

UrbanMMCL: Urban Region Representations via Multi-Modal and Multi-Graph Self-Supervised Contrastive Learning

Jinzhou Cao^a, Jiashi Chen^a, Xiangxu Wang^a, Weiming Huang^c, Dongsheng Chen^d, Tianhong Zhao^{a,*}, Wei Tu^b, Qingquan Li^b

^a*Shenzhen Technology University, School of Artificial Intelligence, Shenzhen, Guangdong, China*

^b*Shenzhen University, Department of Urban Informatics & Guangdong Key Laboratory of Urban Informatics & Guangdong-Hong Kong-Macau Joint Laboratory for Smart Cities, Shenzhen, Guangdong, China*

^c*School of Geography, University of Leeds, Leeds, UK*

^d*Technical University of Munich, Chair of Cartography and Visual Analytics, Munich, Germany*

Abstract

Urban region representation learning has emerged as a fundamental approach for diverse urban analytics tasks, where each neighborhood is encoded as a dense embedding vector for effective downstream applications. However, existing approaches suffer from insufficient multi-modal alignment and inadequate spatial relationship modeling, limiting their representation quality and generalizability. To address these challenges, we propose UrbanMMCL, a novel self-supervised framework that integrates multi-modal multi-view contrastive pre-training with unified fine-tuning for comprehensive urban representation learning. UrbanMMCL employs a dual-stage architecture. First, cross-modal contrastive learning aligns diverse data modalities

*Corresponding author

Email address: zhaotianhong@sztu.edu.cn (Tianhong Zhao)

including remote sensing imagery, street view imagery, location encodings, and Vision-Language Model (VLM)-generated textual descriptions. Second, multi-view adaptive graph contrastive learning captures complex spatial relationships across human mobility, functional similarity, and geographic distance perspectives. The framework then integrates the learned representations through a dedicated fusion mechanism for effective adaptation to downstream tasks. Comprehensive experiments demonstrate that UrbanMMCL consistently outperforms state-of-the-art methods across pollutant emission prediction, population density estimation, and land use classification with minimal fine-tuning requirements, thereby advancing foundation model development for diverse Geo-AI applications.

Keywords: Urban Region Representation Learning, Contrastive Learning, Graph Learning, Multimodal Fusion, Urban Foundation Model

1. Introduction

Urban region representation learning extracts compact features from heterogeneous data to capture spatial, social, and economic characteristics essential for urban tasks like pollution prediction (He and Huang, 2025), socioeconomic estimation (Cao et al., 2025b), and land-use classification (Cao et al., 2025a). These applications provide valuable contributions to urban planning and environmental management in increasingly complex urban environments driven by the rapid urbanization process.

Conventional region representation learning approaches encounter substantial difficulties in efficiently incorporating multi-source multi-modal data to generate unified representations (Zhang et al., 2025). Urban regions are in-

herently heterogeneous, characterized by diverse physical attributes, dynamic socioeconomic activities, and complex interregional interactions. Thus, urban region representations require sophisticated modeling to capture this multidimensional nature (Wang et al., 2024, 2025).

Fusing multi-perspective visual data has emerged as a promising approach to comprehensively reveal heterogeneous urban characteristics. Remote sensing images (RSIs) provide foundational macroscopic perspectives with extensive coverage (Bai et al., 2023; Zhou et al., 2021), capturing urban morphology and land use patterns (Bai et al., 2025), while street view images (SVIs) offer complementary micro-level details by documenting street environments and building facades (Zhao et al., 2025; Zhang et al., 2019, 2024b). Nevertheless, the fusion of multi-perspective visual data presents unique challenges (Zou et al., 2025), as traditional approaches have treated these data independently or employed simplistic fusion mechanisms, failing to preserve complementary information (Gao et al., 2020).

Despite detailed physical characteristics, visual data alone lack semantic depth for complete regional description. While existing representation learning approaches rely on Point-of-Interest (POI) data for textual semantics (Li et al., 2023a), POI data’s spatial sparsity and uneven distribution frequently result in inconsistent representation quality (Qin et al., 2025). Vision-Language Models (VLMs) offer promising alternatives by generating comprehensive semantic descriptions from visual content (Huang et al., 2024), transforming RSIs and SVIs into rich textual annotations such as ‘high-density residential area’ and ‘busy transportation hub with commercial activities’. However, current methodologies underutilise these descriptions,

37 treating them as rudimentary labels rather than leveraging their semantic
38 intricacy for profound visual-textual alignment (Liu et al., 2024).

39 The modeling inherent spatial relationships between urban regions is
40 imperative for the effective representation learning. Urban regions exhibit
41 multi-faceted spatial interdependencies characterised by geographical adja-
42 cency, mobility patterns, and functional similarity (Wang et al., 2024). While
43 recent multi-view frameworks incorporate these dependencies (Li et al., 2019;
44 Wu et al., 2022), they typically process views independently or use simple
45 aggregation strategies (Zhang et al., 2020; Chan and Ren, 2023), missing syn-
46 ergistic information across relational perspectives. The utilization of graph
47 contrastive learning in urban spatial modeling represents a potentially fruit-
48 ful yet underexplored research avenue (Zhang et al., 2023d; Liu et al., 2025).

49 To address these challenges, we propose **UrbanMMCL**, a **Urban Multi-**
50 **Modal and Multi-View dual Contrastive Learning** framework that estab-
51 lishes a self-supervised pre-training and fine-tuning paradigm for compre-
52 hensive region representation learning. **Pre-training Stage** consists of two
53 synergistic components: (1) **multi-modal vision-language contrastive**
54 **learning** that aligns RSI, SVI, location encodings, and semantic textual
55 descriptions through specialized encoders and multi-level contrastive objec-
56 tives; (2) **adaptive multi-view graph contrastive learning** that models
57 complex spatial relationships through dynamic graph structure optimization
58 across multiple relational views. This stage learns generalizable urban repre-
59 sentations from unlabeled multi-modal data without requiring task-specific
60 annotations. **Fine-tuning Stage** integrates the pre-trained multi-modal
61 and multi-view representations through dedicated fusion mechanisms, en-

62 abling effective knowledge transfer to diverse downstream urban analytics
63 tasks including population estimation, pollutant emission monitoring, and
64 land use classification with minimal labeled data requirements.

65 Our key innovations are fourfold:

- 66 1. A systematic dual-stage framework that simultaneously addresses multi-
67 modal data and multi-view relationships, overcoming prior works’ single-
68 focus limitation in urban representation learning.
- 69 2. An comprehensive multi-modal alignment mechanism that unifies RSI-
70 SVI-Location-Text data through triple contrastive learning, establish-
71 ing deep semantic alignment while preserving semantic richness and
72 spatial context.
- 73 3. Adaptive multi-view spatial modeling that captures complex urban de-
74 pendencies (proximity, mobility, demographic similarity) through dy-
75 namic graph structure learning, enabling effective integration of multi-
76 ple relational perspectives without requiring predefined graph topolo-
77 gies.
- 78 4. A domain-specific self-supervised pre-training paradigm with superior
79 transferability across diverse urban analytics tasks, providing extensive
80 analysis of how different modalities, fusion strategies, and training ap-
81 proaches contribute to representation quality in resource-constrained
82 deployment scenarios.

83 Section 2 reviews related work on multimodal contrastive learning, graph
84 contrastive learning, and urban representation learning. Section 3 details the
85 UrbanMMCL framework. Section 4 presents experiments and evaluations.
86 Section 5 analyzes model components, training paradigms, and limitations.

87 Section 6 concludes the study.

88 2. Related works

89 2.1. Multimodal contrastive learning

90 Self-supervised learning (SSL) has emerged as a powerful paradigm that
91 leverages unlabeled data to learn generalizable representations, eliminating
92 the need for costly manual annotations. Among SSL approaches, contrastive
93 learning stands out as a particularly effective technique that learns repre-
94 sentations by maximizing similarity between positive pairs while minimiz-
95 ing similarity with negative samples (Dai et al., 2025; Zhang et al., 2023c).
96 Methods such as InstDis (Wu et al., 2018), SimCLR (Chen et al., 2020), and
97 MoCo series (He et al., 2020; Chen et al., 2021) have proven to be effective
98 in learning robust representations from unlabeled data.

99 Multimodal contrastive learning extends this paradigm by integrating
100 information from different data modalities to create unified representations
101 that capture complementary cross-modal correspondences (Wang et al., 2025;
102 Yong and Zhou, 2024). Vision-language contrastive learning represents a
103 particularly promising approach, combining rich spatial information from
104 imagery with semantic descriptions (Bao et al., 2022). CLIP (Radford et al.,
105 2021b) demonstrates the power of joint image-text representations through
106 contrastive training, enabling enhanced cross-modal understanding. Similar
107 approaches such as ALIGN (Jia et al., 2021) have expanded to billion-level
108 image-text pairs.

109 In urban analytics, multimodal approaches are particularly crucial due
110 to the inherently complex nature of urban environments, which generate di-

verse data types including RSIs, SVIs, POIs, and textual descriptions (Zhou et al., 2023b; Shen et al., 2023). Recent works have explored this direction in geospatial domains (Weng et al., 2025). GeoCLIP (Cepeda et al., 2023) applies contrastive learning for image-based geolocalization, while SatCLIP (Klemmer et al., 2025) extends CLIP to RSIs, learning representations that bridge RSIs with natural language descriptions. UrbanCLIP (Huang et al., 2024; Yan et al., 2024) specifically targets urban region understanding by integrating satellite imagery or street-view images with textual descriptions, and other works (Liu et al., 2023) have explored vision-language modeling and knowledge-infused contrastive frameworks for enhanced geographic understanding.

However, multimodal contrastive learning for urban region representation remains underexplored. Existing methods typically focus on single visual modalities with limited integration and lack effective adaptation of vision-language models for urban contexts. They treat geographical coordinates as auxiliary features rather than fundamental organizing principles for multimodal alignment. These highlight the need for specialized frameworks tailored to urban representation requirements.

2.2. Graph contrastive learning

Graph Neural Networks (GNNs) have revolutionized urban analysis by modeling urban regions as graph-structured data (Khoshraftar and An, 2024; Cao et al., 2025c). However, most GNN models rely on supervised training requiring substantial labeled data (Ju et al., 2024), which may be unavailable in many urban scenarios. To address these limitations, self-supervised graph contrastive learning (GCL) has emerged as a promising alternative that can

136 learn meaningful representations without labeled supervision.

137 GCL integrates both structural and attribute information by maximizing
138 agreement between disparate versions of the same graph while contrasting
139 with negative samples through the implementation of sophisticated architec-
140 tures and augmentation strategies (Wu et al., 2023; Sun et al., 2020a). Two
141 primary paradigms have emerged: global-local methods such as Deep Graph
142 Infomax (DGI) (Veličković et al., 2018), MVGRL (Hassani and Khasahmadi,
143 2020), and InfoGraph (Sun et al., 2020b) that contrast node-level with graph-
144 level representations, and local-local approaches such as GRACE (Zhu et al.,
145 2020), GraphCL (You et al., 2020) with its variants (You et al., 2021; Suresh
146 et al., 2021), and GCA (Zhu et al., 2021) that maximize agreement between
147 node embeddings across augmented graph views.

148 Multi-view graph contrastive learning integrates multiple graph perspec-
149 tives to capture diverse urban relationships (He et al., 2025). Urban appli-
150 cations have constructed complementary views including POI co-occurrence
151 networks (Huang et al., 2023; Zhang et al., 2023a), trajectory-based mobility
152 graphs (Zhang et al., 2024a), and spatial adjacency graphs (Luo et al., 2022).
153 However, contemporary multi-view GCL methods encounter critical limita-
154 tions. Zhang et al. (Zhang et al., 2023d) propose a multi-view framework
155 using triplet loss, but their node-level approach with static view construc-
156 tion misses subgraph-level patterns that characterize urban functional areas.
157 Their method relies on fixed topologies and simple augmentation strategies
158 that cannot adapt to dynamic urban spatial relationships. This highlights the
159 need for sophisticated multi-view GCL frameworks that integrate heteroge-
160 neous urban data through adaptive augmentation strategies while preserving

161 semantic coherence of urban functional regions.

162 *2.3. Urban representation learning*

163 Urban region representation learning aims to generate low-dimensional
164 embeddings that reflect urban regional attributes and interregional relation-
165 ships while preserving spatial and semantic structures. A effective learning
166 requires mining intrinsic correlations among heterogeneous data sources, in-
167 cluding geographic topology, urban visual imagery and human mobility (Wang
168 et al., 2026; Guan et al., 2024). This paradigm enables effective analysis
169 across diverse urban applications from sociodemographic prediction to land
170 use classification.

171 Early methods primarily relied on single modalities such as POI features
172 (Zhai et al., 2019; Sun et al., 2021), human mobility patterns (Zhou and
173 Huang, 2018), or visual imagery (Li et al., 2023b). While achieving task-
174 specific success, single-modal approaches fail to capture multi-dimensional
175 urban characteristics (Zou et al., 2025). Recent advancements focus on multi-
176 modal fusion, integrating spatial, visual and textual data for comprehensive
177 regional characterization (Zou et al., 2025). Representative works include
178 RegionEncoder (Jenkins et al., 2019) for joint encoding of POIs, mobility
179 flows, and RSIs, and Urban2Vec (Wang et al., 2020) combining SVIs with
180 POI descriptions. However, existing approaches predominantly rely on sim-
181 ple concatenation and attention mechanisms, lacking sophisticated semantic
182 alignment and hierarchical adaptive fusion strategies.

183 The field has evolved from traditional techniques including matrix factor-
184 ization (Belkin and Niyogi, 2001) and network embedding methods such as
185 DeepWalk and Node2Vec (Perozzi et al., 2014; Grover and Leskovec, 2016)

186 to GNNs (Xu et al., 2022). Traditional approaches heavily depends on task-
 187 customized supervised paradigms (Gao et al., 2020) targeting specific objec-
 188 tives like poverty assessment (Jean et al., 2016; Yeh et al., 2020) and urban
 189 function classification (Cao et al., 2020). Recent advances embrace self-
 190 supervised learning (Chen et al., 2025), with notable approaches including
 191 ReCP (Li et al., 2024) and GraphST (Zhang et al., 2023b), and multiview
 192 graph learning such as MVURE (Zhang et al., 2020) and CGAP (Xu and
 193 Zhou, 2024). Notwithstanding the advances that have been made, the de-
 194 sign of self-supervised pre-training tasks for universal urban representation
 195 remains a critical challenge.

196 3. Methodology

197 3.1. Preliminaries

198 We formalize the urban region representation problem through the fol-
 199 lowing key components.

200 **Definition 1. Urban Spatial Partitioning.** Given a city divided into
 201 N non-overlapping grid regions $\mathcal{R} = \{r_i\}_{i=1}^N$, each region r_i is associated with
 202 multi-modal urban data.

203 **Definition 2. Remote Sensing Imagery.** Remote sensing imagery
 204 \mathcal{I}^{RS} captures aerial views of the earth’s surface, providing insights into build-
 205 ing distributions and land use patterns. For each region r_i , an orthorectified
 206 image patch $\mathcal{I}_i^{\text{RS}} \in \mathbb{R}^{h \times w}$ is used, where h and w are the dimensions of the
 207 grid.

208 **Definition 3. Street View Imagery.** Street view imagery \mathcal{I}^{SV} provides
 209 ground-level views of urban areas. For each region r_i , multi-directional street

view images are collected as:

$$\mathcal{I}_i^{\text{SV}} = \bigcup_{j=1}^n I_{i,j}^{\theta}, \quad (1)$$

where $\{s_{i,j}\}_{j=1}^n$ represents the uniformly distributed n sampling points along the road network within the region r_i , and $I_{i,j}^{\theta}$ denotes the image captured at the point $s_{i,j}$. This collection approach ensures comprehensive coverage of urban streetscapes from multiple viewpoints.

Definition 4. VLM-Enhanced Textual Description. Textual descriptions of a region r_i include satellite-derived text $\mathcal{T}_i^{\text{RS}}$ and street-view-derived text $\mathcal{T}_i^{\text{SV}}$. These descriptions are generated through advanced visual language models (VLMs). They provide contextual insights into urban morphology, infrastructure, and functional attributes, complementing visual data.

Definition 5. Multi-view Urban Graph. The urban system is modeled as a collection of multiple view graphs $\mathcal{G} = \{\mathcal{G}^{(k)}\}_{k=1}^K$, where each view $\mathcal{G}^{(k)} = (\mathcal{V}, \mathbf{A}^{(k)})$ shares the set of common nodes $\mathcal{V} = \{v_i\}_{i=1}^N$ representing the regions of the urban grid, but has distinct adjacency matrices $\mathbf{A}^{(k)} \in \mathbb{R}^{N \times N}$. Each view graph captures a specific type of urban relationship (e.g. POI-based functional similarity, mobility flow, or spatial proximity). This structure enables comprehensive modeling of the urban system through complementary perspectives while maintaining consistent regional representation across views.

