32B	0.495	0.834	14.422	0.216	27.4
Qwen3-VL-8B	0.491	0.834	14.514	0.217	<u>20.6</u>
GPT-4o-mini	0.499	0.843	14.148	0.211	28.0
Gemini-Flash-Lite	0.499	0.838	14.285	0.214	18.0

592 5.5. Practical Applications

593 The proposed framework enables several practical applications in urban
594 management. For urban planning, the predicted 24-hour mobility profiles
595 allow planners to evaluate how land-use configurations shape daily activ-
596 ity patterns, informing zoning and development decisions without requiring
597 dynamic trajectory data [75]. For transit operations, the hourly temporal res-
598 olution supports demand-responsive scheduling, enabling agencies to adjust
599 service frequency and routes based on predicted ridership [76]. More broadly,
600 because the framework relies solely on static geospatial data—which is widely
601 available through open platforms (e.g., OpenStreetMap, Google Earth)—it
602 can be readily generalized to cities worldwide that lack proprietary mobil-
603 ity datasets, as evidenced by consistent performance across four cities in
604 the Greater Bay Area (Table 3), offering a scalable and privacy-preserving
605 solution for global urban mobility modeling [77].

606 6. Conclusion

607 This study proposes LMEMR, a framework that predicts hourly grid-level
608 mobility solely from static geospatial data by combining VLM-based semantic
609 enhancement, dual-level contrastive learning, and graph-attention spatiotem-
610 poral modeling. Experiments on four cities confirm consistent generalizabil-
611 ity; on Shenzhen, LMEMR achieves an 18.07% MAE reduction ($R^2=0.856$)
612 over the best baseline (CLIP), validating the potential of widely accessible
613 static data for scalable and privacy-friendly mobility inference. The main
614 contributions are: (i) a unified multimodal semantic graph framework for
615 hourly mobility prediction; (ii) a dual-level contrastive strategy aligning raw
616 and semantic features across modalities; (iii) an attention-based intermodal
617 fusion and spatiotemporal module; and (iv) extensive multi-city experiments
618 confirming superiority, interpretability, and generalizability.

619 Despite these advances, several limitations remain. (i) Static features
620 cannot capture event-driven perturbations (e.g., extreme weather, public
621 events) [78]; integrating weakly dynamic priors may mitigate this gap. (ii) All
622 four evaluated cities lie within the Greater Bay Area of China and share sim-
623 ilar urban morphologies and data infrastructures. Generalizability to cities
624 with fundamentally different forms—such as North American sprawl or in-
625 formal settlements in developing countries—requires further validation. (iii)
626 The 500 m grid resolution introduces scale sensitivity, and the Modifiable
627 Areal Unit Problem (MAUP) implies that spatial conclusions may vary with

628 the chosen zoning scheme. Future work should examine multi-scale effects
629 and adaptive spatial partitioning strategies.

630 **Acknowledgements**

631 This work was supported in part by the National Natural Science Foun-
632 dation of China under Grant 42401553, in part by the Natural Science
633 Foundation of Top Talent of Shenzhen Technology University under Grant
634 GDRC202415, in part by the Shenzhen Science and Technology Program un-
635 der Grants JCYJ20240813113300001 and 20231127180406001, and in part by
636 the Guangdong Basic and Applied Basic Research Foundation under Grant
637 2026A1515011875.

