URBANMLLM: JOINT LEARNING OF CROSS-VIEW IM AGERY FOR URBAN UNDERSTANDING

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

Paper under double-blind review

ABSTRACT

Multimodal large language models (MLLMs) have exhibited remarkable capabilities for performing complex vision-language tasks in various domains. Currently, MLLMs based on urban imagery in urban studies are only developed focusing on remote sensing imagery. However, except for the macroscopic information from remote sensing imagery, effective urban understanding also requires detailed appearance information of urban zones from street-view imagery, which is largely overlooked by existing MLLMs. The primary challenges of developing such a versatile urban MLLM are twofold. Firstly, it needs a large-scale corpus with well-organized, cross-view urban imagery paired with corresponding text for cross-modal training. Secondly, traditional MLLMs typically learn image-text pairs independently, hard to support joint modeling of cross-view urban imagery. To address these challenges, in this work, we propose UrbanM-LLM, a novel MLLM that jointly learns from remote sensing and street-view imagery to harness their complementary information. We first collect a large-scale dataset containing satellite-view and street-view imagery along with their geotags and annotated texts. Technically, we propose a brand MLLM architecture with a cross-view perceiver to explicitly connect visual information of cross-view urban imagery. We also introduce a novel pre-training paradigm based on structural interleaved urban image-text documents integrating satellite-view, street-view imagery and related textual descriptions. This approach encourages the model to implicitly learn the relationships between different types of urban imagery, enhancing the understanding in each domain. We evaluate our model on a comprehensive benchmark comprising 13 diverse urban understanding tasks across satellite-view, street-view, and cross-view domains. These tasks include scene classification, object reasoning, spatial relationship reasoning, geo-localization, landmark reasoning, and indicator prediction, providing a robust assessment of the model's capabilities. Extensive experiments demonstrate that UrbanMLLM achieves an average of 27.3% and 25.5% performance improvement compared with the best open-sourced and closed-sourced MLLMs, respectively. Moreover, we thoroughly study the impact of different pre-training data choices and model scales on performance, offering practical insights for effective MLLM design. The proposed UrbanMLLM offers a scalable and versatile solution for understanding urban environments.

041 042 043

044

039

040

1 INTRODUCTION

Urban imagery has been widely used for understanding cities in terms of urban spatial structure, functionality, and socio-economic status. The advancement in computer vision and multimodal learning has driven the utilization of multimodal urban data for urban understanding tasks, such as scene classification (Kuckreja et al., 2024; Mall et al., 2024), scene geo-localization (Vivanco Cepeda et al., 2024; Xu et al., 2024), urban indicator prediction (Fan et al., 2023; Hao et al., 2024), etc. More recently, benefiting from the impressive performance and generalizability of large language models (LLMs), multimodal large language models (MLLMs) have shown great potential for effectively solving distinct multimodal tasks in a "one-for-all" manner.

In the urban study area, there has been a line of works using MLLMs to tackle urban understanding tasks, while predominantly focusing on remote sensing (Kuckreja et al., 2024; Luo et al., 2024; Bazi

054 et al., 2024; Muhtar et al., 2024). Remote sensing imagery provides a macroscopic and comprehen-055 sive overview of how the city's functional zones are laid out while lacking some detailed contexts 056 of the urban elements. Therefore, the existing MLLMs in urban studies can only deal with high-057 level urban understanding tasks such as land use classification and region captions. By comparison, 058 street-view imagery captures a more fine-grained appearance of urban zones, providing complementary information such as building facades and heights. However, no existing MLLMs in urban science has explored integrating street-view imagery to enhance urban understanding. To enhance 060 MLLMs' comprehensive understanding of urban environments, satellite imagery, which captures 061 large-scale spatial layouts, and street-view imagery, which provides ground-level details, should be 062 integrated and jointly learned within a unified model. 063

Achieving this goal needs to address two primary challenges. The first is the lack of well-organized multimodal urban data. Currently, publicly available urban datasets do not pair cross-view imagery with corresponding textual annotations, making them unable to support multimodal learning of MLLMs. The second challenge comes from the joint learning paradigm of connecting remote sensing and street-view imagery. In conventional MLLM frameworks (Liu et al., 2024c; Chen et al., 2024), visual features from cross-view imagery are encoded separately and aligned only with their respective annotated texts. These approaches fail to capture the complementary relationships between satellite and street-view imagery, leaving their information isolated.

In this work, we address the above two challenges and introduce a novel urban MLLM jointly 072 learning from remote sensing and street-view imagery and associated textual data. Specifically, our 073 efforts focus on both data contribution and methodology innovation. Given that existing publicly 074 available datasets about urban imagery (Luo et al., 2024; Astruc et al., 2024) commonly lack paired 075 cross-view images and large-scale annotated text information, we first collect a large-scale multi-076 modal urban imagery dataset. Our dataset covers the whole United States, comprising paired-up 077 satellite-view and street-view images, together with geotags and annotated textual descriptions. We propose two key designs to break the visual knowledge isolation to facilitate the mutual learning 079 of cross-view urban imagery. The first one is about the model architecture, where we propose a cross-view perceiver module that bridges the paired-up satellite-view and street-view visual features 081 through a cross-attention mechanism. This design explicitly facilitates the exchange of information between the region-level context of satellite imagery and the fine-grained appearance details of street-view imagery. For example, the injection of street-view information to the satellite-view 083 encoding can provide more detailed urban region context. The second part is a novel interleaved pre-084 training paradigm to enhance the mutual learning between cross-view imagery. In detail, we design 085 coherent image-text documents that interleave the satellite-view image with matched street-view images and associated textual descriptions, forming a comprehensive profile of an urban region. Such 087 interleaved training corpus helps MLLMs implicitly learn the relationship between different-view 088 urban imagery via in-context learning. Through the explicit and implicit mutual learning between 089 cross-view urban imagery, our proposed UrbanMLLM is expected to overcome the visual isolation 090 issue, benefiting the comprehensive understanding of urban environments from diverse views. 091

For a comprehensive evaluation of urban understanding abilities for MLLMs, we build Urban-092 View Benchmark which includes 13 different tasks of urban perception (scene classification, geo-093 localization), reasoning (object reasoning, spatial relationship reasoning, landmark reasoning) and 094 prediction (indicator prediction) based on single-view or cross-view urban imagery. Extensive ex-095 periments on the benchmark validate the noticeable superiority of UrbanMLLM on a wide range of 096 urban understanding tasks, achieving an average 27.3% and 25.5% performance gain compared with 097 the best open-sourced and closed-sourced MLLMs, respectively. Moreover, we study the impact of 098 different pre-training data choices and model scales on the final performance, offering practical insights for effective MLLM design. This work serves as a foundational technique for addressing a wide range of urban-related tasks requiring comprehensive visual understanding capabilities. In 100 brief, our major contributions can be summarized as follows: 101

- 102
- 103 104

• We propose a brand MLLM architecture with a designed cross-perceiver module to facilitate crossfusion of the complementary visual context from satellite-view and street-view imagery.

105

We construct a novel interleaved pre-training corpus that links satellite and street-view imagery through geo-location relationships, and propose a training paradigm that implicitly promotes mutual learning between cross-view imagery.

Mathad Truna	Madal]	Data	Task			
Method Type	WIGGET	Satellite Image	Street View Image	Perception	Reasoning	Predicti	
	RemoteCLIP	✓	×	×	×	 Image: A second s	
CLIP-based	UrbanCLIP	✓	×	×	×	 Image: A second s	
	UrbanVLP	✓	✓	×	×	 Image: A second s	
	GeoChat	✓	×	 Image: A set of the set of the	✓	×	
MLI M based	LHRS-Bot	✓	×	 Image: A second s	 Image: A second s	×	
WILLWI-Dased	SkysenseGPT	✓	×	✓	 Image: A second s	×	
	UrbanMLLM	1	✓	1	1	 Image: A second s	

....

• We establish a comprehensive benchmark including 13 different urban understanding tasks based on single-view or cross-view urban imagery. Extensive experiments verify that our model achieves substantial improvement in urban understanding over both open-source and closedsourced MLLMs.

122 123 124

125 126

127

119

120

121

108

RELATED WORK 2

2.1 MULTIMODAL LARGE LANGUAGE MODELS (MLLMS)

128 With the rapid development and significant success of large language models (LLMs), recent re-129 search has focused on developing multimodal large language models (MLLMs) with the ability 130 to comprehend both visual and textual knowledge, enabling them to address complex visual rea-131 soning and understanding tasks. Existing MLLMs can be roughly divided into closed-source and 132 open-source models. Closed-source MLLMs are mostly built based on corresponding commercial LLMs, such as GPT-40 (Achiam et al., 2023), Gemini (Reid et al., 2024) and Qwen-VL (Bai et al., 133 2023). These models benefit from the large-scale and extensive training corpus, which have been 134 shown to exhibit powerful general multimodal understanding capabilities. By comparison, open-135 source MLLMs are usually smaller-scale, established by aligning a visual encoding branch to an 136 off-the-shelf LLM. Following an earlier work LLaVA (Liu et al., 2024c), there are mainly two di-137 rections to advance the performance of MLLMs. The first line of works explores more advanced 138 architectures such as introducing dual vision encoders (Li et al., 2024b), more sophisticated visual 139 adapters (Cha et al., 2024) and the mixture-of-expert (MoE) strategy (Li et al., 2024c). Another line 140 tries to uplift the performance of MLLMs with more beneficial pre-training data, such as interleaved 141 multimodal data (Lin et al., 2024) and synthetic data (McKinzie et al., 2024). By comparison, 142 our work introduces a novel MLLM architecture that integrates a cross-view perceiver module to 143 enhance cross-view information fusion and contribute a unique interleaved pre-training corpus for 144 MLLMs in urban areas.