Definition 6. Urban Region Representation Learning. Given a set \mathcal{R} of urban regions and K modal feature matrices $\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^K$ derived from multi-modal data sources (e.g. imagery \mathcal{I} , textual descriptions \mathcal{T}), we

233 aim to learn a mapping function $F : (r_i, \mathbf{x}_i^1, \mathbf{x}_i^2, \dots, \mathbf{x}_i^V) \rightarrow \mathbf{h}_i$ that transforms
 234 a region $r_i \in \mathcal{R}$, described by its feature vectors $\mathbf{x}_i^k \in \mathbf{X}^k$ ($1 \leq k \leq K$), into
 235 a d -dimensional representation $\mathbf{h}_i \in \mathbb{R}^d$, where d is a small constant. The
 236 resulting region embeddings $\mathcal{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N\}$ should preserve essential
 237 urban characteristics across all modalities, enabling their effective application
 238 to a wide range of downstream tasks $\mathbf{Y} \in \mathbb{R}^{N \times K}$ across N regions for K
 239 different socioeconomic and environmental attributes.

240 3.2. Overview

241 The proposed UrbanMMCL framework (Figure 1) enriches urban region
 242 representations through a dual-stage contrastive learning approach that es-
 243 tablishes a self-supervised pre-training paradigm for urban tasks.

244 **Stage 1: Multi-Modal Multi-View Contrastive Pre-training**
 245 combines two complementary learning paradigms to establish comprehensive
 246 urban representations. The *cross-modal contrastive learning* leverages VLMs
 247 to generate semantic descriptions for both RSIs and SVIs, employing special-
 248 ized encoders (textual, visual, and location) with multiple contrastive objec-
 249 tives including RSI-text alignment, SVI-text alignment, and location-image
 250 correspondence. Simultaneously, the *multi-view graph contrastive learning*
 251 captures complex spatial dependencies through three distinct view graphs
 252 representing mobility patterns (Mob-view), functional similarities based on
 253 POI attributes (Fun-view), and spatial distance relationships (Dis-view). Us-
 254 ing adaptive graph encoders with independent processing pathways, this com-
 255 ponent dynamically learns optimized graph structures while capturing both
 256 intra-view dependencies and inter-view correlations for comprehensive spatial
 257 relationship modeling.

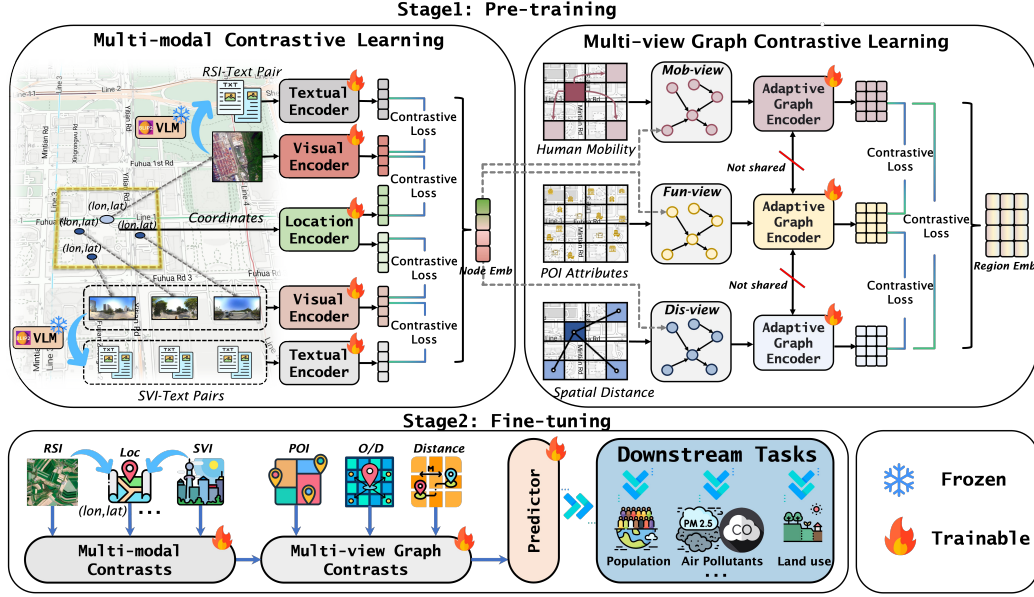


Figure 1: Overview of the UrbanMMCL framework. The framework operates in two stages: (1) **Pre-training Stage** employs multi-modal contrastive learning (aligning RSI, SVI, VLM-generated textual descriptions, and geographical coordinates) alongside multi-view graph contrastive learning across heterogeneous graph views (mobility, functional, and distance) using adaptive graph encoders; (2) **Fine-tuning Stage** integrates the pre-trained multi-modal contrasts and multi-view graph contrasts through dedicated fusion mechanisms, employing trainable predictors for efficient knowledge transfer to downstream urban analytics tasks.

258 **Stage 2: Unified Fine-tuning** integrates the pre-trained multi-modal
 259 and multi-view representations through dedicated fusion mechanisms, en-
 260 abling efficient knowledge transfer to diverse downstream urban analytics
 261 tasks including population estimation, pollutant emission prediction (PM_{2.5},
 262 CO), and land use classification with minimal computational overhead and
 263 labeled data requirements. Details are provided in Sections 3.3, 3.4, and 3.5.

264 3.3. Multi-modal contrastive learning

265 3.3.1. VLM knowledge distillation for text generation

266 We employ BLIP-2, a state-of-the-art vision-language model, to perform
267 knowledge distillation from large-scale pre-trained models, extracting rich
268 semantic information from urban imagery. While advanced models like
269 GPT-4V or Gemini possess extensive knowledge capabilities, their prohibitive
270 costs make them impractical for our dataset of tens of thousands of images.
271 BLIP-2 provides an efficient alternative through knowledge distillation via
272 its lightweight Querying Transformer (Q-Former) architecture, which bridges
273 a frozen image encoder and a frozen language model without requiring end-
274 to-end fine-tuning, significantly reducing computational demands while dis-
275 tilling comprehensive knowledge into high-quality textual descriptions. For
276 each RSI or SVI, BLIP-2 processes the input with prompts to generate de-
277 scriptive text that distills general knowledge into urban-specific semantic
278 representations. Figure 2 illustrates these pairs of prompt descriptions.

279 3.3.2. Vision-language-location feature encoders

280 Using VLM-enhanced text generation, we create a dataset of visual-
281 textual pairs $(\mathcal{I}, \mathcal{T})$, where \mathcal{I} represents RSIs I^{RS} or SVIs \mathcal{I}^{SV} , and \mathcal{T} in-
282 cludes the corresponding textual descriptions \mathcal{T}^{RS} or \mathcal{T}^{SV} . We implement a
283 factorized encoder architecture with dedicated visual, textual, and location
284 encoders, enabling each to capture modality-specific characteristics while es-
285 tablishing the foundation for multimodal alignment.

286 **Visual encoder.** We deploy the Vision Transformer (ViT) architecture
287 (Dosovitskiy et al., 2021) to process urban RSIs and SVIs. Recognizing
288 that standard pre-trained models are optimized for general scenes rather



Figure 2: Examples of prompts and corresponding BLIP-2 generated descriptions for RSI and SVI.

289 than urban environments, we fine-tune this encoder to better capture the
 290 unique structural patterns and spatial relationships characteristic of urban
 291 landscapes.

292 Our encoding process begins by dividing each input image \mathcal{I}_i into p non-
 293 overlapping patches. Each patch \mathbf{P}_j is flattened and projected into a d -
 294 dimensional embedding space with positional encodings:

$$\mathbf{z}_j^{\text{vis}} = \mathbf{E}^{\text{vis}} \cdot \text{Flatten}(\mathbf{P}_j) + \mathbf{p}_j^{\text{vis}}, \quad j = 1, 2, \dots, p, \quad (2)$$

295 where $\mathbf{E}^v \in \mathbb{R}^{(h \cdot w) \times d}$ is a learnable projection matrix and $\mathbf{p}_j^v \in \mathbb{R}^d$ is the

296 positional embedding.

297 These patch embeddings $\{\mathbf{z}_j^{\text{vis}}\}_{j=1}^p$, prepended with a [CLS] token, are
 298 processed through L Transformer layers. Each layer applies multi-head self-
 299 attention (MSA) followed by a multi-layer perceptron (MLP):

$$\mathbf{z}^{\text{vis}'} = \text{MSA}(\text{LN}(\mathbf{z}^{\text{vis}})) + \mathbf{z}^{\text{vis}} \quad (3)$$

$$\mathbf{z}^{\text{vis}''} = \text{MLP}(\text{LN}(\mathbf{z}^{\text{vis}'})) + \mathbf{z}^{\text{vis}'} \quad (4)$$

300 The self-attention mechanism allows each patch to attend to all others.

$$\text{Attention}(\mathbf{Q}^{\text{vis}}, \mathbf{K}^{\text{vis}}, \mathbf{V}^{\text{vis}}) = \text{softmax}\left(\frac{\mathbf{Q}^{\text{vis}}(\mathbf{K}^{\text{vis}})^\top}{\sqrt{d_k}}\right) \mathbf{V}^{\text{vis}} \quad (5)$$

301 where $\mathbf{Q}^{\text{vis}} = \mathbf{W}_Q^{\text{vis}} \mathbf{z}^{\text{vis}}$, $\mathbf{K}^{\text{vis}} = \mathbf{W}_K^{\text{vis}} \mathbf{z}^{\text{vis}}$, $\mathbf{V}^{\text{vis}} = \mathbf{W}_V^{\text{vis}} \mathbf{z}^{\text{vis}}$ are linear projec-
 302 tions.

303 After processing through all transformer layers, we obtain the following:

$$\mathbf{X}_i^{\text{vis}} = \phi_{\text{vis}}(\mathcal{I}_i) = \text{Transformer}(\{\mathbf{z}_j^{\text{vis}}\}_{j=1}^p), \quad (6)$$

304 where $\mathbf{X}_i^{\text{vis}} = [\mathbf{x}_{\text{CLS}}^{\text{vis}_i}, \mathbf{x}_1^{\text{vis}_i}, \mathbf{x}_2^{\text{vis}_i}, \dots, \mathbf{x}_p^{\text{vis}_i}] \in \mathbb{R}^{d \times (p+1)}$ contains embeddings
 305 for the [CLS] token and all image patches. We extract the [CLS] token
 306 embedding $\mathbf{x}_{\text{CLS}}^{\text{vis}_i}$ as the global image representation:

$$\mathbf{v}_i := \mathbf{x}_{\text{CLS}}^{\text{vis}_i}, \quad \mathbf{v}_i \in \mathbb{R}^d. \quad (7)$$

307 The resulting vector \mathbf{v}_i^{rs} or \mathbf{v}_i^{sv} serves as our visual feature for subsequent
 308 cross-modal alignment.

309 **Textual encoder.** Concurrently, we employ a transformer encoder ar-
 310 chitecture (Vaswani et al., 2017) to process textual descriptions generated by
 311 BLIP-2 for our urban imagery.

312 Given a text sequence \mathcal{T}_j with tokens n , we tokenize it and map each token
 313 to an embedding vector using a learnable embedding matrix $\mathbf{E}^t \in \mathbb{R}^{V \times d}$,
 314 where V is the vocabulary size and d is the embedding dimension. Positional
 315 embeddings are added to preserve sequential information:

$$\mathbf{z}_i^{\text{text}} = \mathbf{E}^{\text{text}}[t_i] + \mathbf{p}_i^{\text{text}}, \quad i = 1, 2, \dots, n. \quad (8)$$

316 The sequence passes through L transformer layers:

$$\mathbf{X}_j^{\text{text}} = \phi_{\text{text}}(\mathcal{T}_j) = \text{Transformer}(\{\mathbf{z}_i^{\text{text}}\}_{i=1}^n), \quad (9)$$

317 where $\mathbf{X}_j^{\text{text}} = [\mathbf{h}_1^{\text{text}_j}, \mathbf{h}_2^{\text{text}_j}, \dots, \mathbf{h}_{\text{EOS}}^{\text{text}_j}] \in \mathbb{R}^{d \times (n+1)}$. We extract the [EOS]
 318 token embedding as the global text representation:

$$\mathbf{t}_j = \mathbf{h}_{\text{EOS}}^{\text{text}_j}, \quad \mathbf{t}_j \in \mathbb{R}^d. \quad (10)$$

319 This global text embedding \mathbf{t}_j^{rs} or \mathbf{t}_j^{sv} captures the semantic content of
 320 the entire description, enabling alignment with visual features in subsequent
 321 stages.

322 **Locational encoder.** We integrate geospatial context by encoding the
 323 precise coordinates associated with each visual image. For RSIs, we encode
 324 the center coordinates, while for SVIs, we encode the exact sampling loca-
 325 tions. Inspired by GeoCLIP (Vivanco et al., 2023), our location encoder
 326 transforms geographic coordinates (lon, lat) into meaningful semantic rep-
 327 resentations, producing high-dimensional embeddings that capture spatial
 328 context.

$$\mathbf{l}_i = \phi_l(\text{lon}, \text{lat}) \quad (11)$$

329 These locational embeddings are also further associated with correspond-
 330 ing outputs derived from RSI (\mathbf{l}_i^{rs}) and SVI (\mathbf{l}_i^{sv}). They can capture crucial

geographic context that complements visual and textual features, enabling our model to understand location-specific urban patterns within a unified representation that preserves both semantic content and spatial relationships.

3.3.3. Cross-modality alignment

We implement a cross-modality alignment framework that integrates visual, textual, and location features through specialized contrastive objectives. For each RSI, we aggregate features from all corresponding SVIs via average pooling, creating comprehensive representations that preserve both spatial correspondence and semantic coherence across diverse urban data modalities.

Vision-language alignment. We establish bidirectional alignment between visual and textual modalities through contrastive learning, optimizing the relationships between visual features and their corresponding textual features. The principle is that representations of the same region across modalities converge in semantic space while remaining distinct from other regions.

For each visual-textual pair, we jointly optimize encoders by contrasting matched pairs against others within the batch through a dual-directional contrastive loss:

$$\mathcal{L}_{\text{con}}^{\text{vis-text}} = -\frac{1}{N} \left(\sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{v}_i, \mathbf{t}_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{v}_i, \mathbf{t}_j)/\tau)} + \sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{t}_i, \mathbf{v}_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{t}_i, \mathbf{v}_j)/\tau)} \right), \quad (12)$$

where \mathbf{v}_i represents visual features (either \mathbf{v}_i^{sv} or \mathbf{v}_i^{rs}), \mathbf{t}_i denotes the corresponding textual features \mathbf{t}^{sv} and \mathbf{t}^{rs} , $\text{sim}(\cdot, \cdot)$ represents cosine similarity and

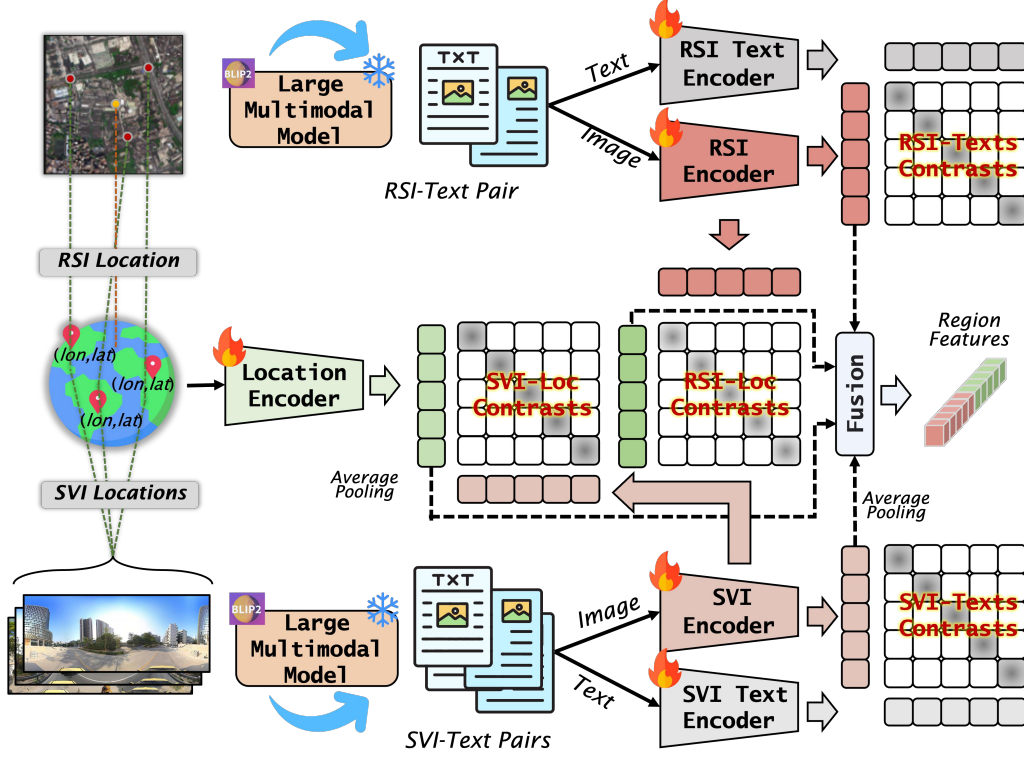


Figure 3: Cross-modal contrastive learning framework integrating visual, textual, and location features through RSI-text, SVI-text, and location-image contrastive objectives.