638 **References**

- 639 [1] E. Chen, Y. Liu, M. Yang, Revealing senior mobility patterns and activ-
640 ities in urban transit systems, *IEEE Transactions on Intelligent Trans-*
641 *portation Systems* 24 (2023) 11424–11437. [doi:10.1109/TITS.2023.](https://doi.org/10.1109/TITS.2023.3275389)
642 [3275389](https://doi.org/10.1109/TITS.2023.3275389).
- 643 [2] M. Chen, Q. Yuan, C. Yang, Y. Zhang, Decoding urban mobility: Appli-
644 cation of natural language processing and machine learning to activity
645 pattern recognition, prediction, and temporal transferability examina-
646 tion, *IEEE Transactions on Intelligent Transportation Systems* 25 (2024)
647 7151–7173. [doi:10.1109/TITS.2023.3339772](https://doi.org/10.1109/TITS.2023.3339772).
- 648 [3] L. G. Azolin, A. N. R. da Silva, N. Pinto, Incorporating public transport
649 in a methodology for assessing resilience in urban mobility, *Transporta-*
650 *tion Research Part D-transport and Environment* 85 (2020) 102386.
651 [doi:10.1016/j.trd.2020.102386](https://doi.org/10.1016/j.trd.2020.102386).
- 652 [4] E. Suryani, R. A. Hendrawan, P. F. E. Adipraja, R. Indraswari, System
653 dynamics simulation model for urban transportation planning: a case
654 study, *International Journal of Simulation Modelling* (2020). [doi:10.](https://doi.org/10.2507/ijstimm19-1-493)
655 [2507/ijstimm19-1-493](https://doi.org/10.2507/ijstimm19-1-493).
- 656 [5] K. Zhang, Q. Jin, K. Pelechris, T. Lappas, On the importance of
657 temporal dynamics in modeling urban activity, in: *Proceedings of the*
658 *2nd ACM SIGKDD International Workshop on Urban Computing*, 2013,
659 pp. 7:1–7:8. [doi:10.1145/2505821.2505825](https://doi.org/10.1145/2505821.2505825).

- 660 [6] K. Qin, Y. Xu, C. Kang, S. Sobolevsky, M. Kwan, Modeling spatio-
661 temporal evolution of urban crowd flows, *ISPRS Int. J. Geo Inf.* 8 (2019)
662 570. [doi:10.3390/ijgi8120570](https://doi.org/10.3390/ijgi8120570).
- 663 [7] C. Liu, L. Chen, Q. Yuan, H. Wu, W. Huang, Revealing dynamic spa-
664 tial structures of urban mobility networks and the underlying evolu-
665 tionary patterns, *ISPRS Int. J. Geo Inf.* 11 (2022) 237. [doi:10.3390/ijgi11040237](https://doi.org/10.3390/ijgi11040237).
- 667 [8] L. B. Santos, L. Carvalho, W. Seron, F. Coelho, E. Macau, M. G. Quiles,
668 A. Monteiro, How do urban mobility (geo)graph's topological properties
669 fill a map?, *Applied Network Science* 4 (2019) 1–14. [doi:10.1007/s41109-019-0211-7](https://doi.org/10.1007/s41109-019-0211-7).
- 671 [9] Z. Ye, Construction and application of urban mobility diagnosis model,
672 *Frontiers of Urban and Rural Planning* 1 (2023). [doi:10.1007/s44243-023-00016-9](https://doi.org/10.1007/s44243-023-00016-9).
- 674 [10] X. Shi, F. Lv, D. Seng, B. Xing, B. Chen, Visual exploration of mo-
675 bility dynamics based on multi-source mobility datasets and poi infor-
676 mation, *Journal of Visualization* 22 (2019) 1209–1223. [doi:10.1007/s12650-019-00594-1](https://doi.org/10.1007/s12650-019-00594-1).
- 678 [11] R. Jiang, X. Song, Z. Fan, T. Xia, Z. Wang, Q. Chen, Z. Cai,
679 R. Shibasaki, Transfer urban human mobility via POI embedding over
680 multiple cities, *ACM Transactions on Data Science* 2 (2021) 1–26.
681 [doi:10.1145/3416914](https://doi.org/10.1145/3416914).
- 682 [12] M. Pastorino, F. Gallo, A. D. Febbraro, G. Moser, N. Sacco, S. Serpico,
683 Multimodal fusion of mobility demand data and remote sensing imagery
684 for urban land-use and land-cover mapping, *Remote. Sens.* 14 (2022)
685 3370. [doi:10.3390/rs14143370](https://doi.org/10.3390/rs14143370).
- 686 [13] W. Tu, Z. Hu, L. Li, J. Cao, J. Jiang, Q. Li, Q. Li, Portraying urban
687 functional zones by coupling remote sensing imagery and human sensing
688 data, *Remote. Sens.* 10 (2018) 141. [doi:10.3390/rs10010141](https://doi.org/10.3390/rs10010141).
- 689 [14] M. Quintana, Y. Gu, X. Liang, Y. Hou, K. Ito, Y. Zhu, M. Abdelrahman,
690 F. Biljecki, Global urban visual perception varies across demographics
691 and personalities, *Nature Cities* 2 (11) (2025) 1092–1106. [doi:10.1038/s44284-025-00330-x](https://doi.org/10.1038/s44284-025-00330-x).
- 692