145 146

147

2.2 MULTIMODAL MODELS FOR URBAN UNDERSTANDING

Understanding the urban environment usually requires multimodal information from diverse 148 sources, such as satellite-view images, street-view images, POI information and geo-locations, 149 etc. Existing methods in urban study can be categorized into two types: CLIP-based methods and 150 MLLM-based methods, as shown in Table 1. From the data aspect, existing methods based urban 151 imagery in urban study all focus on satellite images while overlooking using the street-view imagery 152 for urban understanding. CLIP-based methods are mostly developed based on the contrastive learn-153 ing strategy used in CLIP (Radford et al., 2021), such as training with satellite image-text pairs (Liu 154 et al., 2024a; Yan et al., 2024), street-view image-text pairs (Hao et al., 2024) and satellite-view and 155 street-view image pairs (Mall et al., 2024). These works can only deal with prediction tasks such as 156 indicator prediction via end-to-end fine-tuning, but fail to conduct perception and reasoning tasks. 157 Another line of research focuses on developing specialized MLLMs for problem-solving in the ur-158 ban domain. Existing models, such as GeoChat (Kuckreja et al., 2024), SkysenseGPT (Luo et al., 159 2024), H²RSVLM Pang et al. (2024), and EarthGPT (Zhang et al., 2024) only leverage remote sensing data including satellite images and annotated text for model learning. These models are capable 160 of handling remote sensing perception and reasoning tasks but fail to deal with prediction tasks such 161 as indicator predictions. However, relying solely on region-level knowledge is insufficient to capture the complexities of urban environments, thereby limiting their applications for a wide range of urban understanding tasks. In contrast, our work proposes a novel learning paradigm based on cross-view urban image-text data which is capable of solving both remote sensing and street-view tasks.

- 166 167 3 METHODOLOGY
- 168 169 3.1 OVERVIEW

170 Traditional MLLMs are usually trained with paired-up separate image-text data, resulting in the 171 knowledge isolation between different images. Such limitation constrains their ability of urban un-172 derstanding, which requires a holistic comprehension of urban imagery from diverse perspectives. 173 To overcome this issue, we propose two key designs to facilitate the comprehensive understanding on cross-view urban imagery. Firstly, we introduce a cross-view perceiver module in the MLLM 174 architecture, explicitly enabling satellite-view and street-view visual contexts to complement each 175 other. Secondly, we propose a novel interleaved pre-training paradigm leveraging structurally in-176 terleaved urban image-text contexts, integrating satellite-view and street-view imagery with corre-177 sponding textual descriptions. Training on such interleaved data enables MLLMs to learn relation-178 ships between cross-view urban imagery, leading to a more comprehensive understanding of urban 179 environments. We elaborate the MLLM architecture enhanced with cross-view fusion in Section 3.2, followed by the designed pre-training paradigm based on interleaved urban contexts in Section 3.3. 181

182 183

3.2 CROSS-VIEW FUSION-ENHANCED URBANMLLM

Current MLLMs in urban studies primarily focus on remote sensing tasks. These models are typically developed by directly fine-tuning general-purpose MLLMs (e.g., LLaVA) on satellite imagetext pairs. However, the effective urban understanding not only requires comprehending region-level knowledge from satellite-view imagery but also detailed contexts from street-view imagery. Unfortunately, this objective is unpromising to be achieved with the classical MLLM architecture, where images are individually encoded, failing to receive visual knowledge from relevant images.

190 Aiming to address the visual knowledge isolation issue, we introduce a cross-view perceiver module 191 $g_{\zeta}(\cdot)$ to promote the awareness of urban imagery from other views during the visual encoding pro-192 cess. The cross-view perceiver is shown in the Figure 1. It performs 4 steps: (1) cross-attention from satellite-view image embedding (as queries) to the street-view image embedding; (2) cross-attention 193 from street-view image embedding (as queries) to the satellite-view image embedding; (3) gating 194 module before the residual connection; (4) MLP for aligning the semantic space of text. When both 195 satellite-view and street-view images exist in the multimodal input, let I_{st} denote a satellite-view 196 image and $\{I_{sv}^i\}_{i=1}^n$ represent n paired street-view images. The satellite-view and street-view im-197 ages are firstly encoded by a pre-trained visual encoder $f_{\phi}(\cdot)$, resulting in visual features f_{st} and $\{f_{sv}^i\}_{i=1}^n$, respectively. When encoding each street-view image I_{sv}^i , we inject the matched satellite-199 view features, giving rise to the fused feature $e^i_{st \to sv}$ at the step of SI2SVI attn. For the satellite-view 200 image, we first conduct an average-pooling on the visual features of n paired street-view images then 201 fuse it with the satellite-view feature, obtaining the fused feature $e_{sv \to st}$ at the step of SVI2SI attn. 202 Next, the fused visual feature is adaptively combined with the original feature using a gating strategy 203 implemented with a one-layer MLP. The final visual embeddings \mathbf{V}_{sv}^{i} and \mathbf{V}_{st} are then obtained after 204 a visual adapter (two-layer MLP). The whole operation of the cross-view perceiver is following:

$$\boldsymbol{e}_{st \to sv}^{i} = \text{MLP}(\text{Softmax}(\frac{\boldsymbol{f}_{sv}^{i} \boldsymbol{f}_{st}}{\sqrt{d_{k}}}) \boldsymbol{f}_{st}), \tag{1}$$

$$\mathbf{W}_{sv}^{i} = \mathrm{MLP}(\mathrm{Gating}(\boldsymbol{e}_{st \to sv}^{i}) + \boldsymbol{f}_{sv}^{i}), \qquad (2)$$

$$e_{sv \to st} = \text{MLP}(\text{Softmax}(\frac{f_{st}f_{sv}}{\sqrt{d_k}})\tilde{f}_{sv}), \tag{3}$$

205 206

- $\mathbf{V}_{st} = \mathrm{MLP}(\mathrm{Gating}(\boldsymbol{e}_{sv \to st}) + \boldsymbol{f}_{st}), \tag{4}$
- where $\tilde{f}_{sv} = \text{Pooling}(\{f_{sv}^i\}_{i=1}^n).$
- 215 If there's only single-view imagery in the multimodal input, the cross-view perceiver module receives two identical single-view images as input. In this way, the visual embedding fed to the LLM



Figure 1: Architecture of the proposed UrbanMLLM. UrbanM- Figure 2: The pipeline of in-LLM employs cross-view perceiver to learn cross-view visual rep- terleaved image-text data conresentations.

struction.

backbone is enhanced by the visual context from another view of the same urban region, which possesses more comprehensive urban context.

INTERLEAVED URBAN CONTEXT-BASED PRE-TRAINING 3.3

Existing explorations on general MLLMs have demonstrated that training on interleaved data yields 242 superior performance than the traditional image-text pairs (Lin et al., 2024; McKinzie et al., 2024). 243 The interleaved structure fosters semantic connections between multiple images, enabling MLLMs 244 to better capture contextual relationships across images. This advantage aligns well with our objec-245 tive of jointly learning from cross-view urban imagery to enhance the comprehensive understanding 246 of urban visual knowledge. For instance, when predicting the geo-location of a satellite image, 247 region-level visual information alone may be insufficient. Supplementing detailed street-view information can provide the necessary contextual knowledge for more accurate predictions. 248

249 Motivated by this, we introduce an urban context-based interleaved training paradigm tailored for 250 urban understanding tasks. The core of this part is the construction of multimodal interleaved ur-251 ban data as the training corpus. We first collect a large scale satellite-view and street-view im-252 agery individually across the United States, and perform cross-view matching based on geotags 253 (including located county, longitude and latitude), creating a paired cross-view urban imagery set 254 $S = \{(I_{st}, I_{sv}^1, I_{sv}^2, ..., I_{sv}^n) | n \in \mathbb{Z}^+\}$. We then employ an advanced MLLM InternVL (Chen et al., 2024) with carefully crafted prompts to efficiently generate textual descriptions for each image. 255 Next, we link the cross-view images based on their geographical relationships and integrate their 256 corresponding textual descriptions and geotags, forming a comprehensive urban profile for each el-257 ement in S. An illustrative example is shown in Figure 16. Training on such interleaved multimodal 258 urban data benefits the MLLM to capture the relational knowledge between cross-view imagery, 259 facilitating comprehensive urban understanding by fully integrating contextual information. As-260 suming that the interleaved document contains K ordered urban images $I = \{I_k\}_{k=1}^K$ interleaved 261 with a T-length word sequence $\mathbf{w} = \{w_t\}_{t=1}^T$ tokenized by a θ -parameterized LLM. The k-th im-262 age is successively processed by a frozen visual encoder $f_{\phi}(\cdot)$ and the cross-view perceiver $g_{\zeta}(\cdot)$ 263 into L-length image tokens $\mathbf{V}_k = {\mathbf{v}_l}_{l=1}^L$. Denote K(t) as the image index before the t-th word 264 token. The pre-training objective of UrbanMLLM is to accurately predict the next word token with 265 preceding image and word tokens:

266 267

231

232

233 234 235

236

237 238 239

240 241

$$\mathcal{L}(\Theta = \{\theta, \zeta\}, \mathbf{w}, \mathbf{I}) = -\mathbb{E}_t[\log p_\Theta(\mathbf{w}_t | \mathbf{w}_{< t}), \mathbf{V}_{< K(t)})], \quad \mathbf{V}_{K(t)} = g_\zeta \circ f_\phi(I_{K(t)}).$$
(5)



Figure 3: Examples of satellite, street view, and cross-view tasks in instruct tuning dataset. Diverse task categories include Scene Classification (SC), Object Reasoning (OR), Landmark Recognition (LR), Spatial Relationship Reasoning (SRR), Geo-Localization (GL) and Indicator Prediction (IP).