352 τ is a temperature hyperparameter that controls the similarity distribution.
 353 The first term optimizes image-to-text retrieval, while the second addresses
 354 text-to-image retrieval, creating a unified semantic space for effective cross-
 355 modal understanding.

356 **Visual-location alignment.** To incorporate spatial context, we align

visual characteristics \mathbf{v}_i with location features \mathbf{l}_i using contrastive learning:

$$\begin{aligned} \mathcal{L}_{\text{con}}^{\text{vis-loc}} = & -\frac{1}{N} \left(\sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{v}_i, \mathbf{l}_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{v}_i, \mathbf{l}_j)/\tau)} \right. \\ & \left. + \sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{l}_i, \mathbf{v}_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{l}_i, \mathbf{v}_j)/\tau)} \right), \end{aligned} \quad (13)$$

where \mathbf{l}_i denotes the corresponding location features (either \mathbf{l}_i^{sv} or \mathbf{l}_i^{rs}), and τ is the temperature hyperparameter. Bidirectional loss optimizes both location retrieval from visual features and visual content retrieval from coordinates, associating street-level features with precise locations while capturing broader spatial relationships in overhead views.

Overall objective. The complete multi-modal contrastive learning loss integrates all cross-modal alignments:

$$\mathcal{L}_{\text{mmcl}} = \mathcal{L}_{\text{con}}^{\text{rsi-text}} + \mathcal{L}_{\text{con}}^{\text{rsi-loc}} + \mathcal{L}_{\text{con}}^{\text{svi-text}} + \mathcal{L}_{\text{con}}^{\text{svi-loc}}. \quad (14)$$

This unified optimization creates a shared embedding space where visual, textual, and spatial information are semantically coherent and mutually reinforcing.

Fusion of multi-modal region features. To construct comprehensive region representations, we integrate features from multiple modalities. Given that each region’s RSI typically encompasses multiple SVIs, we first consolidate the SVI features using an averaging operation:

$$\bar{\mathbf{v}}_i^{\text{sv}} = \frac{1}{k} \sum_{j=1}^k \mathbf{v}_j^{\text{sv}}, \quad \bar{\mathbf{l}}_i^{\text{sv}} = \frac{1}{k} \sum_{j=1}^k \mathbf{l}_j^{\text{sv}} \quad (15)$$

where k is the number of SVIs within the region i .

373 The final region feature representation is formed by concatenating RSI vi-
 374 sual features, RSI location features, aggregated SVI features, and aggregated
 375 SVI location features:

$$\mathbf{X}_i = [\mathbf{v}_i^{\text{rs}}; \mathbf{l}_i^{\text{rs}}; \bar{\mathbf{v}}_i^{\text{sv}}; \bar{\mathbf{l}}_i^{\text{sv}}] \quad (16)$$

376 This unified embedding seamlessly combines aerial and ground-level vi-
 377 sual information with their corresponding geographic contexts, providing a
 378 holistic representation that serves as the initial node feature for subsequent
 379 region graph construction.

380 3.4. Multi-view graph contrastive learning

381 To enhance urban region representations with complex spatial relation-
 382 ships, we employ adaptive multi-view graph contrastive learning that cap-
 383 tures intricate interdependencies between urban regions. This component
 384 leverages heterogeneous graph views that represent spatial proximity, func-
 385 tional similarity, and mobility flows, using adaptive graph encoders with
 386 VGAEs and GCNs to dynamically learn optimized graph structures rather
 387 than relying on fixed topologies. The framework incorporates random walk-
 388 based subgraph sampling and employs inter-view contrastive learning to
 389 model spatial dependencies while preserving view-specific characteristics. This
 390 multi-view approach transcends single-view limitations by learning optimal
 391 graph representations that demonstrate enhanced robustness against noise
 392 and data sparsity in urban spatial interactions.

393 3.4.1. Multi-view region graph construction

394 We model geographic regions through a multi-view graph $\mathcal{G} = \{\mathcal{G}^{(k)} =$
 395 $(\mathcal{V}, \mathbf{A}^{(k)}) | k \in \mathcal{K}\}$, where \mathcal{V} represents region nodes, $\mathbf{A}^{(k)}$ denotes the adja-

396 cency matrix of the relationship type k , and $\mathcal{K} = \{P, M, D\}$ specifies rela-
 397 tionship type sets. Our multi-view representation integrates three comple-
 398 mentary urban relationships: **human mobility** flows between regions, **POI**
 399 **category similarity** reflecting functional characteristics, and **geographical**
 400 **distance** capturing spatial proximity. The adjacency matrix $\mathbf{A}^{(k)} \in \mathbb{R}^{N \times N}$
 401 for each relationship type $k \in \mathcal{K}$ defines the pairwise connections between
 402 the regions, where $\mathbf{A}_{ij}^{(k)}$ quantifies the strength of the connection between the
 403 regions v_i and v_j under the corresponding relationship.

404 **Function-aware Region Graph** $\mathcal{G}^{(P)} = (\mathcal{V}, \mathbf{A}^{(P)})$: We characterize ur-
 405 ban functionality through POIs using a distribution matrix $\mathcal{P} \in \mathbb{R}^{N \times C}$, where
 406 C denotes POI categories (restaurants, hotels, hospitals, etc.). Each element
 407 $p_{i,c}$ counts places in the region r_i belonging to the c -th POI category, with the
 408 functionality of each region encoded as a vector $\mathbf{p}_i \in \mathbb{R}^{1 \times C}$. The adjacency
 409 matrix $\mathbf{A}^{(P)} = [a_{ij}^p] \in \mathbb{R}^{N \times N}$ encodes functional similarity through cosine
 410 similarity: $a_{ij}^p = \text{sim}(\mathbf{p}_i, \mathbf{p}_j)$, allowing information flow between functionally
 411 similar areas regardless of geographical distance.

412 **Mobility-based Region Graph** $\mathcal{G}^{(M)} = (\mathcal{V}, \mathbf{A}^{(M)})$: Human movement
 413 patterns are captured through trajectory records in format (r_s, r_d, m_{sd}) , cap-
 414 turing source/destination regions and departure/arrival times. These trajec-
 415 tories are aggregated into an origin-destination flow matrix $\mathcal{M} = [m_{ij}] \in$
 416 $\mathbb{R}^{N \times N}$, where m_{ij} measures trips from region r_i to r_j . The adjacency matrix
 417 $\mathbf{A}^{(M)} = [a_{ij}^m] \in \mathbb{R}^{N \times N}$ is defined as $a_{ij}^m = \frac{\log(1+m_{ij})}{\sum_{k=1}^N \log(1+m_{ik})}$, using logarithmic
 418 normalization to balance flow variations in regions of different populated
 419 regions while preserving movement patterns.

420 **Distance-based Region Graph** $\mathcal{G}^{(D)} = (\mathcal{V}, \mathbf{A}^{(D)})$: Spatial proximity

relationships are encoded in a distance matrix $\mathcal{D} \in \mathbb{R}^{N \times N}$ based on Euclidean distances between the centroids of the region. The adjacency matrix $\mathbf{A}^{(D)} = [a_{ij}^d] \in \mathbb{R}^{N \times N}$ is calculated as $a_{ij}^d = 1/d_{ij}$, creating stronger connections between the physically proximate regions. This structure facilitates the propagation of information between adjacent or nearby areas that typically share similar urban characteristics due to their spatial proximity.

3.4.2. Variational graph auto-encoder

We employ Variational Graph Auto-Encoder (VGAE) as the first component of our multi-view graph contrastive learning framework. VGAE’s probabilistic nature effectively models variability across all three views, while providing regularized latent representations that prevent overfitting and enable meaningful interpolation between region embeddings.

For each view-specific region graph $\mathcal{G}^{(k)}$ with adjacency matrix $\mathbf{A}^{(k)}$ and feature matrix \mathbf{X} , VGAE employs graph convolutional networks (GCNs) to encode graph structure into latent space parameters, specifically the mean vector $\boldsymbol{\mu}$ and the diagonal covariance vector $\boldsymbol{\sigma}$. For notational simplicity, we omit the superscript (k) :

$$\boldsymbol{\mu} = \text{GCN}_{\mu}(\mathbf{A}, \mathbf{X}), \quad \boldsymbol{\sigma} = \text{GCN}_{\sigma}(\mathbf{A}, \mathbf{X}). \quad (17)$$

Latent representations are sampled using the reparameterization trick. For each node i , we have:

$$\mathbf{z}_i = \boldsymbol{\mu}_i + \boldsymbol{\sigma}_i \odot \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(0, \mathbf{I}), \quad (18)$$

where $\mathbf{z}_i \in \mathbb{R}^d$ is the embedding for node i , and $\boldsymbol{\epsilon}_i$ is random noise. The complete embedding matrix $\mathbf{Z} = [\mathbf{z}_1; \mathbf{z}_2; \dots; \mathbf{z}_N] \in \mathbb{R}^{N \times d}$ contains all node embeddings.

443 The decoder reconstructs the adjacency matrix from the latent node rep-
 444 resentations. For each pair of nodes (i, j) , the edge probability is:

$$\tilde{A}_{ij} = \text{Softplus}(\mathbf{z}_i^\top \mathbf{z}_j), \quad (19)$$

445 where $\text{Softplus}(\cdot)$ is the activation function.

446 The VGAE is trained by optimizing both reconstruction loss and KL-
 447 divergence. The reconstruction loss minimizes the difference between original
 448 and reconstructed adjacency matrices:

$$\mathcal{L}_{\text{recon}} = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \left(A_{ij} - \tilde{A}_{ij} \right)^2. \quad (20)$$

449 The KL divergence loss regularizes the latent space by minimizing diver-
 450 gence between the learned latent distribution $q(\mathbf{z}_i | \mathbf{X}, \mathbf{A}) = \mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i^2))$
 451 and a prior distribution $p(\mathbf{z}_i) = \mathcal{N}(0, \mathbf{I})$:

$$\mathcal{L}_{\text{KL}} = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^d \left(1 + \log \sigma_{ij}^2 - \mu_{ij}^2 - \sigma_{ij}^2 \right), \quad (21)$$

452 This process is applied independently to each view $\mathcal{G}^{(k)}$, yielding respec-
 453 tive reconstructed adjacency matrices $\tilde{\mathbf{A}}^{(k)}$. These view-specific embeddings
 454 capture different aspects of urban region relationships for subsequent multi-
 455 view contrastive learning.

456 3.4.3. Random walk-based subgraph generation

457 To efficiently handle complex urban graph structures and enhance ro-
 458 bustness against data skewness, we employ adaptive random walks on the
 459 reconstructed view-specific graphs $\tilde{\mathcal{G}}^{(P)}$, $\tilde{\mathcal{G}}^{(M)}$, and $\tilde{\mathcal{G}}^{(D)}$. For each node i in a
 460 given graph $\tilde{\mathcal{G}}$, we perform a single random walk of fixed length L to capture

both local and global structural properties while ensuring equal contribution from each node.

Starting from node i , the walker transitions from the current node v_t to neighboring node v_{t+1} based on adaptive transition probabilities:

$$P(v_{t+1}|v_t) = \frac{\tilde{A}_{v_t, v_{t+1}}}{\sum_{v_k \in \mathcal{N}(v_t)} \tilde{A}_{v_t, v_k}}, \quad (22)$$

where $\mathcal{N}(v_t)$ denotes neighbors of v_t , and $\tilde{A}_{v_t, v_{t+1}}$ is the learned edge weight. These adaptive allow transition probabilities to prioritize stronger connections while reducing influence of weaker relationships.

The walker continues for L steps, with probabilistic rather than deterministic sampling to introduce variability and enhance robustness. If a node lacks neighbors, the walk terminates early. Upon completion, the visited node sequence $[v_0, v_1, \dots, v_L]$ forms subgraph $\hat{\mathcal{G}}_i$, including all visited nodes and their interconnecting edges.

This sampling strategy effectively addresses noisy edges by utilizing VGAE-learned weights to guide walks toward relevant connections. For each subgraph, we learn node-level representations through:

$$\mathbf{H}_i = \text{GCN}(\hat{\mathcal{G}}_i, \mathbf{X}[\mathcal{V}_i]), \quad (23)$$

where $\mathbf{H}_i \in \mathbb{R}^{|\mathcal{V}_i| \times d}$ contains learned node embeddings, and $\mathbf{X}[\mathcal{V}_i]$ represents features restricted to the subgraph nodes.

We then aggregate subgraph information into a single comprehensive vector through readout operations:

$$\mathbf{s}_i = \text{READOUT}(\mathbf{H}_i), \quad (24)$$

where $\mathbf{s}_i \in \mathbb{R}^d$ encapsulates structural and semantic information of the neighborhood centered on node i .

The process is repeated independently for each view, generating three sets of subgraph representations $\mathbf{s}_i^{(k)}$, where $k \in \mathcal{K} = \{P, M, D\}$, each capturing view-specific structural properties. Unlike rule-based methods, this approach adaptively captures region dependencies through learnable parameters, providing enhanced resilience against data skewness and noise.

3.4.4. Inter-view contrastive learning

After generating view-specific subgraphs and their embeddings, we design a multi-view contrastive learning objective that maximizes mutual information between different view representations of the same node while minimizing it between different nodes. This core intuition is that subgraphs centered around the same urban region should exhibit semantic similarities despite structural differences across views. This contrastive mechanism achieves automatic feature selection through the InfoNCE loss structure, where the numerator enforces cross-view alignment while the denominator ensures inter-node discrimination. This design amplifies features that are consistently discriminative across multiple views while suppressing view-specific noise that fails to maintain cross-view consistency. Consequently, the model automatically prioritizes meaningful urban relationships while filtering out noisy connections, resulting in robust representations that capture fundamental regional patterns.

For any pair of views (k_1, k_2) , the contrastive loss is formulated as:

$$\mathcal{L}_{\text{con}}^{k_1, k_2} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{s}_i^{(k_1)}, \mathbf{s}_i^{(k_2)})/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{s}_i^{(k_1)}, \mathbf{s}_j^{(k_2)})/\tau)}, \quad (25)$$

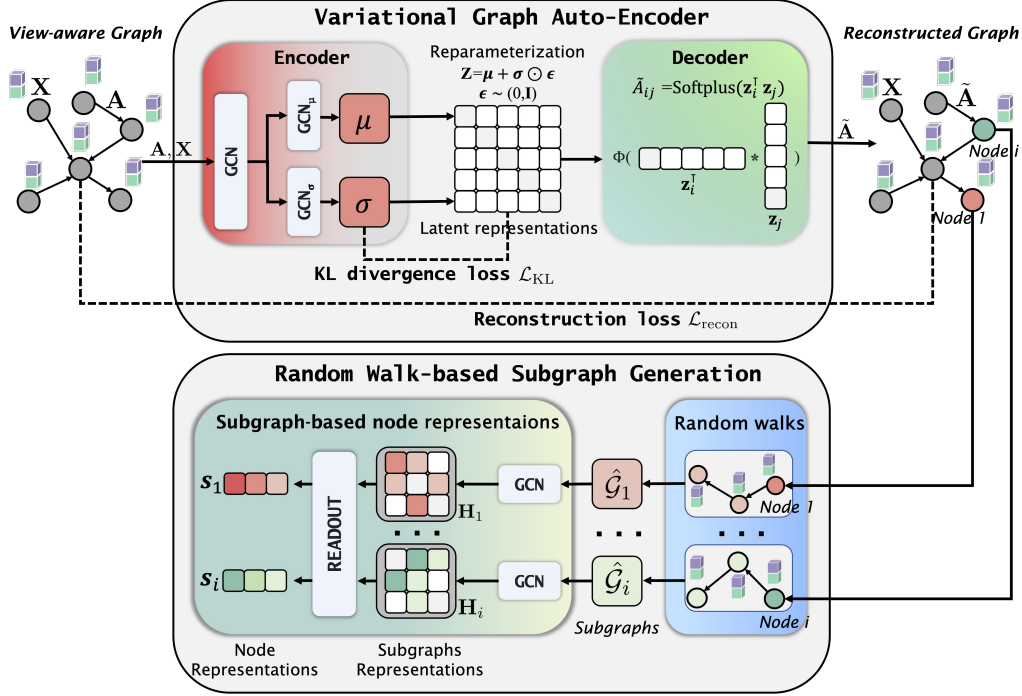


Figure 4: Architecture of the adaptive graph encoder framework. The top component shows the Variational Graph Auto-Encoder (VGAE) that learns latent graph representations through an encoder-decoder structure with reparameterization, incorporating both KL divergence and reconstruction losses for adaptive graph structure learning. The bottom component illustrates the random walk-based subgraph generation process, where random walks from each node create diverse subgraphs that are processed through GCNs to generate enhanced node representations.

where $\mathbf{s}_i^{(k_1)}$ and $\mathbf{s}_i^{(k_2)}$ are the embeddings of node i in views k_1 and k_2 , $\text{sim}(\cdot, \cdot)$ denotes cosine similarity, and τ is a temperature parameter.