- 693 [15] C. Ye, F. Zhang, L. Mu, Y. Gao, Y. Liu, Urban function recognition
694 by integrating social media and street-level imagery, *Environment and*
695 *Planning B: Urban Analytics and City Science* 48 (2020) 1430–1444.
696 [doi:10.1177/2399808320935467](https://doi.org/10.1177/2399808320935467).
- 697 [16] T. Zhao, X. Liang, W. Tu, Z. Huang, F. Biljecki, Sensing urban sound-
698 scapes from street view imagery, *Comput. Environ. Urban Syst.* 99
699 (2023) 101915. [doi:10.1016/j.compenvurbsys.2022.101915](https://doi.org/10.1016/j.compenvurbsys.2022.101915).
- 700 [17] F. Zhang, L. Wu, D. Zhu, Y. Liu, Social sensing from street-level im-
701 agery: A case study in learning spatio-temporal urban mobility patterns,
702 *ISPRS journal of photogrammetry and remote sensing* 153 (2019) 48–58.
703 [doi:10.1016/j.isprsjprs.2019.04.017](https://doi.org/10.1016/j.isprsjprs.2019.04.017).
- 704 [18] Y. Zhang, P. Liu, F. Biljecki, Knowledge and topology: A two layer
705 spatially dependent graph neural networks to identify urban functions
706 with time-series street view image, *ISPRS Journal of Photogrammetry*
707 *and Remote Sensing* (2023). [doi:10.1016/j.isprsjprs.2023.03.008](https://doi.org/10.1016/j.isprsjprs.2023.03.008).
- 708 [19] Y. Liu, X. Zhang, J. Ding, Y. Xi, Y. Li, Knowledge-infused contrastive
709 learning for urban imagery-based socioeconomic prediction, *Proceed-*
710 *ings of the ACM Web Conference 2023* (2023). [doi:10.1145/3543507.](https://doi.org/10.1145/3543507.3583876)
711 [3583876](https://doi.org/10.1145/3543507.3583876).
- 712 [20] Z. Wang, H. Li, R. Rajagopal, Urban2vec: Incorporating street view im-
713 agery and pois for multi-modal urban neighborhood embedding, *ArXiv*
714 *abs/2001.11101* (2020). [doi:10.1609/AAAI.V34I01.5450](https://doi.org/10.1609/AAAI.V34I01.5450).
- 715 [21] M. Chen, Z. Li, H. Jia, X. Shao, J. Zhao, Q. Gao, M. Yang, Y. Yin,
716 MGRL4RE: A multi-graph representation learning approach for urban
717 region embedding, *ACM Transactions on Intelligent Systems and Tech-*
718 *nology* 16 (2025) 1–23. [doi:10.1145/3712698](https://doi.org/10.1145/3712698).
- 719 [22] X. Wang, J. Cao, T. Zhao, B. Zhang, G. Chen, Z. Li, H. Chen, W. Tu,
720 Q. Li, St-camba: A decoupled-free spatiotemporal graph fusion state
721 space model with linear complexity for efficient traffic forecasting, *In-*
722 *formation Fusion* (2025) 103495 [doi:10.2139/ssrn.5206097](https://doi.org/10.2139/ssrn.5206097).
- 723 [23] X. Wang, T. Zhao, W. Tu, B. Zhang, G. Chen, J. Cao, Sat2flow: A
724 structure-aware diffusion framework for human flow generation from