Figure 4: UrbanMLLM consistently improves the downstream task accuracy compared with both open-sourced and closedsourced MLLMs. Abbreviations SI, SVI and CV stand for satellite imagery task, street view imagery task, cross-view task.

³¹⁵ 4 EXPERIMENTS



Figure 5: A performance comparison of UrbanMLLM's 3B, 8B, and 13B models across 13 urban understanding tasks.

In this section, we evaluate UrbanMLLM on three types of task: satellite view domain, street view
 domain and cross-view domain and then present the impact of various design choices on model
 performance.

4.1 EXPERIMENTAL SETUP

Dataset We use over 2 million satellite and street view images to build a large-scale cross-view interleaved pretraining dataset. Street view images offer ground-level details and appearance, while

325	Table 2: Satellite in	nagery-b	ased ur	ban und	erstand	ing resu	ilts on fi	ve tasks.
326	Satellite Imagery Task	S	С	OR	SRR	GL		IP
327	Sub-task	Single	Multi				Pop	Nightlight
328			0.000			0.000		
329	LLavA-N-8B	0.622	0.292	0.616	0.402	0.608	0.597	Failed
330	LLaVA-OV-7B	0.588	0.316	0.602	0.594	0.714	0.572	Failed
331	CogVLM2-19B	0.678	0.122	0.595	0.458	0.455	0.750	Failed
220	LLaVA-N-34B	0.574	0.220	0.629	0.588	0.608	0.597	Failed
332	VILA1.5-40B	0.650	0.152	0.645	0.583	0.475	0.599	Failed
333	InternVL-2-40B	0.664	0.479	0.672	0.593	0.756	0.632	Failed
334			0.101			0.010		
335	Qwen-VL-Plus	0.589	0.191	0.611	0.533	0.810	0.647	Failed
336	GPT-40	0.680	0.513	0.691	0.552	0.745	0.484	Failed
337	GeoChat	0.435	0.214	0.528	0.404	0.591	0.641	Failed
338	LHRS-Bot	0.439	0.128	0.568	0.386	0.243	0.533	0.449
339	UrbanMLLM-3B	0.901	0.816	0.815	0.590	0.909	0.923	0.735
340	UrbanMLLM-8B	0.910	0.825	0.821	0.577	0.924	0.898	0.789
341	UrbanMLLM-13B	0.898	0.810	<u>0.816</u>	0.626	0.906	0.871	0.728
342	Improv	33.8%	60.8%	18.8%	54%	14 1%	23.1%	
343	improvi	100.070	00.070	10.070	10.170	1.170		

Table 2: Satellite imagery-based urban understanding results on five tasks

344 satellite imagery provide top-down views, capturing urban structures for comprehensive understand-345 ing of the entire landscape. Satellite and street view imagery in the same census tract are batched 346 together with descriptive captions generated by MLLM. The corresponding county name and coor-347 dinates of the satellite image are also integrated into the batch. 348

We also construct an instruction tuning dataset for a variety of urban tasks, ranging from perception, 349 reasoning to numerical prediction, as detailed below: Satellite Imagery Tasks (SI): Scene Clas-350 sification (SC), Object Reasoning (OR), Spatial Relationship Reasoning (SRR), Geo-Localization 351 (GL),Indicator Prediction (IP), population density prediction (Pop) and nightlight intensity predic-352 tion (Nightlight) are the sub-task of Indicator Prediction. Single Scene Classification (Single) and 353 Multi-Scene Classification (Multi) are the sub-task of Scene Classification. Street View Imagery 354 Tasks (SVI): Scene Classification (SC), Object Reasoning (OR), Landmark Recognition (LR), Spa-355 tial Relationship Reasoning (SRR), Geo-Localization (GL), Indicator Prediction (IP), predicting the 356 beautiful (BF), wealthy (WE) and depressing (DP) level are the sub-tasks of Indicator Prediction. Cross-View Tasks (CV): Spatial Relationship Reasoning (SRR), Indicator Prediction (IP), predict-357 ing the median income (Med. income), poverty ratio (Pov. ratio), total population (Population) and 358 depression rate (Depr. rate) level are the sub-tasks of Indicator Prediction. More details on the task 359 settings and evaluation can be seen in A.6. 360

361 **Implementation** We initialize our model's weights using the pretrained VILA-1.5 model, and adapt the AdamW optimizer with a cosine learning rate scheduler during training. The training process 362 consists of two stages: in the first stage, we train on the entire interleaved pretraining dataset with a 363 batch size of 8 for one epoch, corresponding to 7200 steps with 8 hours. For the second stage, we 364 fine-tune the model on the instruct tuning dataset at a batch size of 16 for one epoch with 8 hours. 365 More information about baselines can be seen in A.5. 366

4.2 RESULTS

369 We compare the performance of our proposed UrbanMLLM with baselines on three tasks: satellite 370 imagery task, street view imagery task and cross-view task on Table 2, Table 3 and Table 4. Based 371 on these results, we have these noteworthy observations:

372 373 374

375

367

368

324

• UrbanMLLM achieves the best performance across both satellite view and street view tasks. The results showcase that UrbanMLLM achieves state-of-the-art performance, which successfully demonstrates the effectiveness of the proposed model for urban understanding tasks.

We observe that our model achieves over 85% accuracy on simple perception tasks, such as scene 376 classification and geo-localization, significantly outperforming general MLLMs. For more fine-377 grained tasks and object-level reasoning, our model outperformed the optimal baseline by 18.8%

410 411

412

413

414

415

416

417

418

419

379	Street View Task	SC	OR		SRR	GL		IP	
380	Sub-task						BF	WE	DP
381	LL aVA N 8B	0.513	0 102	0.643	0 705	0.575	0.600	0 503	0.413
382	LLaVA-IN-0D	0.515	0.472 0.572	0.045	0.703	0.782	0.000	0.575	0.413
383	CogVLM2-19B	0.004	0.372	0.580	0.742	0.702	0.482	$\frac{0.770}{0.435}$	0.027
384	LLaVA-N-34B	0.870	0.548	0.691	0.775	0.637	0.757	0.727	0.283
385	VILA1.5-40B	0.672	0.473	0.698	0.657	0.670	0.509	0.717	0.216
386	InternVL-2-40B	0.715	0.423	0.734	0.651	0.828	0.662	0.747	0.629
387	Qwen-VL-Plus	0.536	0.434	0.759	0.720	0.914	0.635	0.724	0.762
388	GPT-40	0.662	0.590	0.756	0.709	0.840	0.824	0.723	0.673
389	GeoChat	0.316	0.378	0.282	0.279	0.306	0.577	0.605	0.652
390	LHRS-Bot	0.532	0.221	0.295	0.316	0.242	0.189	0.325	0.255
391	Lishan MLLM 2D	0.020		0.025			0.026	0.775	0.705
392	UrbanMLLM-3B	0.829	0.703	0.814	0.971	0.899	0.830	0.779	0.745
393	UlualiiviLLivi-ob	0.042	0.703	0.014	$\frac{0.974}{0.976}$	$\frac{0.902}{0.002}$	0.041	0.770	0.740
394		0.044	0.702	0.029	0.970	0.902	0.004	0.790	0.795
395	Improv.	-3.0%	19.0%	10.0%	25.9%	-1.3%	4.9%	1.5%	4.3%
396									
397	Table 4: Cross view	v image	ery-base	d urban	underst	anding	results	on two	o tasks.
398	Cross-View Task				IP				SPP
399	Sub-task	Depr	rate M	ed inco	me Po	v ratio	Popula	ation	SKK
400			1410 111	0 5 5 5			1 opun		
401	LLaVA-OV-7B	0.43	87	0.557	(0.521	0.40	52	0.235
402	VILAI.5-40B	0.4	30 20	0.672		0.540 0.570	0.4	(4)	0.304
403	Intern v L-2-40B	0.5	38	0.397	t	0.572	0.40	52	0.280
404	Qwen-VL-Plus	0.5	12	0.648	C).618	0.48	39	0.299
405	GPT-40	Fail	ed	0.684	0).848	0.49	99	0.322
406	UrbanMLLM-3B	0.7	59	0.790	C	0.804	0.58	38	0.389
407	UrbanMLLM-8B	0.6	53	0.773	Ō	0.760	0.59	96	0.421
408	UrbanMLLM-13B	0.7	14	0.798	C	0.762	0.57	71	0.429
409	Improv.	41 1	%	16.7%		5.2%	194	%	33.2%
44.0	P- 0			10., 10		/ -	* / • •		

Table 3: Street view imagery-based urban understanding results on six tasks.

and 19.2%, respectively. This is because our model is pretrained on a large dataset of street-view and satellite images, allowing it to retain highly effective foundational image perception abilities. For more challenging reasoning tasks, such as predicting population density, our model outperforms the best general models by 23.1%, and surpassed the leading specialized models by 44.0% . It is important to note that in the SRR task, due to the limited availability of spatial relationship reasoning datasets in remote sensing, we used a task setup based on SkySenseGPT. As a result, our model has not previously encountered this specific task in the context of satellite imagery. Despite this, our model achieves performance comparable to general models, demonstrating its ability to acquire spatial relationship reasoning skills alongside its target inference capabilities.