Overall objective. The complete multi-view graph contrastive learning loss combines view-specific graph reconstruction with cross-view alignment:

$$\mathcal{L}_{\text{mvgcl}} = \sum_{k \in \mathcal{K}} \left(\mathcal{L}_{\text{recon}}^{(k)} + \mathcal{L}_{\text{KL}}^{(k)} \right) + \sum_{(k_1, k_2) \in \mathcal{K} \times \mathcal{K}} \mathcal{L}_{\text{con}}^{k_1, k_2}. \quad (26)$$

507 where $\mathcal{K} \times \mathcal{K} = \{(P, M), (P, D), (M, D)\}$ represents all view pairs.

508 The first term maintains view-specific structural information through
 509 VGAE reconstruction and regularization losses, ensuring each view preserves
 510 its inherent graph properties. The second term enforces cross-view consis-
 511 tency by aligning representations across pairwise views.

512 3.5. Urban Region Representation Task

513 3.5.1. Pre-training stage

514 Our framework employs a unified pre-training strategy that jointly op-
 515 timizes multi-modal contrastive learning and multi-view graph contrastive
 516 learning. This approach aligns heterogeneous data modalities into a coher-
 517 ent feature space while learning view-invariant representations across differ-
 518 ent urban graph perspectives. The complete pre-training objective combines
 519 both learning stages:

$$\mathcal{L}_{\text{total}} = \alpha \mathcal{L}_{\text{mmcl}} + \beta \mathcal{L}_{\text{mvgcl}}, \quad (27)$$

520 where α and β are hyperparameters to balance cross-modal alignment
 521 and inter-view consistency.

522 This self-supervised framework generates complementary learning signals
 523 through modality fusion and view integration, producing robust representa-
 524 tions that are both semantically meaningful and structurally consistent. Af-
 525 ter pre-training, we obtain region embeddings $\mathbf{Z}^* \in \mathbb{R}^{N \times d}$, providing a strong
 526 foundation for downstream urban analysis tasks.

527 3.5.2. *Fine-tuning for downstream tasks*

528 We employ a task-specific fine-tuning approach to tailor our region em-
529 beddings to various urban prediction tasks. In this process, all model pa-
530 rameters are updated, and the generated embeddings \mathbf{Z}^* are used as input
531 features for a lightweight MLP classifier or regressor: $Y_i = MLP(\mathbf{Z}^*)$. This
532 strategy capitalizes on the rich representations learned during pre-training to
533 demonstrate how effectively our framework captures essential urban patterns
534 across various applications.

535 We evaluated our framework on three downstream tasks: pollutant emis-
536 sion prediction (PEP), population density estimation (PDE), and land use
537 classification (LUC). PEP involves predicting environmental pollutant levels,
538 which tests the model’s ability to capture environmental and spatial factors
539 that influence emissions. PDE is a regression task that estimates regional
540 population density, assessing how well the embeddings capture demographic
541 patterns. LUC is a multi-class classification task categorizing regions into
542 specific land use types (residential, commercial, industrial, etc.), evaluating
543 the model’s capacity to identify distinct urban functional patterns. This
544 fine-tuning approach demonstrates efficient transfer of learned representa-
545 tions to diverse urban applications spanning environmental, demographic,
546 and land-use domains with minimal additional training.

547 4. Experiments and analyses

548 4.1. *Study areas and datasets*

549 As illustrated in Figure 5, our research focuses on Shenzhen, a rapidly
550 developing metropolitan city in China with distinct urban characteristics. he

study leverages diverse datasets for both pre-training and fine-tuning phases to conduct comprehensive experiments and analyses.

The Pre-training Data. The pre-training phase utilizes multi-modal urban data including RSI, SVI, POI data, and human mobility data to learn comprehensive urban representations. RSI, sourced from the GaoFen-2 satellite via Tianditu, features 1.0 meter spatial resolution and three spectral bands (RGB). These images were segmented into $1 \text{ km} \times 1 \text{ km}$ tiles to align with the grid-based analysis framework. **Complementing this, 224,826 high-resolution panoramic SVIs of Shenzhen (4096×1036 pixels), obtained from Baidu Maps at approximately 15-meter intervals along the road network, provide comprehensive 360-degree ground-level visual coverage.** POI data, sourced from AMap¹, includes 1,064,085 points of interest categorized into 23 primary classes such as Life Services, Corporate Entities and Mixed-use Commercial and Residential Areas. Furthermore, human mobility data from China Unicom² consist of 34,960,199 hourly movement records, aggregated to daily origin-destination flows at a $1 \text{ km} \times 1 \text{ km}$ resolution.

The Fine-tuning Data. The fine-tuning phase employs task-specific datasets for downstream urban analysis applications. Population density data from WorldPop³ uses random forest-based dasymetric mapping to deliver high-resolution estimates (30 arc seconds, approximately 1 km at the equator) of people per km^2 , resampled to align with our grid structure. Pollutant emission data include key air pollutants—carbon monoxide (CO,

¹<https://lbs.amap.com>

²<http://www.smartsteps.com/>

³<https://hub.worldpop.org>

573 mg/m^3) and particulate matter ($\text{PM}_{2.5}$, $\mu\text{g}/\text{m}^3$)—sourced from the National
 574 Tibetan Plateau Scientific Data Center (Wei et al., 2023; Wei and Li, 2024),
 575 providing high-quality $1 \text{ km} \times 1 \text{ km}$ resolution raster data. Land use clas-
 576 sification data are derived from SinoLC-1 (Li et al., 2022, 2023c), China’s
 577 first national-scale 1 meter resolution land cover map developed using deep
 578 learning techniques, with the Shenzhen portion specifically utilized for our
 579 analysis.



Figure 5: Overview of the datasets used in this study.

580 4.2. Experiment setup

581 4.2.1. Baselines

582 To comprehensively evaluate our model, we compare with six recent base-
 583 lines in two categories: (1) **Vision-based methods**, including ViT, PG-
 584 SimCLR, and UrbanVLP, where the latter two employ contrastive learn-
 585 ing strategies. (2) **Graph-based methods**, including MVURE, HREP and
 586 ReMVC:

- 587 • **ViT** ([Dosovitskiy et al., 2021](#)). ViT adapts transformers to com-
588 puter vision by partitioning images into fixed-size patches and demon-
589 strates strong performance with sufficient pre-training data. In our
590 experiments, we employ ViT-B as the baseline image encoder, con-
591 catenating extracted features from different modalities for final em-
592 beddings.
- 593 • **PG-SimCLR** ([Xi et al., 2022](#)). A contrastive learning framework
594 that adapts SimCLR([Chen et al., 2020](#)) for urban region representa-
595 tion using satellite imagery by incorporating geographic proximity con-
596 straints and POI category distributions, allowing the model to learn
597 representations that respect both spatial relationships and functional
598 similarities.
- 599 • **UrbanVLP** ([Hao et al., 2025](#)). A multi-granularity vision-language
600 pretraining framework that combines RSI, SVI, and high-quality tex-
601 tual descriptions to predict urban socioeconomic indicators through
602 cross-modal alignment and automatic text calibration.
- 603 • **MVURE** ([Zhang et al., 2020](#)). Leverages human mobility data and
604 urban region attributes (POI and check-in data) to construct multi-
605 view correlations through graph attention networks, enabling cross-
606 view information sharing and adaptive fusion for comprehensive urban
607 region embeddings. In our experiments, we did not use check-in data.
- 608 • **HREP** ([Zhou et al., 2023a](#)). A relation-aware graph-based ap-
609 proach using human mobility, POI information, and geographic neigh-
610 bor data, combined with prompt learning to capture intra-region and

611 inter-region correlations for robust region embeddings.

- 612 • **ReMVC (Zhang et al., 2023a)**. Employs multi-view contrastive
613 learning with POI data and human mobility records to extract robust
614 region embeddings by capturing intra-view distinctions and cross-view
615 correlations.

616 4.2.2. Evaluation metrics

617 To quantitatively evaluate the performance of our method, we employ
618 standard metrics for regression and classification tasks (\downarrow indicates lower is
619 better, and \uparrow indicates higher is better):

620 **Metrics for Regression Tasks.** For regression tasks, the goal is to pre-
621 dict continuous variables. We assess model performance using three comple-
622 mentary metrics: Mean Absolute Error (MAE), Mean Squared Error (MSE),
623 and Coefficient of Determination (R^2).

624 **Metrics for Classification Tasks.** For classification tasks, our objec-
625 tive is to evaluate the ability of the model to correctly classify samples into
626 their respective categories. We employ two complementary metrics: F1 Score
627 and Recall. These metrics provide a comprehensive assessment of the accu-
628 racy and robustness of our method across different task types.

629 4.2.3. Implementation details

630 **Pre-training Setup.** We divide the dataset into 60% training, 20%
631 validation and 20% testing sets. For image inputs, we apply data augmenta-
632 tion techniques that include random cropping, flipping, and normalization,
633 following the methodology described in Radford et al. (2021a). During the
634 multi-modal contrastive learning stage, we generate textual descriptions for

635 each image using the BLIP-2 model, leveraging OPT-2.7b (a large language
 636 model with 2.7 billion parameters). Text descriptions are limited to 77 to-
 637 kens. For the visual encoder, we use ViT-B/32 with a hidden dimension of
 638 768 and an output dimension of 512. For the textual encoder, both the hid-
 639 den and output dimensions are 512. For the location encoder, we adopt the
 640 architecture and configuration settings proposed by Vivanco et al. (2023).
 641 For feature aggregation, SVI features within the same RSI coverage area are
 642 aggregated using average pooling. The aggregated SVI features are then
 643 added element-wise to the RSI features to form unified regional representa-
 644 tions. In the multi-view graph contrastive learning stage, we construct three
 645 graph structures based on population flow, POI similarity, and spatial dis-
 646 tance. VGAE is used to reconstruct these graphs. VGAE uses a one-layer
 647 GCN encoder with an output latent space size of 64, and LeakyReLU ac-
 648 tivation, and Adam optimizer with learning rate 1e-4. Next, random walks
 649 with walk length 20 are performed on the reconstructed graphs to generate
 650 subgraphs. Finally, GCNs with hidden layer size 128, output size 64, and
 651 LeakyReLU activation are applied for graph representation embedding. The
 652 temperature parameter τ for contrastive learning is 0.5. The Adam optimizer
 653 is used with a learning rate of 1e-6 and weight decay of 1e-4. The model is
 654 trained for 1000 epochs with early stopping based on validation loss.

655 **Fine-tuning Setup.** During fine-tuning, we update both the pre-trained
 656 encoders and the task-specific components. For regression tasks, we use mean
 657 squared error loss, while for classification tasks, we use cross-entropy loss
 658 with accuracy as the evaluation metric. The training runs for 1000 epochs
 659 using the Adam optimizer, with learning rate 1e-6 and weight decay 1e-4.

660 Early stopping is applied based on validation loss to prevent overfitting. All
661 experiments were performed on NVIDIA A6000 GPUs with 48GB memory.

662 4.3. Model performance

663

664 4.3.1. Hyperparameter sensitivity analysis

665 To understand the interaction between multi-modal and multi-view graph
666 contrastive learning, we conduct sensitivity analysis of hyperparameters α
667 and β in our unified objective function $\mathcal{L}_{\text{total}} = \alpha\mathcal{L}_{\text{mmcl}} + \beta\mathcal{L}_{\text{mvgcl}}$ with
668 constraint $\alpha + \beta = 1$. Table 1 shows the optimal configuration occurs at
669 $\alpha = 0.5, \beta = 0.5$ across all tasks, indicating that indicating equal weighting
670 between multi-modal contrastive learning and spatial learning yields best
671 performance.

672 The results also reveal asymmetric degradation patterns. Pure graph
673 learning ($\alpha = 0$) causes dramatic performance drops, while pure multi-modal
674 learning ($\alpha = 1.0$) shows moderate decreases. This asymmetry indicates that
675 multi-modal information provides fundamental semantic grounding, while
676 spatial learning offers crucial structural guidance. The consistent optimal
677 ratio ($\alpha : \beta = 1 : 1$) across diverse urban tasks suggests that effective urban
678 representation learning requires balanced integration of semantic richness
679 and urban spatial structure, , rather than over-relying on either component
680 alone.

681 4.3.2. Comparison with baseline methods

682 We conduct comprehensive comparisons with state-of-the-art methods.
683 Table 2 presents the overall results, from which we can derive the following

α	β	PEP(CO)			PEP(PM _{2.5})			PDE			LUC	
		MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	F1 ↑	Recall ↑
1.0	0.0	0.0382	0.0018	0.4025	1.4009	3.2739	0.5019	4134.67	14969395.03	0.4817	0.4106	0.4235
0.8	0.2	0.0248	0.0012	0.6034	1.2431	2.1900	0.6756	3816.31	15647112.25	0.4841	0.4189	0.4673
0.6	0.4	0.0250	0.0012	0.5927	0.8591	1.4778	0.7654	3622.39	13598495.15	0.6556	0.4504	0.4957
0.5	0.5	0.0212	0.0009	0.7417	0.6785	1.0212	0.8481	1931.64	9943831.52	0.7670	0.6058	0.5994
0.4	0.6	0.0247	0.0012	0.6857	0.8140	1.3593	0.7895	3656.10	16647168.90	0.6777	0.4844	0.5229
0.2	0.8	0.0238	0.0011	0.6290	1.0394	1.7856	0.7234	3860.65	19548932.31	0.5648	0.4738	0.5068
0.0	1.0	0.0327	0.0018	0.4341	1.4808	3.6995	0.4824	3824.24	17542438.41	0.5194	0.4223	0.4445

Table 1: Hyperparameter sensitivity analysis for α and β . The best results are highlighted in **boldface**.

Methods	PEP(CO)			PEP(PM _{2.5})			PDE			LUC	
	MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	F1 ↑	Recall ↑
ViT	0.0252	0.0012	0.5569	1.0062	1.7733	0.7225	2086.0500	12767905.19	0.6654	0.5136	0.5235
PG-SimCLR	0.0358	0.0020	0.4158	1.4840	3.6337	0.4280	2194.7627	13386507.27	0.7116	0.4690	0.4555
UrbanVLP	<u>0.0214</u>	<u>0.0011</u>	<u>0.6875</u>	<u>0.7573</u>	1.1847	<u>0.8240</u>	<u>1954.4890</u>	9754879.86	<u>0.7635</u>	<u>0.5811</u>	<u>0.5873</u>
MVURE	0.0239	0.0011	0.5669	0.8121	<u>1.0288</u>	0.7933	2046.9664	11701711.17	0.6933	0.5450	0.5446
HREP	0.0237	0.0012	0.5604	0.8860	1.5067	0.7677	2089.4557	13625186.03	0.6429	0.4299	0.4291
ReMVC	0.0260	0.0013	0.5217	1.3249	3.0701	0.5106	2242.8103	16310952.66	0.5725	0.4821	0.4718
UrbanMMCL	0.0212	0.0009	0.7417	0.6785	1.0212	0.8481	1931.64	<u>9943831.52</u>	0.7670	0.6058	0.5994
Improvement(%)	0.93	18.18	7.88	10.41	0.74	2.92	1.17	-1.94	0.46	4.25	2.06

Table 2: Performance comparison of different methods on pollutant emission prediction (PEP) (CO, PM_{2.5}), population density estimation (PDE), and land use classification (LUC). The best results are in **boldface**, and the second-best results are underlined. Improvement(%) shows the relative improvement of our method over the second-best baseline.

key findings.