- 725 satellite imagery, in: Proceedings of the AAAI Conference on Artificial
726 Intelligence, Vol. 40, 2026, pp. 15886–15894.
- 727 [24] Y. Kang, Short-term passenger flow prediction in urban rail transit
728 based on hybrid deep learning models, Highlights in Science, Engineer-
729 ing and Technology (2024). doi:10.54097/nqxs4628.
- 730 [25] Y. Sun, Prediction of short-term passenger flow in the metro station with
731 cnn-lstm model, 2023 IEEE 3rd International Conference on Electronic
732 Technology, Communication and Information (ICETCI) (2023) 1218–
733 1222doi:10.1109/icetci57876.2023.10176978.
- 734 [26] L. Liu, Z. Qiu, G. Li, Q. Wang, W. Ouyang, L. Lin, Contextualized
735 spatial-temporal network for taxi origin-destination demand prediction,
736 IEEE Transactions on Intelligent Transportation Systems 20 (10) (2019)
737 3875–3887. doi:10.1109/tits.2019.2915525.
- 738 [27] J. Zhang, H. Che, F. Chen, W. Ma, Z. He, Short-term origin-destination
739 demand prediction in urban rail transit systems: A channel-wise at-
740 tentive split-convolutional neural network method, Transportation Re-
741 search Part C: Emerging Technologies 124 (2021) 102928. doi:10.1016/
742 j.trc.2020.102928.
- 743 [28] I. D. P. Arenas, H. Alatrística-Salas, M. N. del Prado Cortez, Discovery
744 of urban mobility patterns, Advances in Data Science and Information
745 Engineering (2021). doi:10.1007/978-3-030-71704-9_33.
- 746 [29] J. Cheng, K. Li, Y. Liang, L. Sun, J. Yan, Y. Wu, Rethinking urban mo-
747 bility prediction: A multivariate time series forecasting approach, IEEE
748 Transactions on Intelligent Transportation Systems 26 (2025) 2543–
749 2557. doi:10.1109/tits.2024.3498054.
- 750 [30] Q. Pan, Y. Chen, G. Shen, Y. Yang, X. Kong, Spatio-temporal knowl-
751 edge embedding via circular correlation: insights into functional urban
752 area travel pattern mining, Neural Comput. Appl. 36 (2024) 19075–
753 19095. doi:10.1007/s00521-024-10167-5.
- 754 [31] T. Li, Q. Feng, B. Niu, B. Chen, F. Yan, J. Gong, J. Liu, Mapping urban
755 villages based on point-of-interest data and a deep learning approach,
756 Cities (2025). doi:10.1016/j.cities.2024.105549.

- 757 [32] R. Chen, S. Jiang, W. Huang, Semob: Semantic synthesis for dynamic
758 urban mobility prediction, in: Proceedings of the 2025 Conference on
759 Empirical Methods in Natural Language Processing, 2025, pp. 15346–
760 15366.
- 761 [33] Y. Xu, S. Gao, Q. Huang, A. Göçmen, Q. Zhu, F. Zhang, Predicting
762 human mobility flows in cities using deep learning on satellite imagery,
763 Nature Communications 16 (1) (2025) 10372.
- 764 [34] J. Cao, X. Wang, G. Chen, W. Tu, X. Shen, T. Zhao, J. Chen, Q. Li,
765 Disentangling the hourly dynamics of mixed urban function: A multi-
766 modal fusion perspective using dynamic graphs, Information Fusion 117
767 (2025) 102832. doi:10.1016/j.inffus.2024.102832.
- 768 [35] F. Zhang, B. Zhou, L. Liu, Y. Liu, H. H. Fung, H. Lin, C. Ratti,
769 Measuring human perceptions of a large-scale urban region using ma-
770 chine learning, Landscape and Urban Planning 180 (2018) 148–160.
771 doi:10.1016/j.landurbplan.2018.08.020.
- 772 [36] M. Wu, Q. Huang, S. Gao, Z. Zhang, Mixed land use measurement and
773 mapping with street view images and spatial context-aware prompts via
774 zero-shot multimodal learning, International Journal of Applied Earth
775 Observation and Geoinformation 125 (2023) 103591. doi:10.1016/j.
776 jag.2023.103591.
- 777 [37] C. Su, X. Hu, Q. Meng, L. Zhang, W. Shi, M. Zhao, A multimodal
778 fusion framework for urban scene understanding and functional identi-
779 fication using geospatial data, International Journal of Applied Earth
780 Observation and Geoinformation 127 (2024) 103696. doi:10.1016/j.
781 jag.2024.103696.
- 782 [38] Y. Zhao, K. Zhao, Z. Tang, X. Lu, Y. Zhang, Y. Du, GraphJCL:
783 A dual-perspective graph-based framework for urban region represen-
784 tation via joint contrastive learning, in: Proceedings of the Euro-
785 pean Conference on Machine Learning and Principles and Practice of
786 Knowledge Discovery in Databases (ECML PKDD), Springer, 2025.
787 doi:10.1007/978-3-032-06066-2_3.
- 788 [39] Y. Chen, W. Huang, K. Zhao, Y. Jiang, G. Cong, Self-supervised repre-
789 sentation learning for geospatial objects: A survey, Information Fusion
790 (2025) 103265doi:10.2139/ssrn.5173896.