420 On street-view tasks, our model achievs a 25.7% improvement in spatial relationship reasoning, 421 a 6.5% increase in average prediction accuracy, and a 10.0% advantage in landmark recognition. 422 Although its performance in geographic location prediction is slightly lower, trailing the closedsource model by 0.012, the model's consistency and strong results across other tasks demonstrate 423 its overall robustness and effectiveness in street-view tasks. Through pre-training, the model de-424 velops a nuanced understanding of spatial relationships and world knowledge, as well as the ability 425 to interpret abstract concepts such as beauty, wealth, or feelings of depression. Our model's per-426 formance on average metrics is comparable to that of the best open-source MLLMs, demonstrating 427 its strong generalization ability in handling complex urban understanding tasks. 428

UrbanMLLM demonstrates greater consistency in cross-view tasks. Table 4 showcases the performance of various models on cross-view tasks. It is evident that UrbanMLLM outperforms the baseline models in most tasks, especially in key areas such as depression rate, poverty ratio, and SRR, demonstrating their effectiveness in handling complex cross-view tasks. For exam-

432 ple, UrbanMLLM-3B achieves the best performance in the depression rate task, outperforming 433 InternVL-2-40B by 41%, which is a significant improvement over other models. UrbanMLLM-434 8B excels in both median income and SRR, with the former showing a 15% improvement over 435 the second-best model, VILA-1.5-40B, and the latter surpassing GPT-40 by 33.2%, highlighting 436 its strong spatial reasoning capability. This indicates that larger models, such as UrbanMLLM-8B, are better suited for tasks that require complex spatial and economic reasoning. In contrast, 437 other MLLMs like VILA-1.5-40B and Qwen-VL-Plus show mixed performance. While VILA 438 performs relatively well on the median income task, it falls behind in other tasks. GPT-40, despite 439 excelling in the poverty ratio task, fails to complete the depression rate task, revealing a lack of 440 consistency. In summary, UrbanMLLM provides more balanced and superior performance across 441 multiple tasks, significantly outperforming the baseline models, which often exhibit strengths in 442 specific areas but lack overall consistency. 443

 Model size enhances performance but complexity of urban understanding task determines 444 optimal gains. Firstly, there is a clear trend that larger MLLMs, such as UrbanMLLM-13B, gener-445 ally outperform smaller models like UrbanMLLM-3B across various tasks. This is demonstrated 446 in Figure 5. For example, in the CV-SRR task, the 13B model achieves a score of 0.429, compared 447 to 0.389 for the 3B model, indicating that increased model size often leads to better performance. 448 Similar patterns are seen in tasks like object reasoning and spatial relationship reasoning, suggest-449 ing that larger MLLMs capture the complexities of image-based tasks more effectively, whether in 450 single- or multi-task settings. However, in cross-view tasks (Table 4), the performance of the 3B 451 and 8B models is nearly identical, with the 3B or 8B model even slightly outperforming the 13B 452 model in the depression rate task. This indicates that MLLMs do not always guarantee superior performance, and that task complexity and data characteristics also significantly influence results. 453

455 4.3 EVALUATION ON URBANVIEW BENCHMARK

We evaluate our model with different open-source and closed-source MLLMs on our benchmark.
We tested various models with the same set of questions on the same dataset. Due to differences in the models' ability to follow instructions, many do not provide answers exactly matching the ground truth but instead include additional explanatory text. Therefore, we consider a response correct as long as it contains the correct answer.

462 As shown in Table 2, 3 4 and Figure 4, our benchmark reveals the key challenges and limitations of current MLLMs in real-world urban environments. The results show that most advanced MLLMs 463 do not perform well in satellite and street view tasks. For satellite view tasks, the top-performing 464 closed-source models, such as GPT-40 and another leading model, InternVL-2-40B, achieved only 465 52.4% and 54.2% on average across various metrics. On street view imagery, their performance is 466 similarly limited, with average scores of 72.2% and 67.4%, respectively. This discrepancy is because 467 most of our images is collected recently, while the training data for these MLLMs generally lacks 468 similar real-world data (Wang et al., 2024). Furthermore, many current MLLMs do not yet support 469 multi-image inputs, and those do rarely handle tasks involving joint cross-view predictions for urban 470 understanding. Consequently, this benchmark clearly highlights the limitations of advanced models, 471 showing their challenges in performing well on urban understanding tasks, especially with street 472 view and remote sensing images, and in joint cross-view prediction tasks.

473 474 4

475

454

456

4.4 ABLATION ANALYSIS

To evaluate the effectiveness of each module in UrbanMLLM, we evaluate the performance of various task of different model variants in Table 5, Table 6, and Table 7. Specifically, we evaluate the UrbanMLLM without cross-view perceiver (w/o Perceiver), satellite imagery in the pretraining stage and cross-view perceiver (w/o SI+Perceiver), street view imagery in the pretraining stage and cross-view perceiver (w/o SI+Perceiver), satellite imagery and street view imagery in the pretraining stage and cross-view perceive (w/o SI+SVI+Perceiver). Note that in the variant without the cross-view Perceiver (w/o Perceiver), a two-layer MLP is implemented as a replacement.

According to the results, cross-view perceiver is the most essential module for explicitly facilitate
 mutual learning of cross-view urban imagery. It brings 2%-81% gains for all tasks, because satel lite and street view images represent two completely different modalities of information, making it
 difficult to directly integrate and interact within LLMs during the pretraining stage to learn cross-

487	Table 5: Ablation stud	ly of Ur	banMl	LLM va	ariants o	on sate	llite im	agery	tasks.
488	Variants	S	SC		SDD	CI		IP	
489	Variants	Single	Mult	i l			Pop	Nigh	tlight
490		10111810		-	<u> </u>	1	- °P		
491	UrbanMLLM-8B	0.910	0.825	5 0.821	1 0.577	0.924	0.898	0.'	789
492	w/o Perceiver	0.749	0.596	5 0.732	2 0.106	0.427	0.869	<u>0.'</u>	747
493	w/o SI+Perciver	0.897	0.819	$\theta 0.81^{2}$	1 0.549	0.921	0.913	0.'	737
494	w/o SVI+Perceiver	0.907	0.822	$\frac{2}{0.818}$	<u>8</u> 0.615	0.919	0.880	0.'	707
495	w/o SI+SVI+Perceiver	0.888	0.806	5 <u>0.818</u>	<u>8 0.604</u>	0.903	6 0.869	0.'	713
496									
497	Table 6: Ablation study	of Urb	anMLl	LM var	iants on	street	view in	nager	y tasks.
498	Variants	SC	OR	тм	SRR	GI		IP	
499	Variants						BF	WE	DP
500	UrbanMI I M 8D	0 842	0 702	0.01/	0.074	0.002	0.941	0 770	0.746
501	UlballiviLLivi-ob	$\frac{0.042}{0.666}$	0.703	0.014	0.974	0.902	0.041	$\frac{0.770}{0.760}$	0.740
502	w/o SL Derectiver	0.000	0.040	0.455	0.800	0.042	0.771	0.700	0.727
503	w/o SVI + Perceiver	0.829	0.099	0.803	0.974	0.891	0.097	0.776	0.762
504	w/o SI+SVI+Perceiver	0.772	0.701 0.696	0.814	$\frac{0.973}{0.964}$	0.888	$\frac{0.880}{0.878}$	0.774	0.734 0.737
505		0.772	0.070	0.012	0.201	0.000	0.070	0.771	0.757
506									
507	Table 7: Ablation s	tudy of	Urban	MLLN	l varian	ts on c	ross-vie	ew tas	ks.
508	Variants				IP				SPP
509	variants	Depr. r	ate M	ed. inc	ome Po	ov. ratio	o Popu	lation	
510		0.650	<u>, </u>	0 770		0 = < 0		-0.6	
511	UrbanMLLM-8B	0.653	5	0.773		0.760	0.3	5 70	0.421
512	w/o Perceiver	0.493) •	0.462		0.520	0.4	+/ð	0.247
513	W/O SI+Perceiver	0.752	4	$\frac{0.792}{0.792}$		0.759	0.3	520	0.372
	w/o SVI+Perceiver	0.714	ł	0.782		0.739	0.3	531	10.419