(1) **UrbanMMCL achieves superior performance across most metrics, demonstrating the effectiveness of our dual contrastive learning approach.** Our framework outperforms the best baselines in 10 of 11 metrics, with an average R² improvement of 3.75% in regression tasks and a 4.25% improvement in the F1 score for classification compared to the second best method (UrbanVLP). The only exception is the PDE MSE metric, where our method shows a marginal difference of 1.94%. **This slight**

692 discrepancy stems from the long-tail distribution of population density data,
693 where extreme values disproportionately influence the squared error metric.
694 Our superior MAE and R^2 scores demonstrate robustness across the majority
695 of urban regions.

696 **(2) Text-enhanced vision-language methods significantly out-**
697 **perform single-modality and POI-enhanced approaches.** UrbanVLP
698 consistently outperforms both the vision-only ViT model and the POI-enhanced
699 PG-SimCLR across all tasks, with notable improvements in $PM_{2.5}$ prediction
700 (R^2 : 0.8240 vs 0.4280 for PG-SimCLR). This confirms that rich textual de-
701 scriptions provide more contextually relevant information than structured
702 POI data alone. UrbanMMCL further advances this paradigm by effectively
703 integrating visual features with geographical coordinates and adaptive graph
704 relationships.

705 **(3) Adaptive graph contrastive learning significantly outper-**
706 **forms static graph-based methods.** Unlike existing graph-based meth-
707 ods (MVURE, HREP and ReMVC) that rely on predetermined region rela-
708 tionships, UrbanMMCL uses VGAE and adaptive random walks to automat-
709 ically learn and refine meaningful region connections. Our method achieves
710 a remarkable 30.8% R^2 improvement in CO prediction over MVURE, high-
711 lighting how our adaptive approach addresses the limitations of fixed graph
712 structures in complex urban environments.

713 **(4) The synergy between multi-modal integration and multi-**
714 **view graph modeling creates generalizable urban representations.**
715 The integration of RSI, SVI, geographical positions, and textual descriptions
716 through dual-stage contrastive learning allows UrbanMMCL to capture both

Methods	PEP(CO)			PEP(PM _{2.5})			PDE			LUC	
	MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	F1 ↑	Recall ↑
RSI-CLIP	0.0313	0.0019	0.4573	1.4185	2.7707	0.5969	2615.4236	16563126.19	0.5844	0.4952	0.5114
SVI-CLIP	0.0288	0.0015	0.5435	1.1461	2.1668	0.6376	2404.1780	14664183.37	0.5596	0.4597	0.4394
<i>w/o</i> Text	0.0246	0.0013	0.6372	0.8683	1.4384	0.7751	1996.8511	10413018.36	0.6957	0.5550	0.5536
<i>w/o</i> MCL	0.0327	0.0018	0.4341	1.4808	3.6995	0.4824	3824.2488	17542438.41	0.5194	0.4223	0.4445
UrbanMMCL	0.0212	0.0009	0.7417	0.6785	1.0212	0.8481	1931.64	9943831.52	0.7670	0.6058	0.5994

Table 3: Ablation on multimodal components. The best results are highlighted in **bold-face**.

Methods	PEP(CO)			PEP(PM _{2.5})			PDE			LUC	
	MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	MAE ↓	MSE ↓	R ² ↑	F1 ↑	Recall ↑
<i>w/o</i> $G^{(P)}$	0.2284	0.0011	0.5342	1.1219	1.9569	0.6983	2527.02	17422814.99	0.5768	0.3852	0.3952
<i>w/o</i> $G^{(M)}$	0.0202	0.0009	0.5910	1.2310	2.5737	0.6111	2629.2102	18889712.73	0.5679	0.4702	0.5124
<i>w/o</i> $G^{(D)}$	0.0276	0.0014	0.3575	1.0027	1.8363	0.6967	2454.5840	17698496.25	0.4292	0.4450	0.4373
<i>w/o</i> VGAE	0.0263	0.0014	0.6020	1.0447	1.9016	0.6771	1935.23	13555443.7517	0.6531	0.4424	0.4428
<i>w/o</i> RW	0.0220	0.0010	0.6554	0.7185	1.2764	0.7947	1912.9005	11225293.52	0.6808	0.4680	0.4851
<i>w/o</i> GCL	0.0335	0.0019	0.4104	1.8679	3.5980	0.5078	2095.8574	13879595.94	0.6071	0.3966	0.4123
UrbanMMCL	0.0212	0.0009	0.7417	0.6785	1.0212	0.8481	1931.64	9943831.52	0.7670	0.6058	0.5994

Table 4: Ablation on multi-view graph components. The best results are highlighted in **boldface**.

717 fine-grained visual details and macro-scale spatial relationships. This com-
718 prehensive modeling creates generalizable features that maintain consistent
719 performance across both regression and classification tasks, from environ-
720 mental monitoring (CO, PM_{2.5}) to socioeconomic analysis (PDE, LUC).

721

722 4.3.3. Cross-city generality

723 To assess the generalization capability of UrbanMMCL, we conducted val-
724 idation studies in Beijing and Chengdu, two cities with contrasting develop-
725 mental and geographic profiles. Our evaluation adopts hierarchical transfer
726 learning leveraging UrbanMMCL’s modular design, where multimodal en-
727 coders trained on Shenzhen are directly transferred while graph components

Cities	Models	PEP(CO)	PEP(PM _{2.5})	PDE	LUC
		$R^2 \uparrow$	$R^2 \uparrow$	$R^2 \uparrow$	$F1 \uparrow$
Beijing	PG-SimCLR	0.4929	0.4301	0.6109	0.4380
	ReMVC	0.5488	0.5157	0.5433	0.4791
	UrbanMMCL	0.7032	0.7811	0.6317	0.5538
	Improvement(%)	+28.15%	+51.46%	+3.40%	+15.59%
Chengdu	PG-SimCLR	0.5171	0.4540	0.5583	0.4581
	ReMVC	0.5728	0.5384	0.5679	0.4904
	UrbanMMCL	0.6912	0.7508	0.6420	0.5395
	Improvement(%)	+20.67%	+39.45%	+13.04%	+10.01%

Table 5: Cross-city transfer learning performance comparison in Beijing and Chengdu.

are re-initialized for city-specific spatial relationships.

Table 5 demonstrates impressive cross-city performance. Both cities achieve strong performance with R^2 scores of 0.63-0.78 across regression tasks and competitive F1 scores of 0.54-0.55 for land use classification, maintaining remarkably consistent results despite diverse urban contexts. UrbanMMCL consistently outperforms baseline methods PG-SimCLR and ReMVC by 10-51%, confirming robust generalization across diverse urban environments.

4.3.4. Ablation studies

To validate our design principles, we conduct comprehensive ablation studies addressing two key questions: (1) What are the essential multi-modal components and integration strategies for effective urban representation learning? (2) What are the essential graph perspectives and learning mechanisms for effective urban spatial relationship modeling?

742 **Ablation on multimodal components.** We design four variants to
743 test specific hypotheses: (1) **RSI-CLIP** and **SVI-CLIP** replace our domain-
744 specific encoders with general-purpose pre-trained CLIP model (ViT-B/32)
745 to assess the necessity of domain specialization; (2) ***w/o* Text** eliminates
746 textual enhancement to quantify semantic information contribution; (3) ***w/o***
747 **MCL removes cross-modal contrastive learning and initializes the encoders**
748 **with their original weights while preserving multimodal fusion to isolate the**
749 **impact of explicit cross-modal alignment.**

750 Table 3 reveals three key findings that validate our design choices. First,
751 replacing specialized components with general CLIP causes substantial degra-
752 dation, demonstrating that urban scene understanding requires architectural
753 adaptations beyond general vision models. Second, removing text consis-
754 tently decreases performance by 14.1% across tasks, with environmental
755 monitoring particularly affected, showing that semantic descriptions cap-
756 ture abstract urban characteristics invisible to visual features alone. Third,
757 eliminating contrastive learning causes the most severe degradation, confirm-
758 ing that explicit cross-modal alignment is essential for coherent multimodal
759 representations. These results demonstrate that each component addresses
760 specific urban representation challenges, and their synergistic integration is
761 critical for optimal performance across diverse urban tasks.

762 **Ablation on multi-view graph components.** We design six variants
763 testing: (1) individual graph view contributions (***w/o* $\mathcal{G}^{(P/M/D)}$**); (2) ***w/o***
764 **VGAE:** replacing the VGAE with standard GCN for graph encoding to as-
765 sess probabilistic graph structure learning; (3) ***w/o* RW:** eliminating random
766 walk-based subgraph generation and using full graphs to evaluate local struc-

ture sampling effectiveness; and (4) *w/o GCL*: removing graph contrastive learning while retaining basic graph encoders to isolate cross-view alignment impact.

Table 4 reveals critical insights into urban spatial modeling. First, distance-based graph removal causes the most severe degradation, particularly impacting the CO prediction, while mobility-based graph and function-aware graph removal results in 24.8% and 23.5% average decrease, demonstrating that geometric relationships serve as fundamental structural foundation with all three perspectives capturing distinct spatial aspects. Second, replacing VGAE with standard GCN decreases performance by 17.9%, while removing random walk sampling causes a 9.7% drop, demonstrating that both probabilistic structure learning and local sampling contribute to effective spatial modeling. Third, eliminating graph contrastive learning results in the largest performance decline with 35.2% average R^2 decrease, confirming that learning coherent multi-perspective representations requires explicit alignment mechanisms. These results demonstrate that effective urban spatial understanding requires integrated design of multiple graph perspectives, adaptive structure learning, and cross-view contrastive alignment.

4.4. Analysis of learned representations

4.4.1. Representation visualizations

To validate that our model learns meaningful representations that effectively distinguish different urban area types, we examine learned region representations through a t-SNE dimensionality reduction, which maps the high-dimensional embeddings learned by UrbanMMCL into an interpretable two-dimensional space. Figure 6 reveals three distinct clusters with high

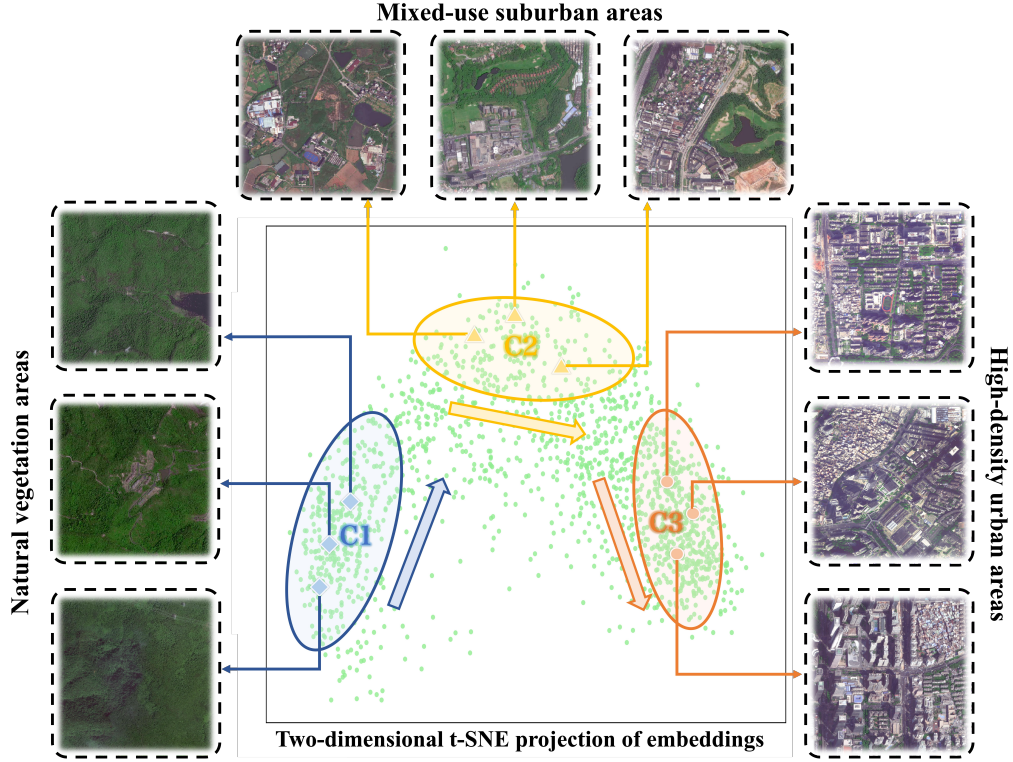


Figure 6: t-SNE visualization of region representations showing three distinct clusters corresponding to different urbanization levels: natural vegetation areas (left), mixed-use suburban regions (center), and high-density urban areas (right).

intra-cluster similarity. The clusters exhibit a progressive urbanization gradient from left to right: areas dominated by natural vegetation, mixed-use suburban regions, and high-density urban areas. This clustering pattern validates our approach successfully captures subtle yet critical geographical differences and maps regions with similar architectural layouts and land use patterns into proximate embedding positions.

To validate the necessity and effectiveness of multi-view fusion over single-view approaches, Figure 7 demonstrates our multi-view approach through

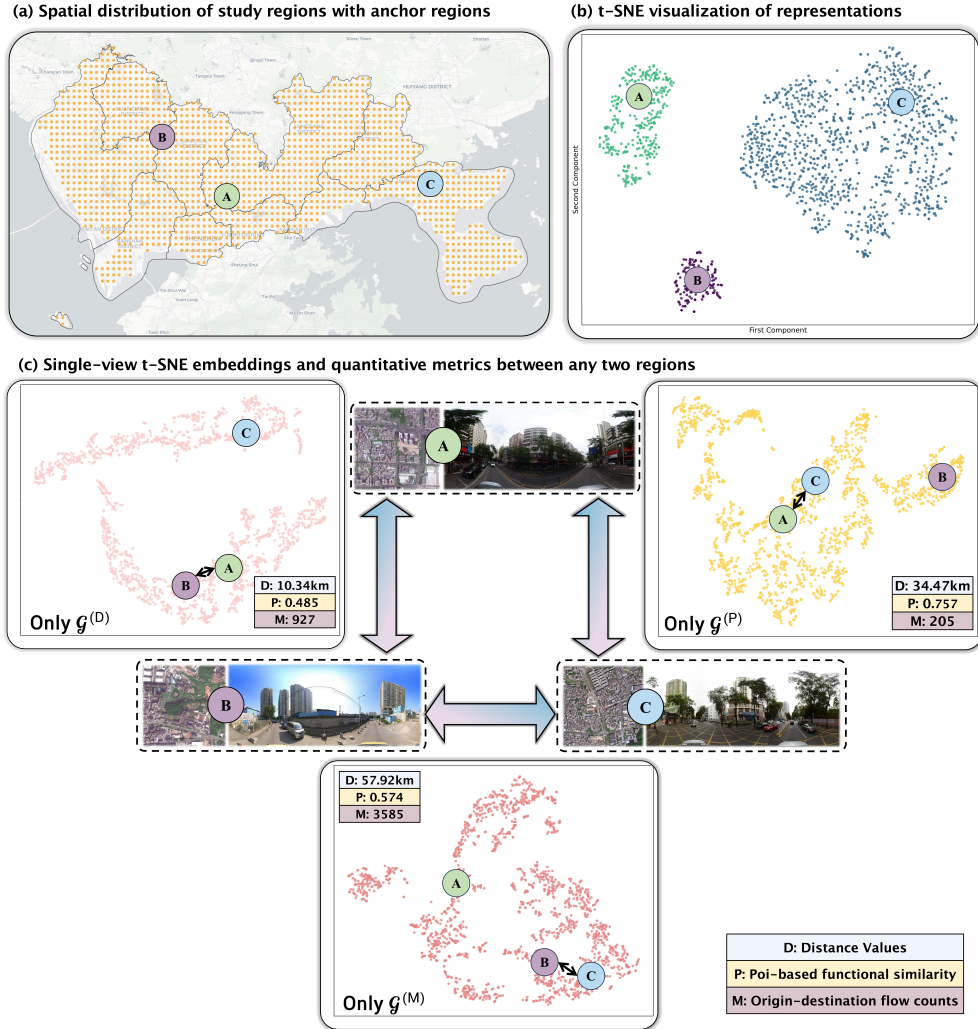


Figure 7: Multi-view region representation analysis. (a) Spatial distribution of study regions with anchor regions A, B, C. (b) UrbanMMCL embedding space showing integrated clustering. (c) Single-view embeddings and quantitative metrics revealing individual graph limitations.

comparative analysis of three anchor regions. Panel (c) reveals limitations of single-view embeddings through individual embedding spaces and quan-

titative relationship metrics (D, P, M values). The only distance-based embedding places regions A-B closely while positioning A-C far apart despite their strong functional similarity. The only POI-based embedding brings functionally similar regions A-C together but inappropriately positions A-B and B-C by neglecting spatial and mobility constraints. The only mobility-based embedding clusters regions B-C closely due to strong movement connections while under-representing A-B and A-C relationships. These positioning biases highlight the the limitations of single-view approaches in capturing comprehensive urban relationships. In contrast, Panel (b) demonstrates our UrbanMMCL embedding space where regions achieve balanced clustering through consensus-based optimization integrating all perspectives. This integrated approach produces robust representations that position regions appropriately by balancing functional similarity, spatial proximity, and mobility connectivity in a unified embedding space.