- 791 [40] Y. Li, W. Huang, G. Cong, H. Wang, Z. Wang, Urban region representa-
792 tion learning with openstreetmap building footprints, in: Proceedings of
793 the 29th ACM SIGKDD Conference on Knowledge Discovery and Data
794 Mining, 2023, pp. 1363–1373. doi:10.1145/3580305.3599538.
- 795 [41] Y. Tao, W. Liu, J. Chen, J. Gao, R. Li, X. Wang, Y. Zhang, J. Ren,
796 S. Yin, X. Zhu, et al., A graph-based multimodal data fusion framework
797 for identifying urban functional zone, International Journal of Applied
798 Earth Observation and Geoinformation 136 (2025) 104353. doi:10.
799 1016/j.jag.2024.104353.
- 800 [42] R. Yang, Y. Zhong, Y. Su, Self-supervised joint representation learning
801 for urban land-use classification with multisource geographic data, IEEE
802 Transactions on Geoscience and Remote Sensing 63 (2025) 1–21.
- 803 [43] W. Xu, J. Wang, Y. Wu, Multi-dimension geospatial feature learning for
804 urban region function recognition, in: IGARSS 2022-2022 IEEE Inter-
805 national Geoscience and Remote Sensing Symposium, IEEE, 2022, pp.
806 5832–5835. doi:10.1109/igarss46834.2022.9884450.
- 807 [44] X. Sun, J. Gao, Y. Yuan, Alignment and fusion using distinct sensor
808 data for multimodal aerial scene classification, IEEE Transactions on
809 Geoscience and Remote Sensing 62 (2024) 1–11.
- 810 [45] Y. Feng, J. Jin, Y. Yin, C. Song, X. Wang, Mcft: Multimodal contrastive
811 fusion transformer for classification of hyperspectral image and lidar
812 data, IEEE Transactions on Geoscience and Remote Sensing 62 (2024)
813 1–17.
- 814 [46] J. Cao, X. Wang, J. Chen, W. Tu, Z. Li, X. Yang, T. Zhao, Q. Li, Ur-
815 ban representation learning for fine-grained economic mapping: A semi-
816 supervised graph-based approach, ISPRS Journal of Photogrammetry
817 and Remote Sensing 226 (2025) 317–331. doi:10.1016/j.isprsjprs.
818 2025.05.007.
- 819 [47] D. J. Mühlematter, L. Che, Y. Hong, M. Raubal, N. Wiedemann, Urban-
820 fusion: Stochastic multimodal fusion for contrastive learning of robust
821 spatial representations, arXiv preprint arXiv:2510.13774 (2025).