0.674

514 515 516

486

517 view semantic information. Therefore, a specialized mechanism is required to fuse these modalities in advance. The use of satellite image data during pretraining has a significant impact on satellite 518 image-related tasks, contributing performance gains ranging from 0.3% to 7.1%. However, its ef-519 fect on various economic indicators in cross-view tasks differs. For example, it resulted in a 15.2% 520 improvement in the depression rate task but caused a 12.8% decrease in accuracy for total popu-521 lation estimation of one region. This difference is due to the depression rate being more closely 522 related to visible green space in satellite images, while population estimation requires a more nu-523 anced understanding of urban environmental factors. Similarly, using street view image data during 524 pretraining has a greater impact on street view-related tasks compared to satellite data, contributing 525 performance gains of 0.3% to 1.7%. This demonstrates that pretraining with data closely aligned 526 to downstream tasks can significantly enhance model performance. Additionally, the interleaved 527 image-text pretraining on satellite and street view images provides a task-agnostic yet semantically rich initialization, contributing a performance gain of 0.3% to 9% in tasks such as scene classifi-528 cation on street view images. Therefore, the interleaved image-text pretraining of the two types of 529 images, along with the cross-view perceiver, are essential components of our approach. 530

0.793

0.697

0.557

0.348

531 532

5 CONCLUSION

w/o SI+SVI+Perceiver

533 534

In this paper, we propose UrbanMLLM, a novel multimodal large language model designed to jointly
 learn from remote sensing and street-view imagery for comprehensive urban understanding. By
 leveraging a large-scale cross-view dataset and a cross-view perceiver architecture, UrbanMLLM
 effectively captures complementary information from satellite-view and street-view. Our model out performs existing MLLMs, achieving significant improvements across various urban understanding tasks.

540 REFERENCES 541

569

570

571

572

573

577

578

579

580

584

585

542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical 543 report. arXiv preprint arXiv:2303.08774, 2023. 544

- Guillaume Astruc, Nicolas Dufour, Ioannis Siglidis, Constantin Aronssohn, Nacim Bouia, Stephanie 546 Fu, Romain Loiseau, Van Nguyen Nguyen, Charles Raude, Elliot Vincent, et al. Openstreetview-547 5m: The many roads to global visual geolocation. In Proceedings of the IEEE/CVF Conference 548 on Computer Vision and Pattern Recognition, pp. 21967-21977, 2024. 549
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang 550 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. 551 arXiv preprint arXiv:2308.12966, 2023. 552
- 553 Yakoub Bazi, Laila Bashmal, Mohamad Mahmoud Al Rahhal, Riccardo Ricci, and Farid Melgani. 554 Rs-llava: A large vision-language model for joint captioning and question answering in remote 555 sensing imagery. Remote Sensing, 16(9):1477, 2024. 556
- U.S. Census Bureau. 2019 cartographic boundary shapefile: Current census tract for 558 united states (1:500,000), 2019. URL https://catalog.data.gov/dataset/ 2019-cartographic-boundary-shapefile-current-census-tract-for-united-states-1-5000 559 Accessed: November 21, 2024. 560
- 561 Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. Honeybee: Locality-enhanced 562 projector for multimodal llm. In Proceedings of the IEEE/CVF Conference on Computer Vision 563 and Pattern Recognition, pp. 13817-13827, 2024. 564
- 565 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qing-566 long Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. arXiv 567 preprint arXiv:2312.14238, 2023. 568
 - Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 24185–24198, 2024.
- 574 Abhimanyu Dubey, Nikhil Naik, Devi Parikh, Ramesh Raskar, and César A. Hidalgo. Deep learning the city : Quantifying urban perception at a global scale, 2016. URL https://arxiv.org/ 575 abs/1608.01769. 576
 - Zhuangyuan Fan, Fan Zhang, Becky PY Loo, and Carlo Ratti. Urban visual intelligence: Uncovering hidden city profiles with street view images. Proceedings of the National Academy of Sciences, 120(27):e2220417120, 2023.
- 581 Jie Feng, Jun Zhang, Junbo Yan, Xin Zhang, Tianjian Ouyang, Tianhui Liu, Yuwei Du, Siqi Guo, 582 and Yong Li. Citybench: Evaluating the capabilities of large language model as world model. 583 arXiv preprint arXiv:2406.13945, 2024.
- Google. Google maps api, 2024. URL https://developers.google.com/maps/ documentation. Accessed: 2024-11-21. 586
- Xixuan Hao, Wei Chen, Yibo Yan, Siru Zhong, Kun Wang, Qingsong Wen, and Yuxuan Liang. 588 Urbanvlp: A multi-granularity vision-language pre-trained foundation model for urban indicator 589 prediction. arXiv preprint arXiv:2403.16831, 2024. 590
- 591 Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. 592 Visual genome: Connecting language and vision using crowdsourced dense image annotations, 2016. URL https://arxiv.org/abs/1602.07332.

Kartik Kuckreja, Muhammad Sohail Danish, Muzammal Naseer, Abhijit Das, Salman Khan, and 595 Fahad Shahbaz Khan. Geochat: Grounded large vision-language model for remote sensing. 596 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 597 27831-27840, 2024. 598 Benjamin Lee. National, state-level, and county-level prevalence estimates of adults aged ≥ 18 years self-reporting a lifetime diagnosis of depression—united states, 2020. MMWR. Morbidity 600 and Mortality Weekly Report, 72, 2023. 601 602 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, 603 Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024a. URL https: 604 //arxiv.org/abs/2408.03326. 605 Xuecao Li, Yuyu Zhou, Min Zhao, and Xia Zhao. A harmonized global nighttime light dataset 606 1992-2018. Scientific data, 7(1):168, 2020. 607 608 Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng 609 Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. 610 arXiv preprint arXiv:2403.18814, 2024b. 611 Yunxin Li, Shenyuan Jiang, Baotian Hu, Longyue Wang, Wanqi Zhong, Wenhan Luo, Lin Ma, and 612 Min Zhang. Uni-moe: Scaling unified multimodal llms with mixture of experts. arXiv preprint 613 arXiv:2405.11273, 2024c. 614 615 Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, 616 Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023. 617 Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On pre-618 training for visual language models. In Proceedings of the IEEE/CVF Conference on Computer 619 Vision and Pattern Recognition, pp. 26689–26699, 2024. 620 621 Fan Liu, Delong Chen, Zhangqingyun Guan, Xiaocong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, 622 and Jun Zhou. Remoteclip: A vision language foundation model for remote sensing. IEEE 623 Transactions on Geoscience and Remote Sensing, 2024a. 624 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 625 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL https:// 626 llava-vl.github.io/blog/2024-01-30-llava-next/. 627 628 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances 629 in neural information processing systems, 36, 2024c. 630 Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei 631 Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for 632 open-set object detection. arXiv preprint arXiv:2303.05499, 2023. 633 634 Junwei Luo, Zhen Pang, Yongjun Zhang, Tingzhu Wang, Linlin Wang, Bo Dang, Jiangwei Lao, Jian 635 Wang, Jingdong Chen, Yihua Tan, et al. Skysensegpt: A fine-grained instruction tuning dataset 636 and model for remote sensing vision-language understanding. arXiv preprint arXiv:2406.10100, 637 2024. 638 Utkarsh Mall, Cheng Perng Phoo, Meilin Kelsey Liu, Carl Vondrick, Bharath Hariharan, and Kavita 639 Bala. Remote sensing vision-language foundation models without annotations via ground remote 640 alignment. In The Twelfth International Conference on Learning Representations, 2024. 641 642 Rohin Manvi, Samar Khanna, Gengchen Mai, Marshall Burke, David Lobell, and Stefano Er-643 mon. Geollm: Extracting geospatial knowledge from large language models. arXiv preprint 644 arXiv:2310.06213, 2023. 645 Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, 646 Dhruti Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights 647 from multimodal llm pre-training. arXiv preprint arXiv:2403.09611, 2024.