4.4.2. Geographic mapping of clustered representations

To validate that our UrbanMMCL framework captures meaningful urban structures, we apply hierarchical clustering to the learned embeddings. Figure 8 presents clustering results for $k=2$ to 6, with the dendrogram distances indicating cluster distinctiveness. Clustering analysis demonstrates a clear hierarchical organization of urban spaces. At $k=2$, a fundamental binary partition emerges: built-up areas (pink) and natural areas (cyan), separating urban development zones from mountainous regions and water bodies. As k increases to 3, the urban domain subdivides into high-density cores (Futian, Luohu, Nanshan) and lower-density periphery (Longgang, Guangming, Pingshan), while natural areas remain cohesive. Higher k values (4-6)

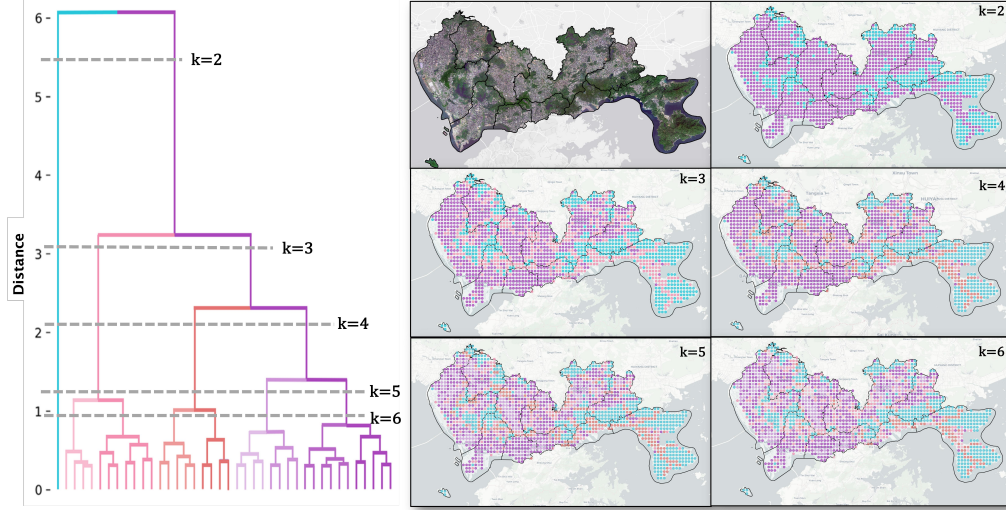


Figure 8: Geographic mapping of clustered representations across different cluster numbers ($k=2$ to $k=6$). The dendrogram (left) shows the hierarchical structure of learned embeddings, with dashed lines indicating cut heights for different k values. The satellite image (top center) provides the geographic context of Shenzhen. The cluster maps (right) visualize the spatial distribution of clusters for each k value.

827 demonstrate progressive refinement within urban areas while maintaining
828 stable natural clusters. This pattern indicates that our embeddings suc-
829 cessfully encode urban heterogeneity, as they capture development intensity
830 variations and functional zones while recognizing the homogeneity of natu-
831 ral landscapes. Such spatially coherent clustering demonstrates the practical
832 utility of our framework for automated urban region categorization.

833 4.4.3. Predictive performance analysis

834 To demonstrate the practical effectiveness of our method, we conduct
835 comprehensive case studies examining prediction performance in representa-
836 tive four urban regions. Our analysis includes: (1) regression analysis with

837 characteristics for CO, PM_{2.5} emissions, and population density prediction,
838 and (2) classification analysis for land use.

839 **Regression Analysis.** We select four representative regions that pro-
840 vide an ideal testbed for evaluating multi-modal and multi-view necessity.
841 Critically, region pairs (A-B and C-D) exhibit similar visual appearances but
842 substantial differences in urban indicators, creating challenging discrimina-
843 tion scenarios.

844 Figure 9 presents regression results across three urban indicators. Our
845 complete UrbanMMCL framework consistently achieves the closest approx-
846 imations to ground truth. When individual modalities are removed, sys-
847 tematic degradation emerges. Eliminating RSI or SVI causes predictions to
848 converge toward averaged values, losing spatial discrimination. For exam-
849 ple, without RSI, CO predictions become nearly uniform (0.33-0.39), failing
850 to capture the actual variation (0.63-0.96). View-specific ablation reveals
851 distinct dependency patterns for different urban indicators. Environmental
852 indicators show greater sensitivity to structural patterns capturing physical
853 processes, while socio-economic indicators correlate more strongly with hu-
854 man behavior modeling and functional interactions. These heterogeneous
855 patterns validate our multi-view approach by demonstrating that different
856 urban processes operate through distinct channels, and no single structural
857 perspective adequately captures urban system complexity.

858 **Land Use Classification Analysis.** We examine prediction perfor-
859 mance across six primary land cover categories: tree cover, building, shrub-
860 land, cropland, traffic route, and grassland. The spatial distribution anal-
861 ysis(Figure 10 a-b) shows that UrbanMMCL predictions closely align with

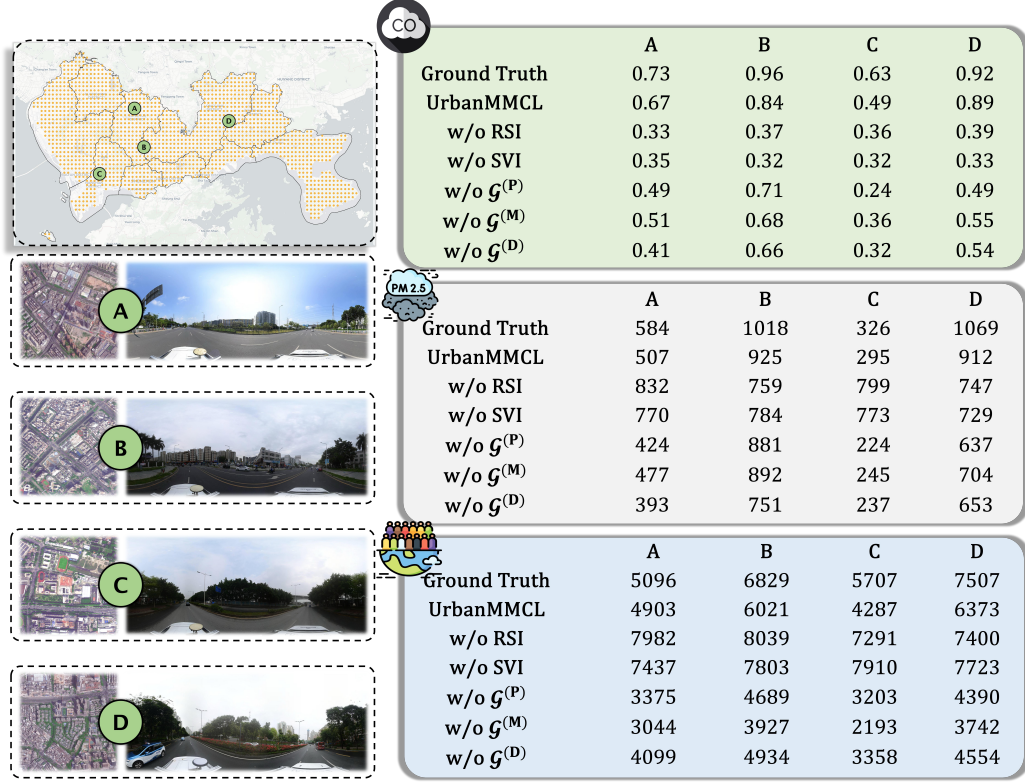


Figure 9: Case Study analysis comparing prediction performance across four representative regions for three urban indicators. The left panel displays the spatial distribution and corresponding SVI/RSI data, while the right panel presents compares ground truth with UrbanMMCL predictions and ablation configurations.

862 ground truth patterns. Our framework successfully captures complex spatial
863 organization and maintains clear boundaries between natural areas and built
864 environments. The confusion matrix (Figure 10 c) reveals varying perfor-
865 mance across categories. Grassland achieves the highest accuracy at 85.3%,
866 followed by shrubland at 75.0%, building at 64.8%, and traffic routes at
867 59.7%. Tree cover and cropland show more challenging classification at 37.7%
868 and 36.4% respectively, likely due to seasonal variations and spectral simi-

869 larity with other vegetation types.

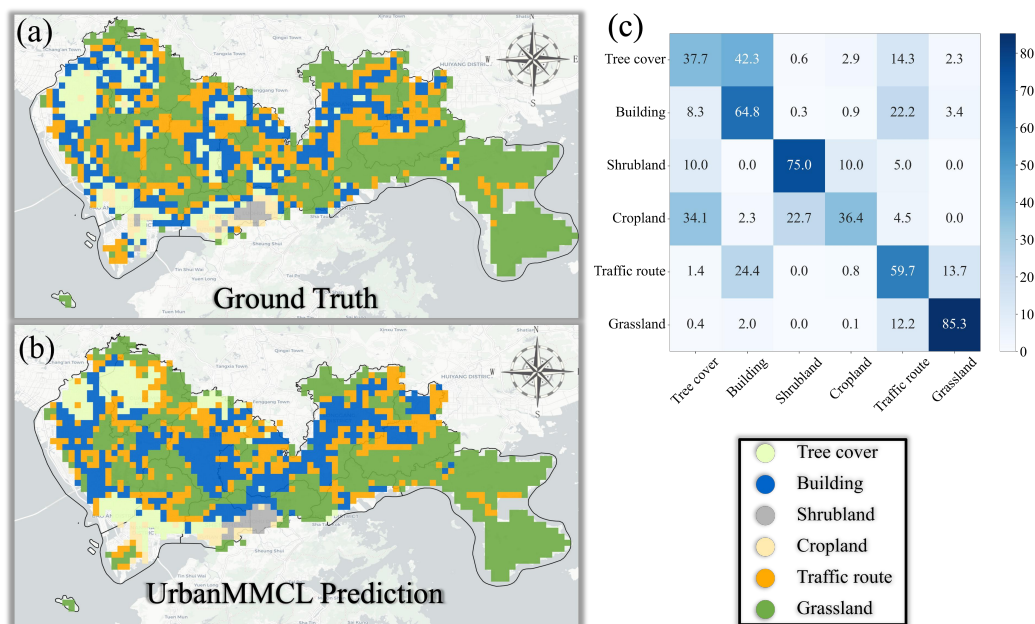


Figure 10: Land use classification analysis comparing UrbanMMCL predictions with ground truth data. (a) Ground truth spatial distribution. (b) UrbanMMCL prediction results. (c) Confusion matrix quantifying classification accuracy for each land cover category.

870 Our complete UrbanMMCL framework consistently achieves the closest
 871 approximations to ground truth across all regions and indicators in both
 872 regression and classification tasks. The multimodal data integration and
 873 multi-view framework ensures that when visual similarities mask functional
 874 differences, complementary perspectives provide the discriminative power
 875 necessary for accurate urban dynamics prediction and land use classification.

5. Discussions

5.1. Multi-modal contribution analysis

Understanding how disparate urban data modalities contribute to representation learning provides insights into feature complementarity and information hierarchies in multi-modal urban analysis. Our ablation experiments reveal distinct roles for each modality in capturing different aspects of urban complexity, with results presented in Figure 11.

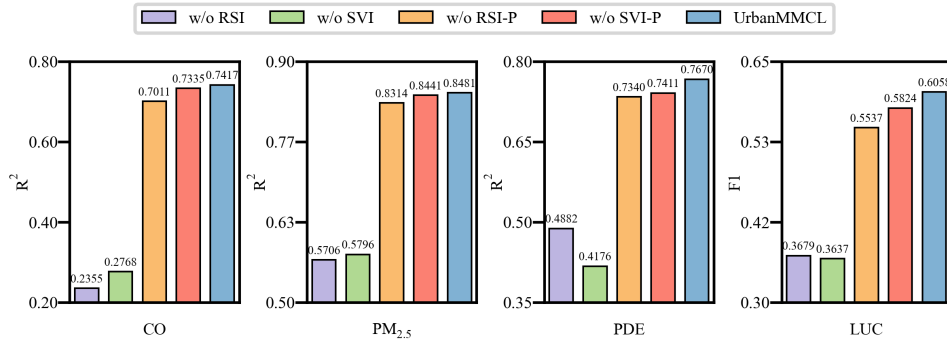


Figure 11: Multi-modal ablation study showing the contribution of different modalities to urban representation learning.

Fundamental Role of Remote Sensing Imagery. RSI emerges as the most fundamental modality, with its removal causing severe performance degradation across all tasks. This dominance stems from RSI’s ability to capture spatial patterns and urban morphological features at scale. For PEP, substantial performance drops reveal that RSI encodes critical spatial dependencies correlating with environmental phenomena. RSI’s high information

density enables learning of rich spatial representations that serve as foundational embeddings for other modalities.

Complementary Value of Street-View Imagery. SVI contributes fine-grained environmental features through local context augmentation. Performance improvements from SVI inclusion demonstrate its role in capturing micro-environmental variations invisible in overhead imagery. SVI functions as local environmental validators that refine broad spatial patterns captured by RSI, particularly evident in LUC where ground-level visual cues help distinguish functionally similar areas.

Spatial Context Enhancement Through Positional Encoding. Geographical coordinates provide modest but consistent contributions as spatial relationship encoders. The relatively small impact when removing positional encoding suggests that visual features carry majority predictive information, while coordinates primarily enhance spatial coherence and topological consistency in learned representations.

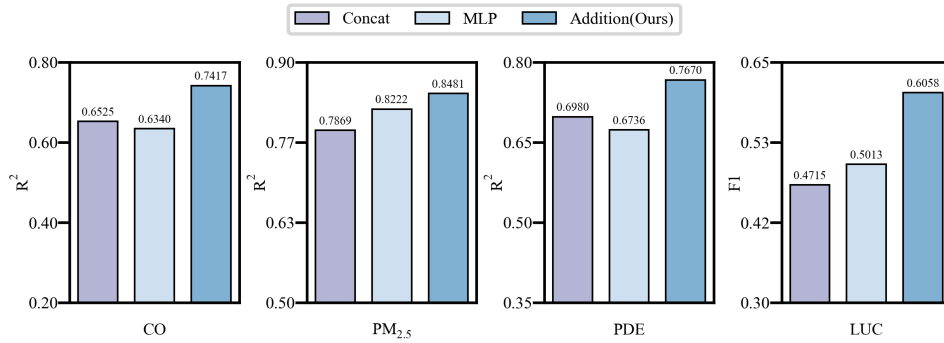


Figure 12: Performance comparison of different multi-modal fusion strategies.

904 **Implications for Urban Representation Learning.** These findings
 905 establish fundamental principles for effective urban AI systems. Perfor-
 906 mance differentials when removing individual modalities demonstrate that
 907 each modality contributes unique, irreplaceable information, establishing a
 908 clear modality hierarchy where RSI provides foundational spatial structure,
 909 SVI adds critical environmental detail, and positional encoding serves as
 910 spatial regularization (Figure 11). More critically, the superior performance
 911 of specialized encoders over generic CLIP-based alternatives (Table 3) and
 912 the critical role of multi-view contrastive learning (Table 4) demonstrate
 913 that urban environments require domain-specific architectures and multi-
 914 perspective integration rather than universal approaches. This advocates for
 915 specialized multi-modal urban AI systems that embrace complexity through
 916 tailored encoders and multi-perspective integration.

917 5.2. Multi-modal feature fusion strategy analysis

918 Multi-modal feature fusion significantly impacts the model’s ability to
 919 leverage complementary information from heterogeneous urban data sources.
 920 The choice of fusion strategy is therefore crucial for maximizing the benefits
 921 of multi-modal urban data integration. We compared three fusion strate-
 922 gies: (1) concatenation of visual and location features from RSI and aggre-
 923 gated SVIs; and (2) MLP-based fusion with multi-layer perceptrons; and (3)
 924 element-wise addition. As shown in Figure 12, our addition method consis-
 925 tently achieves superior performance across all metrics despite its simplicity.

926 The superior performance of element-wise addition can be attributed
 927 to its ability to preserve original feature distributions while enabling direct
 928 correspondence between spatially aligned features from different modalities.