- 822 [48] J. Cao, J. Chen, X. Wang, W. Huang, D. Chen, T. Zhao, W. Tu, Q. Li,
823 Urbanmmcl: Urban region representations via multi-modal and multi-
824 graph self-supervised contrastive learning, *ISPRS Journal of Photogram-*
825 *metry and Remote Sensing* 232 (2026) 75–93.
- 826 [49] Y. Lu, L. Zhang, J. Liu, Q. Tian, Constructing concept lexica with small
827 semantic gaps, *IEEE Transactions on Multimedia* 12 (2010) 288–299.
828 [doi:10.1109/tmm.2010.2046292](https://doi.org/10.1109/tmm.2010.2046292).
- 829 [50] J. Zhang, T. Li, X. Lu, Z. Cheng, Semantic classification of high-
830 resolution remote-sensing images based on mid-level features, *IEEE*
831 *Journal of Selected Topics in Applied Earth Observations and Remote*
832 *Sensing* 9 (2016) 2343–2353. [doi:10.1109/jstars.2016.2536943](https://doi.org/10.1109/jstars.2016.2536943).
- 833 [51] J. Wang, C. Gao, M.-L. Wang, Y. Zhang, Identification of urban func-
834 tional areas and urban spatial structure analysis by fusing multi-source
835 data features: A case study of zhengzhou, china, *Sustainability* (2023).
836 [doi:10.3390/su15086505](https://doi.org/10.3390/su15086505).
- 837 [52] W. Huang, L. zhen Cui, M. Chen, D. Zhang, Y. Yao, Estimating urban
838 functional distributions with semantics preserved poi embedding, *Inter-*
839 *national Journal of Geographical Information Science* 36 (2022) 1905–
840 1930. [doi:10.1080/13658816.2022.2040510](https://doi.org/10.1080/13658816.2022.2040510).
- 841 [53] Q. Qin, S. Xu, M. Du, S. Li, Identifying urban functional zones by
842 capturing multi-spatial distribution patterns of points of interest, *Inter-*
843 *national Journal of Digital Earth* 15 (2022) 2468–2494. [doi:10.1080/](https://doi.org/10.1080/17538947.2022.2160841)
844 [17538947.2022.2160841](https://doi.org/10.1080/17538947.2022.2160841).
- 845 [54] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal,
846 G. Sastry, A. Askell, P. Mishkin, J. Clark, et al., Learning transfer-
847 able visual models from natural language supervision, in: *International*
848 *conference on machine learning*, PMLR, 2021, pp. 8748–8763.
- 849 [55] J. Li, D. Li, C. Xiong, S. Hoi, BLIP: Bootstrapping language-image pre-
850 training for unified vision-language understanding and generation, in:
851 *International conference on machine learning*, PMLR, 2022, pp. 12888–
852 12900.
- 853 [56] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen,
854 X. Liu, J. Wang, W. Ge, et al., Qwen2-VL: Enhancing vision-language

- 855 model’s perception of the world at any resolution, arXiv preprint
856 arXiv:2409.12191 (2024).
- 857 [57] G. Mai, W. Huang, J. Sun, S. Song, D. R. Mishra, N. Liu, S. Gao,
858 T. Liu, G. Cong, Y. Hu, C. Cundy, Z. Li, R. Zhu, N. Lao, On the op-
859 portunities and challenges of foundation models for geoai (vision paper),
860 ACM Transactions on Spatial Algorithms and Systems 10 (2024) 1–46.
861 [doi:10.1145/3653070](https://doi.org/10.1145/3653070).
- 862 [58] R. Mushkani, Do vision-language models see urban scenes as people
863 do? an urban perception benchmark, arXiv preprint arXiv:2509.14574
864 (2025).
- 865 [59] T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework
866 for contrastive learning of visual representations, in: International con-
867 ference on machine learning, PMLR, 2020, pp. 1597–1607. [arXiv:
868 2002.05709](https://arxiv.org/abs/2002.05709).
- 869 [60] C. Liu, Q. Xu, H. Miao, S. Yang, L. Zhang, C. Long, Z. Li, R. Zhao,
870 Timecma: Towards llm-empowered multivariate time series forecasting
871 via cross-modality alignment, in: Proceedings of the AAAI Conference
872 on Artificial Intelligence, Vol. 39, 2025, pp. 18780–18788. [doi:10.1609/
873 aaai.v39i18.34067](https://doi.org/10.1609/aaai.v39i18.34067).
- 874 [61] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image
875 recognition, in: Proceedings of the IEEE conference on computer vision
876 and pattern recognition, 2016, pp. 770–778. [doi:10.1109/cvpr.2016.
877 90](https://doi.org/10.1109/cvpr.2016.90).
- 878 [62] Y. Zhang, M. Li, D. Long, X. Zhang, H. Lin, B. Yang, P. Xie, A. Yang,
879 D. Liu, J. Lin, et al., Qwen3 embedding: Advancing text embedding and
880 reranking through foundation models, arXiv preprint arXiv:2506.05176
881 (2025). [arXiv:2506.05176](https://arxiv.org/abs/2506.05176).
- 882 [63] Y. Yan, H. Wen, S. Zhong, W. Chen, H. Chen, Q. Wen, R. Zim-
883 mermann, Y. Liang, UrbanCLIP: Learning text-enhanced urban re-
884 gion profiling with contrastive language-image pretraining from the
885 web, Proceedings of the ACM Web Conference 2024 (2023). [doi:
886 10.1145/3589334.3645378](https://doi.org/10.1145/3589334.3645378).