648 Dilxat Muhtar, Zhenshi Li, Feng Gu, Xueliang Zhang, and Pengfeng Xiao. Lhrs-bot: Empow-649 ering remote sensing with vgi-enhanced large multimodal language model. arXiv preprint 650 arXiv:2402.02544, 2024. 651 652 Chao Pang, Jiang Wu, Jiayu Li, Yi Liu, Jiaxing Sun, Weijia Li, Xingxing Weng, Shuai Wang, Litong 653 Feng, Gui-Song Xia, and Conghui He. H2rsvlm: Towards helpful and honest remote sensing large vision language model, 2024. 654 655 656 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 657 models from natural language supervision. In International conference on machine learning, pp. 658 8748-8763. PMLR, 2021. 659 660 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-661 baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem-662 ini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint 663 arXiv:2403.05530, 2024. 664 665 SafeGraph. Open census data documentation, 2024. URL https://docs.safegraph.com/ 666 docs/open-census-data. Accessed: November 21, 2024. 667 668 Andrew J Tatem. Worldpop, open data for spatial demography. Scientific data, 4(1):1–4, 2017. 669 670 Vicente Vivanco Cepeda, Gaurav Kumar Nayak, and Mubarak Shah. Geoclip: Clip-inspired align-671 ment between locations and images for effective worldwide geo-localization. Advances in Neural Information Processing Systems, 36, 2024. 672 673 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, 674 Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng 675 Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's 676 perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024. 677 678 Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, 679 Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. 680 Cogvlm: Visual expert for pretrained language models, 2023. 681 682 Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2 – a large-683 scale benchmark for instance-level recognition and retrieval, 2020. URL https://arxiv. 684 org/abs/2004.01804. 685 686 Shixiong Xu, Chenghao Zhang, Lubin Fan, Gaofeng Meng, Shiming Xiang, and Jieping Ye. Ad-687 dressclip: Empowering vision-language models for city-wide image address localization. arXiv 688 preprint arXiv:2407.08156, 2024. 689 Yibo Yan, Haomin Wen, Siru Zhong, Wei Chen, Haodong Chen, Qingsong Wen, Roger Zimmer-690 mann, and Yuxuan Liang. Urbanclip: Learning text-enhanced urban region profiling with con-691 trastive language-image pretraining from the web. In Proceedings of the ACM on Web Conference 692 2024, pp. 4006-4017, 2024. 693 694 Wei Zhang, Miaoxin Cai, Tong Zhang, Yin Zhuang, and Xuerui Mao. Earthgpt: A universal multi-695 modal large language model for multi-sensor image comprehension in remote sensing domain. 696 IEEE Transactions on Geoscience and Remote Sensing, 2024. 697 Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Young-699 jae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal c4: An open, billion-700 scale corpus of images interleaved with text. Advances in Neural Information Processing Systems, 701 36, 2024.

702 A APPENDIX

704 A.1 ETHICS

Our model uses a large amount of satellite and street view images, which poses a potential risk to individual privacy. While the resolution of the satellite imagery is not high enough to identify individuals, it can still detect environmental changes resulting from human activity. The street view images were crawled and downloaded from the Google platform, with key information blurred, ensuring that no private information is compromised.

Our model is designed to better understand cities from cross views by incorporating various data types to enhance its understanding capabilities. To minimize the misuse of our model and data, we will release the dataset and trained model only to those who agree to adhere to ethical guidelines. By following these guidelines, users agree to comply with laws, regulations, and the ASPRS Code of Ethics. It is also important to note that the data used for this training is already freely available and public, so our model does not exacerbate privacy concerns.

717 718

719

704

A.2 LIMITATIONS AND FUTURE WORK

Although our model, UrbanMLLM, covers a wide range of urban understanding tasks and achieves
 state-of-the-art performance, there are still some limitations. Our dataset is limited to the United
 States, and the generalization of the model to other countries may require additional data collection
 and pre-training. Therefore, we plan to extend the dataset to cover more regions and improve the
 general applicability of the model, and explore the use of additional modalities for urban under standing tasks.

- 726
- 727 728

A.3 IMPLEMENTATION DETAILS FOR REPRODUCIBILITY

We perform experiments using Python 3.10 and Pytorch 2.3.0+cu121 with 8× NVIDIA A100 GPUs.
Here we provide detailed values of the hyper-parameters used in the experiments for reproducibility in Table 13 and Table 14 for the training and testing, respectively.

732

736

737

738

739

740 741 742

- 733 A.4 EXPERIMENTAL RESULTS 734
- 735 A.4.1 ABLATION STUDY

As our model supports multi-image as input, we conduct an ablation study on the number of images (N = 2, 4, 6) for indicator prediction tasks. For example, 6 images means one satellite image and five street-view images as input. The results are as follows (8). It can be seen that more images bring certain performance gain for the indicator prediction task.

Table 8: Image number ablation study of UrbanMLLM on cross-view imagery tasks.

Number of input images	IP							
g	Depr. rate	Med. income	Pov. ratio	Population				
2	0.724	0.766	0.696	0.596				
4	0.772	0.789	0.689	0.596				
6	0.764	0.809	0.720	0.614				

749 750

751 A.4.2 DATA SCALE STUDY 752

We also provide the results of UrbanMLLM trained on different dataset scales, including 0.35 million, 0.92 million, and 1.86 million images. The results (Table 9, Table 10, Table 11) show that the performance of UrbanMLLM improves a little with the increase in dataset scale. Because our pre-training data is much less than data size that scaling law requires.

Table 9: Satellite imagery-based urban understanding results with different data scale on five tasks.

			U				
Satellite Imagery Task	S	2	OR	SRR	GL		IP
Sub-task	Single	Multi				Pop	Nightlight
0.35M	0.899	0.816	0.822	0.644	0.921	0.904	0.741
0.92M	0.908	0.825	0.823	0.608	0.935	0.923	0.775
1.86M	0.910	0.825	0.821	0.577	0.924	0.898	0.789

Table 10: Street view imagery-based urban understanding results with different data scale on six tasks.

Street View Imagery Task	SC	OR			GL		IP	
Sub-task						BF	WE	DP
0.35M	0.840	0.706	0.847	0.982	0.902	0.834	0.777	0.746
0.92M	0.835	0.703	0.825	0.970	0.894	0.838	0.761	0.773
1.86M	0.842	0.703	0.814	0.974	0.902	0.841	0.778	0.746

Table 11: Cross view imagery-based urban understanding results with different data scale on two tasks.

Cross-View Imagery Task		IP			
Sub-task	Depr. rate	Med. income	Pov. ratio	Population	
0.35M	0.754	0.781	0.744	0.486	0.415
0.92M	0.701	0.805	0.750	0.555	0.418
1.86M	0.653	0.773	0.760	0.596	0.421

A.4.3 EVALUATION ON CITYBENCH

A.4.4 CASE STUDY

We also provide the results of the proposed dataset on other benchmarks, including Citybench (Feng et al., 2024). We select CityInfer, LocInfer, and Population as tasks to evaluate the performance. We use Accuracy, Accuracy@25km, and R^2 as evaluation metrics. More details can be found in the Citybench. The results are shown in Table 12. The proposed dataset outperforms the state-of-theart models on these benchmarks, demonstrating the effectiveness of the proposed dataset for urban understanding tasks.

Table 12: Best Performance on Close-Source Model, Open-Source Model, and UrbanMLLM on Citybench.

Model	CityInfer	LocInfer	Population
SOTA closed-source model	0.862	0.797	0.122
SOTA open-source model	0.574	0.555	-0.113
UrbanMLLM-8B	0.904	0.840	0.324

We have added a bad case analysis in the revised paper. We show some examples of bad cases in scene classification and indicator prediction tasks. The results are shown in Figure 6, 7, 8. Firstly, in the scene classification task, the model misclassifies the image with a truck parking as a car parking. Although there are a few differences between the two classes, the more granular understanding of the urban environment is required to distinguish them. Secondly, in the indicator prediction task, as shown in Figure 7, 8, the model predicts the population density of an urban area as 6.8 and the actual value is 9.9 using a satellite image. The model fails to capture detailed information with a single-view image, which makes it challenging for the model to learn from the limited dataset. For poverty rate prediction, the model gets a high score of 5.4, but the poverty rate is 2.6. It's may be



	Stage1	Stage2
Optimizer	AdamW	AdamW
Learning Rate	5e-5	1e-4
Batch Size	8	16
Accumulation Step(s)	1	2
Weight Decay	0.	0
Epoch(s)/Step(s)	1 Epoch	1 Epoch
Save Steps	1200	750
Scheduler	Cos	sine
Warmup Ratio	500	100
Model Max Length	20	48



A.5 BASELINES

We evaluated several advanced MLLMs on the UrbanView Benchmark. However, some of the 903 MLLMs are not pretrained on multi-image data or does not support multi-image inference, such 904 as LLaVA-Next (Liu et al., 2024b) and CogVLM2 (Wang et al., 2023), so only single-image 905 tasks are evaluated. For VILA-1.5 (Lin et al., 2023), InternVL2 (Chen et al., 2023), and LLaVA-906 OneVision (Li et al., 2024a), the whole benchmark evaluation is done. In addition to the open-source 907 models, state-of-the-art closed-source models Qwen-VL-Plus and GPT-40 are also fully evaluated 908 on the benchmark. Specifically, we assess a satellite domain-specific model, GeoChat (Kuckreja 909 et al., 2024), to further prove our capability. To ensure fairness, domain-specific models that are not 910 yet open-source are excluded from this evaluation.

911

899 900 901

902

912 A.6 URBANVIEW DATASET AND BENCHMARK 913

As the thriving development of satellite and street view data, a series of publicly available datasets
about urban imagery has been brought up. However, considering the complexity and diveristy of
urban environment, single-view data of satellite-view or street-view is not enough. Therefore, to
enhance MLLM's comprehensive and all-level understanding of cities, we propose UrbanView, a
dataset and benchmark that composes of massive amount of multi-source and cross-view urban



imagery, in combination with various collected labels across geo-locations, grounded objects, spatial relationships, income and health indicators.