929 Unlike concatenation, which introduces feature redundancy and increased
930 dimensionality, or MLP fusion, which adds parameters and optimization
931 complexity, addition fusion maintains the semantic integrity of individual
932 modalities while creating meaningful cross-modal interactions. This vali-
933 dates our design choice and demonstrates that simpler fusion strategies can
934 be more effective.

935

936 5.3. Training paradigms and efficiency analysis

937 To comprehensively evaluate our model’s representation capabilities and
938 training efficiency, we examine two additional training paradigms that rep-
939 resent different approaches to leveraging pre-trained knowledge for urban
940 downstream tasks. **Pretrain-finetune** first optimizes the encoder on a large-
941 scale, task-agnostic urban data to learn general representations, followed by
942 fine-tuning on downstream tasks; **Linear probing** freezes the pretrained
943 encoder and trains only a linear head, providing an efficient assessment of
944 representation quality with minimal computational resources; **End-to-end**
945 **training** initializes with pre-trained weights but allows unrestricted param-
946 eter updates throughout the entire architecture.

947 Figure 13 presents comparative analysis of three training paradigms across
948 prediction accuracy, runtime per epoch, and epochs to convergence. The
949 bars represent average performance while scattered points show individual
950 task values for CO, PM_{2.5}, PDE, and LUC tasks. The results reveals criti-
951 cal insights into the performance-efficiency trade-offs inherent in each train-
952 ing approach. Pretrain-finetune emerges as the optimal strategy, achieving
953 the highest average accuracy (0.7409) with computational efficiency (16.7

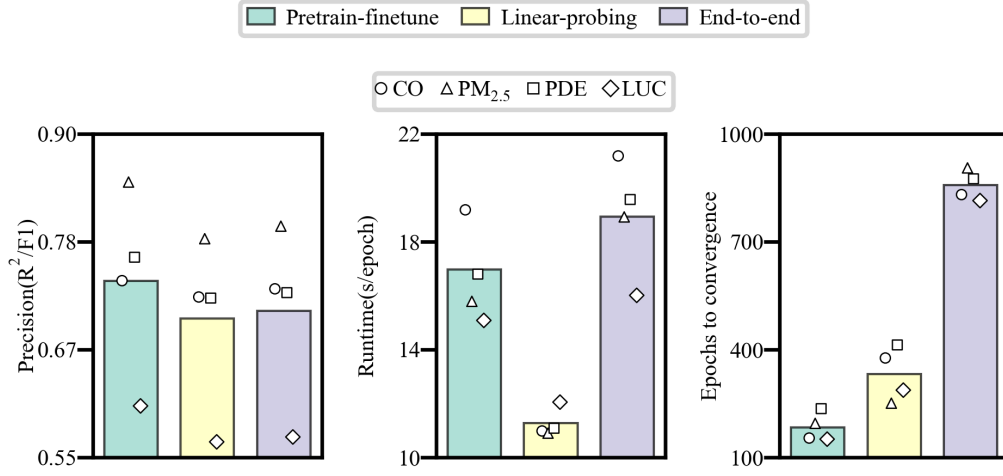


Figure 13: Comparative analysis of three training paradigms across prediction accuracy, runtime per epoch, and epochs to convergence across urban tasks. Bars represent average performance across all urban tasks, while scattered points show individual task performance (circles: CO prediction, triangles: PM_{2.5} prediction, squares: Population Density Estimation, diamonds: Land Use Classification).

s/epoch, 184 epochs). This paradigm preserves lspatial reasoning capabilities through selective parameter adaptation, making it ideal for resource-constrained urban monitoring applications. Linear probing shows the fastest per-epoch computation (11.3s) but suffers from limited representation adaptability. The frozen encoder prevents overfitting but results in systematic accuracy degradation, particularly in complex regression tasks. Despite faster iterations, it requires nearly twice as many epochs to converge (348 vs 184), offsetting its computational advantage. End-to-end training incurs prohibitive costs (18.9s/epoch, 857 epochs to converge) while achieving only marginal improvements over linear probing and falling 2.53% short of pretrain-finetune, making it viable only with abundant resources. Overall, the results validate

965 pretrain-finetune as the optimal training paradigm, effectively balancing pre-
966 dictive performance with computational efficiency.

967 5.4. Limitations and future directions

968 While UrbanMMCL demonstrates significant advances in urban repre-
969 sentation learning, several limitations warrant acknowledgment and present
970 opportunities for future research.

971 First, our framework’s reliance on high-quality textual descriptions gen-
972 erated by BLIP-2 introduces a potential bottleneck, as variations in text
973 generation quality across different urban scenes could lead to inconsistent
974 performance, particularly in challenging scenarios where visual content is
975 ambiguous or degraded. Future work should explore more robust text gen-
976 eration methods or develop alternative approaches to incorporate semantic
977 information less dependent on generative model.

978 Second, the static nature of our graph construction methods may not
979 fully capture dynamic temporal patterns inherent in urban systems, such as
980 daily traffic patterns or seasonal environmental changes. Additionally, the
981 choice of graph construction criteria may not be optimal for all urban tasks.
982 Future directions should explore temporal modeling capabilities, dynamic
983 graph learning approaches, and task-adaptive graph construction strategies.

984 Finally, while our framework demonstrates cross-city transferability from
985 Shenzhen to Beijing and Chengdu, complete zero-shot generalization remains
986 limited. The graph structure components require re-initialization and adap-
987 tation for city-specific spatial relationships, indicating that spatial modeling
988 still needs localized fine-tuning. Future research should investigate develop-
989 ing fully generalizable urban foundation models that can achieve complete

990 zero-shot inference without requiring any component re-training, potentially
991 through learning universal spatial relationship patterns or developing city-
992 agnostic graph construction strategies that can adapt automatically to new
993 urban environments.

994 6. Conclusions

995 This paper presents UrbanMMCL, a novel self-supervised dual-stage con-
996 trastive learning framework that advances urban representation learning through
997 innovative integration of multi-modal fusion and adaptive graph learning.
998 Our approach establishes a comprehensive pre-training paradigm that learns
999 generalizable urban representations without requiring task-specific labels, ad-
1000 dressing the critical challenge of limited annotated data in urban analysis.

1001 Comprehensive experimental validation demonstrates that UrbanMMCL
1002 consistently outperforms state-of-the-art methods across environmental mon-
1003 itoring, population estimation, and land use classification tasks. Cross-city
1004 transfer experiments further validate the generalizability of our learned rep-
1005 resentations across different urban environments. The framework’s success
1006 stems from its principled integration of RSI, SVI, textual descriptions, and
1007 geographical coordinates through contrastive learning, while adaptive graph
1008 learning captures dynamic inter-regional relationships that static approaches
1009 cannot model.

1010 UrbanMMCL represents a significant advancement toward urban founda-
1011 tion models by demonstrating how multi-modal pre-training can learn trans-
1012 ferable urban knowledge that generalizes across different tasks and cities.
1013 This work bridges the gap between domain-specific urban analysis tools and

1014 the broader vision of unified urban AI systems, laying the groundwork for
1015 more comprehensive urban foundation models that can support evidence-
1016 based urban planning, sustainable development, and smart city initiatives at
1017 unprecedented scale and sophistication.

1018 **CRedit authorship contribution statement**

1019 **Jinzhou Cao:** Methodology, Conceptualization, Investigation, Funding
1020 acquisition, Writing - original draft, Writing – review & editing. **Jiashi**
1021 **Chen:** Formal analysis, Data curation, Visualization, Writing - original
1022 draft, Writing – review & editing. **Xiangxu Wang:** Writing - original
1023 draft, Writing – review & editing. **Weiming Huang:** Writing – review &
1024 editing. **Dongsheng Chen:** Writing – review & editing. **Tianhong Zhao:**
1025 Writing – review & editing. **Wei Tu:** Investigation, Writing – review &
1026 editing. **Qingquan Li:** Resources, Writing – review & editing.

1027 **Acknowledgements**

1028 This research was supported in part by Shenzhen Science and Technology
1029 Program (No. JCYJ20240813113300001, 20231127180406001); the National
1030 Natural Science Foundation of China (No. 42401553); Natural Science Foun-
1031 dation of Top Talent of SZTU (No. GDRC202415).

1032 **References**

1033 Bai, L., Huang, W., Zhang, X., Du, S., Cong, G., Wang, H., Liu, B., 2023.
1034 Geographic mapping with unsupervised multi-modal representation learn-

1035 ing from VHR images and POIs. *ISPRS Journal of Photogrammetry and*
1036 *Remote Sensing* 201, 193–208.

1037 Bai, L., Zhang, X., Wang, H., Du, S., 2025. Integrating remote sensing
1038 with OpenStreetMap data for comprehensive scene understanding through
1039 multi-modal self-supervised learning. *Remote Sensing of Environment* 318,
1040 114573.

1041 Bao, H., Wang, W., Dong, L., Liu, Q., Mohammed, O.K., Aggarwal, K.,
1042 Som, S., Piao, S., Wei, F., 2022. VLMo: Unified vision-language pre-
1043 training with mixture-of-modality-experts, in: *Advances in Neural Infor-*
1044 *mation Processing Systems*, pp. 32897–32912.

1045 Belkin, M., Niyogi, P., 2001. Laplacian eigenmaps and spectral techniques for
1046 embedding and clustering, in: *Advances in Neural Information Processing*
1047 *Systems*.

1048 Cao, J., Wang, X., Chen, G., Tu, W., Shen, X., Zhao, T., Chen, J., Li, Q.,
1049 2025a. Disentangling the hourly dynamics of mixed urban function: A
1050 multimodal fusion perspective using dynamic graphs. *Information Fusion*
1051 117, 102832.

1052 Cao, J., Wang, X., Chen, J., Tu, W., Li, Z., Yang, X., Zhao, T., Li, Q., 2025b.
1053 Urban representation learning for fine-grained economic mapping: A semi-
1054 supervised graph-based approach. *ISPRS Journal of Photogrammetry and*
1055 *Remote Sensing* 226, 317–331.

1056 Cao, J., Wang, X., Chen, J., Zhang, B., Ma, Y., Zhao, T., 2025c.
1057 SemiGPS: GraphGPS-based Semi-supervised Graph Learning for Sector-

1058 Specific GDP Mapping, in: ICASSP 2025 - 2025 IEEE International Con-
1059 ference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5.

1060 Cao, R., Tu, W., Yang, C., Li, Q., Liu, J., Zhu, J., Zhang, Q., Li, Q., Qiu,
1061 G., 2020. Deep learning-based remote and social sensing data fusion for
1062 urban region function recognition. *ISPRS Journal of Photogrammetry and*
1063 *Remote Sensing* 163, 82–97.

1064 Cepeda, V.V., Nayak, G.K., Shah, M., 2023. GeoCLIP: Clip-inspired align-
1065 ment between locations and images for effective worldwide geo-localization,
1066 in: *Proceedings of the 37th International Conference on Neural Informa-*
1067 *tion Processing Systems*, Curran Associates Inc., Red Hook, NY, USA. pp.
1068 8690–8701.

1069 Chan, W., Ren, Q., 2023. Region-Wise Attentive Multi-View Representa-
1070 tion Learning For Urban Region Embedding, in: *Proceedings of the 32nd*
1071 *ACM International Conference on Information and Knowledge Manage-*
1072 *ment*, Association for Computing Machinery, New York, NY, USA. pp.
1073 3763–3767.

1074 Chen, T., Kornblith, S., Norouzi, M., Hinton, G., 2020. A simple frame-
1075 work for contrastive learning of visual representations, in: *International*
1076 *conference on machine learning*, PmLR. pp. 1597–1607.

1077 Chen, X., Xie, S., He, K., 2021. An empirical study of training self-supervised
1078 vision transformers, in: *Proceedings of the IEEE/CVF international con-*
1079 *ference on computer vision*, pp. 9640–9649.

1080 Chen, Y., Huang, W., Zhao, K., Jiang, Y., Cong, G., 2025. Self-supervised
1081 representation learning for geospatial objects: A survey. *Information Fu-*
1082 *sion* 123, 103265.

1083 Dai, G., Yi, W., Cao, J., Gong, Z., Fu, X., Zhang, B., 2025. CRRL: Con-
1084 trastive Region Relevance Learning Framework for Cross-city Traffic Pre-
1085 diction. *Information Fusion* 122, 103215.

1086 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Un-
1087 terthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkor-
1088 eit, J., Houlsby, N., 2021. An image is worth 16x16 words: Transformers
1089 for image recognition at scale, in: *International Conference on Learning*
1090 *Representations*.

1091 Gao, J., Li, P., Chen, Z., Zhang, J., 2020. A Survey on Deep Learning for
1092 Multimodal Data Fusion. *Neural Computation* 32, 829–864.

1093 Grover, A., Leskovec, J., 2016. Node2vec: Scalable feature learning for net-
1094 works, in: *Proceedings of the 22nd ACM SIGKDD International Confer-*
1095 *ence on Knowledge Discovery and Data Mining*, New York, NY, USA. pp.
1096 855–864.

1097 Guan, Q., Wang, J., Ren, S., Gao, H., Liang, Z., Wang, J., Yao, Y., 2024.
1098 Predicting short-term pm2. 5 concentrations at fine temporal resolutions
1099 using a multi-branch temporal graph convolutional neural network. *Inter-*
1100 *national Journal of Geographical Information Science* 38, 778–801.

1101 Hao, X., Chen, W., Yan, Y., Zhong, S., Wang, K., Wen, Q., Liang, Y.,
1102 2025. UrbanVLP: Multi-granularity vision-language pretraining for urban

1103 socioeconomic indicator prediction. Proceedings of the AAAI Conference
1104 on Artificial Intelligence 39, 28061–28069.

1105 Hassani, K., Khasahmadi, A.H., 2020. Contrastive multi-view representa-
1106 tion learning on graphs, in: International conference on machine learning,
1107 PMLR. pp. 4116–4126.

1108 He, J., Huang, B., 2025. Estimating global anthropogenic CO2 emissions
1109 through satellite observations. Environmental Research 279, 121767.

1110 He, K., Fan, H., Wu, Y., Xie, S., Girshick, R., 2020. Momentum contrast
1111 for unsupervised visual representation learning, in: Proceedings of the
1112 IEEE/CVF conference on computer vision and pattern recognition, pp.
1113 9729–9738.

1114 He, L., Cheng, D., Zhang, G., Zhang, S., 2025. Leveraging long-range nodes
1115 in multi-view graph contrastive learning. Information Fusion 122, 103186.

1116 Huang, W., Wang, J., Cong, G., 2024. Zero-shot urban function inference
1117 with street view images through prompting a pretrained vision-language
1118 model. International Journal of Geographical Information Science 38,
1119 1414–1442.

1120 Huang, W., Zhang, D., Mai, G., Guo, X., Cui, L., 2023. Learning urban
1121 region representations with POIs and hierarchical graph infomax. ISPRS
1122 Journal of Photogrammetry and Remote Sensing 196, 134–145.

1123 Jean, N., Burke, M., Xie, M., Davis, W.M., Lobell, D.B., Ermon, S., 2016.
1124 Combining satellite imagery and machine learning to predict poverty. Sci-
1125 ence 353, 790–794.

1126 Jenkins, P., Farag, A., Wang, S., Li, Z., 2019. Unsupervised representation
1127 learning of spatial data via multimodal embedding, in: Proceedings of
1128 the 28th ACM International Conference on Information and Knowledge
1129 Management, Association for Computing Machinery, New York, NY, USA.
1130 pp. 1993–2002.

1131 Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung,
1132 Y.H., Li, Z., Duerig, T., 2021. Scaling up visual and vision-language repre-
1133 sentation learning with noisy text supervision, in: International conference
1134 on machine learning, PMLR. pp. 4904–4916.

1135 Ju, W., Fang, Z., Gu, Y., Liu, Z., Long, Q., Qiao, Z., Qin, Y., Shen, J., Sun,
1136 F., Xiao, Z., Yang, J., Yuan, J., Zhao, Y., Wang, Y., Luo, X., Zhang, M.,
1137 2024. A Comprehensive Survey on Deep Graph Representation Learning.
1138 Neural Networks 173, 106207.

1139 Khoshraftar, S., An, A., 2024. A survey on graph representation learning
1140 methods. ACM Transactions on Intelligent Systems and Technology 15,
1141 19:1–19:55.

1142 Klemmer, K., Rolf, E., Robinson, C., Mackey, L., Rußwurm, M., 2025. Sat-
1143 CLIP: Global, general-purpose location embeddings with satellite imagery,
1144 in: Proceedings of the 39th Annual AAAI Conference on Artificial Intelli-
1145 gence: AAAI-25 Technical Tracks, pp. 4347–4355.