- 887 [64] T. Zhao, Z. Huang, W. Tu, F. Biljecki, L. Chen, Developing a multiview
888 spatiotemporal model based on deep graph neural networks to predict
889 the travel demand by bus, *International Journal of Geographical In-*
890 *formation Science* 37 (7) (2023) 1555–1581. [doi:10.1080/13658816.](https://doi.org/10.1080/13658816.2023.2203218)
891 [2023.2203218](https://doi.org/10.1080/13658816.2023.2203218).
- 892 [65] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, Y. Bengio,
893 Graph attention networks, *arXiv preprint arXiv:1710.10903* (2017).
- 894 [66] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,
895 Ł. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural*
896 *information processing systems* 30 (2017).
- 897 [67] D. Maulud, A. M. Abdulazeez, A review on linear regression compre-
898 hensive in machine learning, *Journal of applied science and technology*
899 *trends* 1 (2) (2020) 140–147. [doi:10.38094/jastt1457](https://doi.org/10.38094/jastt1457).
- 900 [68] A. Liaw, M. Wiener, et al., Classification and regression by randomfor-
901 est, *R news* 2 (3) (2002) 18–22.
- 902 [69] Z. Li, F. Liu, W. Yang, S. Peng, J. Zhou, A survey of convolutional
903 neural networks: analysis, applications, and prospects, *IEEE transac-*
904 *tions on neural networks and learning systems* 33 (12) (2021) 6999–7019.
905 [doi:10.1109/tnnls.2021.3084827](https://doi.org/10.1109/tnnls.2021.3084827).
- 906 [70] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural com-*
907 *putation* 9 (8) (1997) 1735–1780. [doi:10.1162/neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735).
- 908 [71] B. Yu, H. Yin, Z. Zhu, Spatio-temporal graph convolutional net-
909 works: A deep learning framework for traffic forecasting, *arXiv preprint*
910 *arXiv:1709.04875* (2017). [doi:10.24963/ijcai.2018/505](https://doi.org/10.24963/ijcai.2018/505).
- 911 [72] S. Yang, S. Qian, Understanding and predicting travel time with spatio-
912 temporal features of network traffic flow, weather and incidents, *IEEE*
913 *Intelligent Transportation Systems Magazine* 11 (2019) 12–28. [doi:](https://doi.org/10.1109/mits.2019.2919615)
914 [10.1109/mits.2019.2919615](https://doi.org/10.1109/mits.2019.2919615).
- 915 [73] Y. Chen, Z. Zhang, T. Liang, Assessing urban travel patterns: An analy-
916 sis of traffic analysis zone-based mobility patterns, *Sustainability* (2019).
917 [doi:10.3390/su11195452](https://doi.org/10.3390/su11195452).

- 918 [74] A. Chatzimparmpas, R. M. Martins, A. Kerren, T-viSNE: Interactive
919 assessment and interpretation of t-SNE projections, *IEEE transactions*
920 *on visualization and computer graphics* 26 (8) (2020) 2696–2714. doi:
921 [10.1109/tvcg.2020.2986996](https://doi.org/10.1109/tvcg.2020.2986996).
- 922 [75] J. Zhang, Y. Zheng, D. Qi, R. Li, X. Yi, Predicting citywide crowd flows
923 using deep spatio-temporal residual networks, *Artificial Intelligence* 259
924 (2018) 147–166. doi:[10.1016/j.artint.2018.03.002](https://doi.org/10.1016/j.artint.2018.03.002).
- 925 [76] W. Wang, F. Zong, B. Yao, A proactive real-time control strategy
926 based on data-driven transit demand prediction, *IEEE Transactions*
927 *on Intelligent Transportation Systems* 22 (4) (2021) 2404–2416. doi:
928 [10.1109/TITS.2020.3028415](https://doi.org/10.1109/TITS.2020.3028415).
- 929 [77] C. Rong, X. Zhang, Y. Xi, H. Sui, J. Ding, Y. Li, Satellites reveal mo-
930 bility: A commuting origin-destination flow generator for global cities,
931 *arXiv preprint arXiv:2505.15870* (2025).
- 932 [78] T. Fontes, R. Correia, J. Ribeiro, J. L. Borges, A deep learning ap-
933 proach for predicting bus passenger demand based on weather condi-
934 tions, *Transport and Telecommunication* 21 (4) (2020) 255–264. doi:
935 [10.2478/ttj-2020-0020](https://doi.org/10.2478/ttj-2020-0020).