950

A.6.1 DATA COLLECTION

951 The images are primarily collected from two sources: Google Maps API (Google, 2024) for street 952 view images, and ESRI for satellite imagery. For street view imagery collection, we randomly gen-953 erate 2,000 random points in each census tract polygon and use their coordinates to query Google 954 Maps API, returning street view image patches and real coordinates. We scrape over 2 million street 955 view images and all satellite imagery of zoom level 15 across the United States. We further gather 956 a variety of socioeconomic data of census tract and grid level from world pop, NIH and US government. We also collect a series of open datasets, such as Google Landmarks Recognition Weyand 957 et al. (2020), Visual Genome Krishna et al. (2016) and Place Pules Dubey et al. (2016) etc. By apply-958 ing some domain-specific adaptation to the original ones, we build a more well-rounded UrbanView 959 dataset. 960

961 The data used for dataset construction have a sparse yet overall coverage of the United States as 962 shown in Figure 9. We use census tract boundary data in 2019 (Bureau, 2019), and gather street 963 view and satellite images in 71,433 out of all 73,868 census tracts in the United States, which is about 96.7%. The Google street view images can be acquired using coordinate queries, however, 964 we don not know the exact coordinate of where the street view exists. We randomly generate 2,000 965 points in each census tract and use these points to query street view images. This is a random 966 process, thus we are not able to sample all the images in a census tract considering the time cost. In 967 fact, in some less populated areas, it is quite hard to get street view images because the randomly 968 generated query points in these areas are always off-road, which is also the main reason for the 969 missing 3.3% coverage. In the end, we randomly sample about 200 images in each census tract, 970 which have been proved to be effective in indicator prediction tasks.

971

The data size for each task of the UrbanView dataset and benchmark is listed below in Table 15:

989 990 991

992 993

Source	Task	Dataset Size	Benchmark Siz
	Scene Classification (SC)	30,000	1,000
	Object Reasoning (OR)	90,000	1,000
Charles Miner	Landmark Recognition (LR)	30,000	1,000
Street view	Spatial Relationship Reasoning (SRR)	30,000	1,000
	Geo-Localization (GL)	30,000	1,000
	Indicator Prediction (IP)	90,000	3,000
	Scene Classification (SC)	51,759	8,668
	Object Reasoning (OR)	115,115	5,556
Satellite	Spatial Relationship Reasoning (SRR)	90,000	8,250
	Geo-Localization (GL)	29,629	1,000
	Indicator Prediction (IP)	60,000	2,000
Cross View	Spatial Relationship Reasoning (SRR)	30,000	1,000
CIUSS-VIEW	Indicator Prediction (IP)	120,000	4,000

11.00 . . .

DATASET STRUCTURE AND CONSTRUCTION A.6.2

994 We first build a large-scale cross-view interleaved pretraining dataset. For each census tract, we 995 match the coordinate between satellite and street view imagery. Since we collect all satellite images in the United States, each street view image can find a match. However, in order to control the size 996 of inputs, at most 5 street view images are matched with single satellite image. For the next step, we 997 use a powerful open-source MLLM, InternVL2-40B, to generate detailed descriptive captions for 998 them. Using a similar data structure in MMC4 (Zhu et al., 2024), the county name and coordinates 999 of the satellite image are also embedded to the interleaved pre-training data, together with imagery 1000 and caption embeddings. 1001

We use a Human-AI mixture method for pre-training caption quality validation. We first use two 1002 powerful open-source MLLM, VILA-1.5-40B and LLaVA-Next-34B to judge if the caption matches 1003 with the image. If either of them thinks it is not a match, we will proceed to send this case to GPT-40, 1004 which has state-of-the-art comprehension ability, but not quite affordable for large-scale deployment. 1005 If GPT-40 also thinks there is a problem with the case, graduate-level human-being will manually check this case to give the final judgement. In order to quantify the caption quality improvement, 1007 we further use GPT-40 to regenerate captions for excluded images with human assistance to test the 1008 quality improvement. We use CLIP-Score as the evaluation metric and calculate our original caption 1009 score and cleaned caption score of 10,000 samples, resulting 29.99535 and 29.99571 respectively. 1010 As a matter of fact, the original caption quality is good enough and only about 1.3 out of a thousand images is marked as unmatched by two-stage MLLM verification, and the regeneration process 1011 enhances the caption quality only by 0.0012%. 1012

1013 Then we construct the instruct tuning dataset, which is categorized into 3 major types: satellite-1014 only, street view-only, and cross-view data. Street view images offer ground-level data of various 1015 environments, including urban and rural areas. These images provide detailed features of the envi-1016 ronment. In contrast, satellite imagery complements this by providing top-down views, capturing a 1017 overall perception of the entire landscape. In UrbanView dataset, not only satellite and street view images linked respectively with diverse tasks such as scene classification, object reasoning, spatial 1018 relationship reasoning, cross-view combinational tasks of socioeconomic prediction and image re-1019 trieval are delicately designed to further enhance MLLM's comprehensive understanding of urban 1020 environment. 1021

In the construction of UrbanView dataset, a lot of street view tasks are in lack of groundtruth labels. 1023 Therefore, we use light-weight object detection specialized model to generate groundtruth bounding boxes for object reasoning tasks, and use powerful open-source MLLMs to identify the scene class 1024 and spatial relationship of the street view image. For geo-localization tasks, we simply use latitude 1025 and longitude of the images to match the boundaries of census tracts in the United States.

We construct an instruction tuning dataset for a variety of urban tasks, ranging from perception, reasoning to numerical prediction, as detailed below:

• Satellite Tasks:

The satellite Scene Classification (SC), Object Reasoning (OR), and Spatial Relationship Reasoning (SRR) dataset are the same as FIT-RS dataset in SkySenseGPT, since they have built a high-quality dataset and proved to be effective on satellite tasks. Scene Classification (SC): Select which scene or scenes does this image conform to.

- ¹⁰³³ **Object Reasoning** (**OR**): Respond the location, presence or count of a specific object.
- Spatial Relationship Reasoning (SRR): Select the correct object relationship displayed in the image.

Geo-Localization (GL): Select which county this image belongs to. We use the latitude and longitude coordinate of each image to find the corresponding county it belongs to. Only the most populated 100 counties in the United States are taken account of. We also use multiple choices format for this task type, and the distraction choices are randomly chose from the 100 counties.

- Indicator Prediction (IP): Predict population density or nightlight intensity from 0.0 to 9.9. We
 follow the normalization method used in GeoLLM (Manvi et al., 2023), scaling down the population and nightlight density to the range of 0.0 to 9.9, and ask the MLLMs to give a direct estimation in this range. The population density data are sourced from WorldPop (Tatem, 2017) and the nightlight data are sourced from VIIRS (Li et al., 2020).
- Street View Tasks:
- **Scene Classification (SC)**: Select which scene does this image conform to. The ground-truth is obtained using LLaVA-Next-34B, which have been verified to generate a pretty reasonable result on scene classification task.
- 1048 object Reasoning (OR): Respond the location, presence or count of a specific object. There are three kinds of sub-tasks in object reasoning, all ground-truth annotations are generated by Ground-ing DINO (Liu et al., 2023), which has shown state-of-the-art ability on open vocabulary object detection. We further process the bounding box and object name results given by Grounding DINO to build object presence and counting dataset.
- Landmark Recognition (LR): Select the correct landmark name shown in the image. We use images from google landmarks dataset v2, and select the street view images in the dataset via LLaVA-Next-34B. Multiple choice questions are made based on the correct landmark name and three distraction landmark names.
- Spatial Relationship Reasoning (SRR): Select the correct object relationship displayed in the image. This is a multiple choice question with four choices. The correct choice is the ground-truth object relationship in Visual Genome dataset, and we format the question by using the object and subject name, such as "What is the relationship between girl and computer?". The three distraction choices are generated by InternVL2-40B based on the image provided with factually incorrect relationships. We also attempt to use other MLLM for this distractor generation task, including LLaVA-Next-34B and Vila-1.5, but InternVL2-40B is the one that generates the most reasonable distractors.
- **Geo-Localization (GL)**: Select which county this image belongs to. We use the latitude and longitude coordinate of each image to find the corresponding county it belongs to. Only the most populated 100 counties in the United States are taken account of. We also use multiple choices format for this task type, and the distraction choices are randomly chose from the 100 counties.
- Indicator Prediction (IP): Predict the beautiful, wealthy and depressing level of the image from a level of 0.0 to 9.9. We use Place Pulse 2.0 dataset, which let human to make comparison between two images in multiple dimensions. Then a ranking algorithm is used to assign ground-truth labels for the images, and we ask the MLLMs to give a direct estimation in this range.
- 1072 Cross-View Tasks:

Indicator Prediction (IP): Predict the median income, poverty ratio, total population (SafeGraph, 2024) and depression rate level (Lee, 2023) of a set of images in the same census tract or 1 kilometer map grid from a level of 0.0 to 9.9. We follow the normalization method used in GeoLLM (Manvi et al., 2023), scaling down the indicators to the range of 0.0 to 9.9, and ask the MLLMs to give a direct estimation in this range.

- **Spatial Relationship Reasoning (SRR)**: Figure out which part of the satellite image does the street
- view image under the same area belong to. The answer should be selected from 'top-left', 'top-right', 'bottom-left', and 'bottom-right'.