1146 Li, T., Xi, Y., Wang, H., Li, Y., Tarkoma, S., Hui, P., 2023a. Learning Rep-
1147 resentations of Satellite Imagery by Leveraging Point-of-Interests. ACM
1148 Transactions on Intelligent Systems and Technology 14, 1–32.

- 1149 Li, Y., Huang, W., Cong, G., Wang, H., Wang, Z., 2023b. Urban Region Rep-
1150 resentation Learning with OpenStreetMap Building Footprints, in: Pro-
1151 ceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery
1152 and Data Mining, Association for Computing Machinery, New York, NY,
1153 USA. pp. 1363–1373.
- 1154 Li, Y., Yang, M., Zhang, Z., 2019. A Survey of Multi-View Representation
1155 Learning. *IEEE Transactions on Knowledge and Data Engineering* 31,
1156 1863–1883.
- 1157 Li, Z., He, W., Cheng, M., Hu, J., Yang, G., Zhang, H., 2023c. Sinolc-1: the
1158 first 1 m resolution national-scale land-cover map of china created with
1159 a deep learning framework and open-access data. *Earth System Science*
1160 *Data* 15, 4749–4780.
- 1161 Li, Z., Huang, W., Zhao, K., Yang, M., Gong, Y., Chen, M., 2024. Urban
1162 region embedding via multi-view contrastive prediction, in: Proceedings of
1163 the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-
1164 Sixth Conference on Innovative Applications of Artificial Intelligence and
1165 Fourteenth Symposium on Educational Advances in Artificial Intelligence,
1166 AAAI Press. pp. 8724–8732.
- 1167 Li, Z., Zhang, H., Lu, F., Xue, R., Yang, G., Zhang, L., 2022. Breaking the
1168 resolution barrier: A low-to-high network for large-scale high-resolution
1169 land-cover mapping using low-resolution labels. *ISPRS Journal of Pho-*
1170 *togrammetry and Remote Sensing* 192, 244–267.
- 1171 Liu, F., Chen, D., Guan, Z., Zhou, X., Zhu, J., Ye, Q., Fu, L., Zhou, J.,

1172 2024. RemoteCLIP: A vision language foundation model for remote sens-
1173 ing. *IEEE Transactions on Geoscience and Remote Sensing* 62, 1–16.

1174 Liu, S., Zhang, T., Fu, N., Huang, Y., 2025. Fine-grained graph representa-
1175 tion learning for heterogeneous mobile networks with attentive fusion and
1176 contrastive learning. *Proceedings of the AAAI Conference on Artificial*
1177 *Intelligence* 39, 18933–18942.

1178 Liu, Y., Zhang, X., Ding, J., Xi, Y., Li, Y., 2023. Knowledge-infused Con-
1179 trastive Learning for Urban Imagery-based Socioeconomic Prediction, in:
1180 *Proceedings of the ACM Web Conference 2023*, Association for Computing
1181 Machinery, New York, NY, USA. pp. 4150–4160.

1182 Luo, Y., Chung, F.I., Chen, K., 2022. Urban Region Profiling via Multi-
1183 Graph Representation Learning, in: *Proceedings of the 31st ACM Interna-*
1184 *tional Conference on Information & Knowledge Management*, Association
1185 for Computing Machinery, New York, NY, USA. pp. 4294–4298.

1186 Perozzi, B., Al-Rfou, R., Skiena, S., 2014. DeepWalk: Online Learning
1187 of Social Representations, in: *Proceedings of the 20th ACM SIGKDD*
1188 *International Conference on Knowledge Discovery and Data Mining*, pp.
1189 701–710.

1190 Qin, Q., Ai, T., Xu, S., Zhang, Y., Huang, W., Du, M., Li, S., 2025. Learning
1191 dual context aware POI representations for geographic mapping. *Inter-*
1192 *national Journal of Applied Earth Observation and Geoinformation* 142,
1193 104683.

1194 Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S.,
1195 Sastry, G., Aspell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.,
1196 2021a. Learning transferable visual models from natural language super-
1197 vision, in: Proceedings of the 38th International Conference on Machine
1198 Learning, PMLR. pp. 8748–8763.

1199 Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S.,
1200 Sastry, G., Aspell, A., Mishkin, P., Clark, J., et al., 2021b. Learning
1201 transferable visual models from natural language supervision, in: Interna-
1202 tional conference on machine learning, PmLR. pp. 8748–8763.

1203 Shen, X., Wang, H., Wei, B., Cao, J., 2023. Real-time scene classification of
1204 unmanned aerial vehicles remote sensing image based on Modified Ghost-
1205 Net. PLOS ONE 18, e0286873.

1206 Sun, F.Y., Hoffman, J., Verma, V., Tang, J., 2020a. Infograph: Unsuper-
1207 vised and semi-supervised graph-level representation learning via mutual
1208 information maximization, in: International Conference on Learning Rep-
1209 resentations.

1210 Sun, F.Y., Hoffmann, J., Verma, V., Tang, J., 2020b. InfoGraph: Unsuper-
1211 vised and semi-supervised graph-level representation learning via mutual
1212 information maximization. [arXiv:1908.01000](https://arxiv.org/abs/1908.01000).

1213 Sun, Z., Jiao, H., Wu, H., Peng, Z., Liu, L., 2021. Block2vec: An Ap-
1214 proach for Identifying Urban Functional Regions by Integrating Sentence
1215 Embedding Model and Points of Interest. ISPRS International Journal of
1216 Geo-Information 10, 339.

- 1217 Suresh, S., Li, P., Hao, C., Neville, J., 2021. Adversarial graph augmentation
1218 to improve graph contrastive learning, in: Proceedings of the 35th Inter-
1219 national Conference on Neural Information Processing Systems, Curran
1220 Associates Inc., Red Hook, NY, USA. pp. 15920–15933.
- 1221 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N.,
1222 Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need, in: Advances in
1223 Neural Information Processing Systems.
- 1224 Veličković, P., Fedus, W., Hamilton, W.L., Liò, P., Bengio, Y., Hjelm, R.D.,
1225 2018. Deep graph infomax, in: International Conference on Learning Rep-
1226 resentations.
- 1227 Vivanco, V., Nayak, G.K., Shah, M., 2023. Geoclip: Clip-inspired alignment
1228 between locations and images for effective worldwide geo-localization, in:
1229 Advances in Neural Information Processing Systems.
- 1230 Wang, X., Cao, J., Zhao, T., Zhang, B., Chen, G., Li, Z., Chen, H., Tu, W.,
1231 Li, Q., 2026. ST-camba: A decoupled-free spatiotemporal graph fusion
1232 state space model with linear complexity for efficient traffic forecasting.
1233 Information Fusion 126, 103495.
- 1234 Wang, X., Chen, H., Liu, Y., 2024. Learning place representations from
1235 spatial interactions. International Journal of Geographical Information
1236 Science 38, 1065–1090.
- 1237 Wang, X., Cheng, T., Law, S., Zeng, Z., Yin, L., Liu, J., 2025. Multi-
1238 modal contrastive learning of urban space representations from POI data.
1239 Computers, Environment and Urban Systems 120, 102299.

- 1240 Wang, Z., Li, H., Rajagopal, R., 2020. Urban2Vec: Incorporating Street View
1241 Imagery and POIs for Multi-Modal Urban Neighborhood Embedding, in:
1242 Proceedings of the AAAI Conference on Artificial Intelligence, pp. 1013–
1243 1020.
- 1244 Wei, J., Li, Z., 2024. GlobalHighCO: Global daily seamless 1 km ground- level
1245 CO dataset over land (2018–present). doi:[10.5281/zenodo.14207363](https://doi.org/10.5281/zenodo.14207363).
- 1246 Wei, J., Li, Z., Lyapustin, A., Wang, J., Dubovik, O., Schwartz, J., Sun,
1247 L., Li, C., Liu, S., Zhu, T., 2023. First close insight into global daily
1248 gapless 1 km PM2.5 pollution, variability, and health impact. Nature
1249 Communications 14, 8349.
- 1250 Weng, X., Pang, C., Xia, G.S., 2025. Vision-language modeling meets remote
1251 sensing: Models, datasets, and perspectives. IEEE Geoscience and Remote
1252 Sensing Magazine , 2–50.
- 1253 Wu, L., Lin, H., Tan, C., Gao, Z., Li, S.Z., 2023. Self-supervised learning
1254 on graphs: Contrastive, generative, or predictive. IEEE Transactions on
1255 Knowledge and Data Engineering 35, 4216–4235.
- 1256 Wu, S., Yan, X., Fan, X., Pan, S., Zhu, S., Zheng, C., Cheng, M., Wang,
1257 C., 2022. Multi-Graph Fusion Networks for Urban Region Embedding, in:
1258 Proceedings of the Thirty-First International Joint Conference on Artificial
1259 Intelligence.
- 1260 Wu, Z., Xiong, Y., Yu, S.X., Lin, D., 2018. Unsupervised feature learning
1261 via non-parametric instance discrimination, in: Proceedings of the IEEE
1262 conference on computer vision and pattern recognition, pp. 3733–3742.

- 1263 Xi, Y., Li, T., Wang, H., Li, Y., Tarkoma, S., Hui, P., 2022. Beyond the
1264 First Law of Geography: Learning Representations of Satellite Imagery by
1265 Leveraging Point-of-Interests, in: Proceedings of the ACM Web Conference
1266 2022, Association for Computing Machinery, New York, NY, USA. pp.
1267 3308–3316.
- 1268 Xu, Y., Jin, S., Chen, Z., Xie, X., Hu, S., Xie, Z., 2022. Application of a
1269 graph convolutional network with visual and semantic features to classify
1270 urban scenes. *International Journal of Geographical Information Science*
1271 36, 2009–2034.
- 1272 Xu, Z., Zhou, X., 2024. CGAP: Urban Region Representation Learn-
1273 ing with Coarsened Graph Attention Pooling, in: Proceedings of the
1274 Thirty-Third International Joint Conference on Artificial Intelligence, In-
1275 ternational Joint Conferences on Artificial Intelligence Organization, Jeju,
1276 South Korea. pp. 7518–7526.
- 1277 Yan, Y., Wen, H., Zhong, S., Chen, W., Chen, H., Wen, Q., Zimmermann,
1278 R., Liang, Y., 2024. UrbanCLIP: Learning Text-enhanced Urban Region
1279 Profiling with Contrastive Language-Image Pretraining from the Web, in:
1280 Proceedings of the ACM Web Conference 2024, Association for Computing
1281 Machinery, New York, NY, USA. pp. 4006–4017.
- 1282 Yeh, C., Perez, A., Driscoll, A., Azzari, G., Tang, Z., Lobell, D., Ermon, S.,
1283 Burke, M., 2020. Using publicly available satellite imagery and deep learn-
1284 ing to understand economic well-being in Africa. *Nature Communications*
1285 11, 2583.

1286 Yong, X., Zhou, X., 2024. MuseCL: Predicting Urban Socioeconomic Indi-
1287 cators via Multi-Semantic Contrastive Learning, in: Proceedings of the
1288 Thirty-Third International Joint Conference on Artificial Intelligence, In-
1289 ternational Joint Conferences on Artificial Intelligence Organization. pp.
1290 7536–7544.

1291 You, Y., Chen, T., Shen, Y., Wang, Z., 2021. Graph Contrastive Learn-
1292 ing Automated, in: Proceedings of the 38th International Conference on
1293 Machine Learning, PMLR. pp. 12121–12132.

1294 You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., Shen, Y., 2020. Graph
1295 contrastive learning with augmentations. *Advances in neural information*
1296 *processing systems* 33, 5812–5823.

1297 Zhai, W., Bai, X., Shi, Y., Han, Y., Peng, Z.R., Gu, C., 2019. Beyond
1298 Word2vec: An approach for urban functional region extraction and identi-
1299 fication by combining Place2vec and POIs. *Computers, Environment and*
1300 *Urban Systems* 74, 1–12.

1301 Zhang, F., Wu, L., Zhu, D., Liu, Y., 2019. Social sensing from street-level
1302 imagery: A case study in learning spatio-temporal urban mobility patterns.
1303 *ISPRS Journal of Photogrammetry and Remote Sensing* 153, 48–58.

1304 Zhang, L., Long, C., Cong, G., 2023a. Region Embedding With Intra and
1305 Inter-View Contrastive Learning. *IEEE Transactions on Knowledge and*
1306 *Data Engineering* 35, 9031–9036.

1307 Zhang, M., Li, T., Li, Y., Hui, P., 2020. Multi-View Joint Graph Repre-
1308 sentation Learning for Urban Region Embedding, in: Proceedings of the

- 1309 Twenty-Ninth International Joint Conference on Artificial Intelligence, In-
1310 ternational Joint Conferences on Artificial Intelligence Organization, Yoko-
1311 hama, Japan. pp. 4431–4437.
- 1312 Zhang, Q., Huang, C., Xia, L., Wang, Z., Yiu, S., Han, R., 2023b. Spatial-
1313 temporal graph learning with adversarial contrastive adaptation, in: Pro-
1314 ceedings of the 40th International Conference on Machine Learning, pp.
1315 41151–41163.
- 1316 Zhang, W., Han, J., Xu, Z., Ni, H., Lyu, T., Liu, H., Xiong, H., 2025.
1317 Towards urban general intelligence: A review and outlook of urban foun-
1318 dation models. [arXiv:2402.01749](https://arxiv.org/abs/2402.01749).
- 1319 Zhang, X., Gong, Y., Zhang, C., Wu, X., Guo, Y., Lu, W., Zhao, L., Dong,
1320 X., 2023c. Spatio-temporal fusion and contrastive learning for urban flow
1321 prediction. *Knowledge-Based Systems* 282, 111104.
- 1322 Zhang, Y., Huang, W., Yao, Y., Gao, S., Cui, L., Yan, Z., 2024a. Ur-
1323 ban region representation learning with human trajectories: A multi-view
1324 approach incorporating transition, spatial, and temporal perspectives. *GI-
1325 Science & Remote Sensing* 61, 2387392.
- 1326 Zhang, Y., Li, Y., Zhang, F., 2024b. Multi-level urban street representation
1327 with street-view imagery and hybrid semantic graph. *ISPRS Journal of
1328 Photogrammetry and Remote Sensing* 218, 19–32.
- 1329 Zhang, Y., Xu, Y., Cui, L., Yan, Z., 2023d. Multi-view graph contrastive
1330 learning for urban region representation, in: 2023 International Joint Con-
1331 ference on Neural Networks (IJCNN), pp. 1–8.

- 1332 Zhao, T., Liang, X., Biljecki, F., Tu, W., Cao, J., Li, X., Yi, S., 2025. Quan-
1333 tifying seasonal bias in street view imagery for urban form assessment: A
1334 global analysis of 40 cities. *Computers, Environment and Urban Systems*
1335 120, 102302.
- 1336 Zhou, S., He, D., Chen, L., Shang, S., Han, P., 2023a. Heterogeneous Region
1337 Embedding with Prompt Learning, in: *Proceedings of the AAAI Confer-*
1338 *ence on Artificial Intelligence*, pp. 4981–4989.
- 1339 Zhou, W., Sun, F., Jiang, Q., Cong, R., Hwang, J.N., 2023b. Wavenet:
1340 Wavelet network with knowledge distillation for rgb-t salient object detec-
1341 tion. *IEEE Transactions on Image Processing* 32, 3027–3039.
- 1342 Zhou, W., Zhu, Y., Lei, J., Wan, J., Yu, L., 2021. Ccafnet: Crossflow and
1343 cross-scale adaptive fusion network for detecting salient objects in rgb-d
1344 images. *IEEE Transactions on Multimedia* 24, 2192–2204.
- 1345 Zhou, Y., Huang, Y., 2018. DeepMove: Learning Place Representations
1346 through Large Scale Movement Data, in: *2018 IEEE International Con-*
1347 *ference on Big Data (Big Data)*, pp. 2403–2412.
- 1348 Zhu, Y., Xu, Y., Yu, F., Liu, Q., Wu, S., Wang, L., 2020. Deep graph
1349 contrastive representation learning. *arXiv preprint arXiv:2006.04131* .
- 1350 Zhu, Y., Xu, Y., Yu, F., Liu, Q., Wu, S., Wang, L., 2021. Graph contrastive
1351 learning with adaptive augmentation, in: *Proceedings of the web confer-*
1352 *ence 2021*, pp. 2069–2080.
- 1353 Zou, X., Yan, Y., Hao, X., Hu, Y., Wen, H., Liu, E., Zhang, J., Li, Y., Li, T.,
1354 Zheng, Y., Liang, Y., 2025. Deep learning for cross-domain data fusion in

1355 urban computing: Taxonomy, advances, and outlook. Information Fusion
1356 113, 102606.