1080 A.6.3 URBANVIEW BENCHMARK AND EVALUATION

1082	We propose UrbanView Benchmark and construct corresponding evaluation methods. One thousand
1083	data points are sampled from our dataset for each street view, cross-view, and satellite indicator
1084	prediction task for the benchmark evaluation, while the original data size are kept for all the satellite
1085	tasks based on FII-RS. For the Benchmark for evaluation, all metrics are in format of accuracy,
1086	except for the satellite multi scene classification task, which uses F1-score as a evaluation metric.
1087	The ground-truth labels in benchmark also remain identical with our dataset.
1088	
1089	
1090	
1091	
1092	
1093	
1094	
1095	
1096	
1097	
1098	
1099	
1100	
1101	
1102	
1103	
1104	
1105	
1106	
1107	
1108	
1109	
1110	
1111	
1112	
1113	
1114	
1115	
1116	
1117	
1118	
1119	
1120	
1121	
1122	
1123	
1124	
1125	
1126	
1127	
1128	
1129	
1130	
1131	
1132	
1133	

5	prompt ='''	
6	You are a powerful street-view image captioner.	
7	image.	
8	The caption annotation procedure follows the principles of:	Captioner: InternVL2-40B
9	(1): Describing object attributes, including object quantity, color, material, shape, size, and spatial	
0	position (including absolute position in the image and	
н	relative position between objects);	
2	specific object;	
3	(3): Instead of describing the imaginary content, only	ST. TT.II
4	describing the content one can determine confidently from the image.	
5	Do not describe the contents by itemizing them in list form.	
6	Minimize aesthetic descriptions as much as possible; (4): Please output less 35 words	
7		
8	Answer - 111	
9	A three-story residential building with light green and	And the second se
50	beige exterior, white trim, and multiple windows. A red SUV	
51	and a DIACK car parked on the street. A small white garage and a tree in the background.	
52		
3		
54	Figure 11: Satellite image captioning for pr	etraining dataset.
5		
56		
57	prompt -111	
58	You are a powerful remote sensing and aerial image	
59	captioner.	
50	Please create SHURI captions describing the contents of the given image.	
61	The caption annotation procedure follows the principles of:	Captioner: InternVL2-40B
51 52	 (1): Describing object attributes, including object (a): concerning object attributes, and spatial 	Captioner: InternVL2-40B
51 52 53	 (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and 	Captioner: InternVL2-40B
51 52 53 54	The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the	Captioner: InternVL2-40B
51 52 53 54 55	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object;</pre>	Captioner: InternVL2-40B
51 52 53 54 55 56	 The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only 	Captioner: InternVL2-40B
51 52 53 54 55 56 66 57	 The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. 	Captioner: InternVL2-40B
51 52 53 54 55 56 66 57 58	 The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. 	Captioner: InternVL2-40B
51 52 53 54 55 55 56 57 58 59	 The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words. 	Captioner: InternVL2-40B
51 52 53 54 55 56 57 58 59 59 70	The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.	Captioner: InternVL2-40B
51 52 53 55 55 56 66 57 58 59 70 70 71	The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.	Captioner: InternVL2-40B
61 62 63 64 65 66 67 68 69 70 70 71 72	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words. '''' Aerial view of a suburban area with a mix of commercial</pre>	Captioner: InternVL2-40B
51 52 53 55 56 57 58 59 70 71 72 73	<pre>he caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words. Answer =''' Aerial view of a suburban area with a mix of commercial buildings, parking lots, and green spaces. A major road commercial the content with a large of the suburban area with a large of the scene with a large of the suburban area with a large of the scene with a large of the suburban area with a large of the scene with a large of the suburban area with a large of the scene with a large of the scene of the s</pre>	Captioner: InternVL2-40B
51 52 53 55 55 56 57 58 59 70 71 72 73 74	<pre>he caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words. Answer =''' Aerial view of a suburban area with a mix of commercial buildings, parking lots, and green spaces. A major road curves through the scene, with a large white building near the center.</pre>	Captioner: InternVL2-40B
51 52 53 55 56 57 58 59 70 71 72 73 74 75	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B
51 52 53 54 55 56 57 58 59 70 71 72 73 74 75 76	<pre>he caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words. ''' Answer =''' Aerial view of a suburban area with a mix of commercial buildings, parking lots, and green spaces. A major road curves through the scene, with a large white building near the center. '''</pre>	Captioner: InternVL2-40B
61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77	The caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words. ''' Answer =''' Aerial view of a suburban area with a mix of commercial buildings, parking lots, and green spaces. A major road curves through the scene, with a large white building near the center. '''	Captioner: InternVL2-40B
51 52 53 54 55 56 57 58 59 70 71 72 73 74 75 76 77 78	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B Output: Description: InternVL2-40B Output: Descr
51 52 53 54 55 56 57 58 59 70 71 72 73 74 75 76 77 78 79	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B Output <poutput< p=""> Output <po< td=""></po<></poutput<>
51 52 53 54 55 56 57 58 59 70 71 72 73 74 75 76 77 78 79 30	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B For the second secon
61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B For the second secon
61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 99 60 81 82	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B For the second secon
31 32 33 34 55 36 37 38 39 31 32 33 34 35 36 37 37 37 37 38 39 31 32 33	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B For the second secon
31 32 33 34 35 36 37 38 37 38 384	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words.</pre>	Captioner: InternVL2-40B For the second secon
31 32 33 34 35 36 37 38 39 31 32 33 34 35 36 37 38 39 30 31 32 33 34 35	<pre>the caption annotation procedure follows the principles of: (1): Describing object attributes, including object quantity, color, material, shape, size, and spatial position (including absolute position in the image and relative position between objects); (2): The annotation process involves just describing the overall scene of the image and some specific object; (3): Instead of describing the imaginary content, only describing the content one can determine confidently from the image. Do not describe the contents by itemizing them in list form. Minimize aesthetic descriptions as much as possible; (4): Please output within 25 words. ''' Answer =''' Aerial view of a suburban area with a mix of commercial buildings, parking lots, and green spaces. A major road curves through the scene, with a large white building near the center. ''' Figure 12: Street view image captioning for p</pre>	Captioner: InternVL2-40B Provide the set of

1188 A.7 TASK EXAMPLES

1190 Satellite Image-Scene Classification (SI-SC)1191

\frown		
p	prompt ='''	
c	lassify the given image into the following classes.	
C +	lasses: smoke, taxiway, cooling_tower, goods_yard,	
f	Coundation_pit, tower_crane, coal_yard, airplane,	Groundtruth:
s	torehouse, cement_concrete_pavement, car, substation, tank,	ship
b	woarding_bridge, apron, untinished_building, breakwater,	
s	hip_lock, chimney, arch_dam, ship, roundabout,	
b	aseball_diamond. \nAnswer with all applicable classes	
s v	''	
a	inswers =	
1	lava_next_llama3: ship	
I	InternVL2-40B: ship	
C	ogvlm2-llama3-chat-19B: smoke, runway,	
1	lava_next_y1_340: snip lava_onevision_gwen2_7b_ov:_ship	
g	pt-40-2024-08-06: ship	WEIGTING HER STOLAND, SAMPLE SPONSOUCH, SPORT, WEIGHT HER
Q	wen-VL-Plus: ship, boat, water	
G	IrbanMLLM: ship	
ree	et View Image-Scene Classification (SVI-SC)	
reo	et View Image-Scene Classification (SVI-SC)	Groundtruth:
red	et View Image-Scene Classification (SVI-SC) prompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings'	Groundtruth: Transportation and utility
	et View Image-Scene Classification (SVI-SC) which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and	Groundtruth: Transportation and utility
	et View Image-Scene Classification (SVI-SC) which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown', Benly	Groundtruth: Transportation and utility
	et View Image-Scene Classification (SVI-SC) which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options.	Groundtruth: Transportation and utility
	<pre>et View Image-Scene Classification (SVI-SC) prompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. '''</pre>	Groundtruth: Transportation and utility
	<pre>et View Image-Scene Classification (SVI-SC) prompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility</pre>	Groundtruth: Transportation and utility
	<pre>et View Image-Scene Classification (SVI-SC) et View Image-Scene Classification (SVI-SC) which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings</pre>	Groundtruth: Transportation and utility
	<pre>prompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings InternVL2-408: 'Transportation and utility' </pre>	Groundtruth: Transportation and utility
	<pre>brompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' cogvlm2-llama3-chat-19B: Mixed residential and commercial buildings</pre>	Groundtruth: Transportation and utility
	<pre>brompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' cogvlm2-llama3-chat-19B: Mixed residential and commercial buildings llava_next_yi_34b: 'Transportation and utility'</pre>	Groundtruth: Transportation and utility
	<pre>brownet =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' cogvlm2-llama3-chat-19B: Mixed residential and commercial buildings llava_next_yi_34b: 'Transportation and utility' llava_onevision_qwen2_7b_ov: Transportation and utility gpt-40-2024-08-06: "Family buildings"</pre>	Groundtruth: Transportation and utility
	<pre>brownet =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILAL.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' cogvIm2-llama3-chat-19B: Mixed residential and commercial buildings llava_next_yi_34b: 'Transportation and utility' llava_onexision_qwen2_7b_ov: Transportation and utility gpt-40-2024-08-06: "Family buildings" Dwen-VL-Plus: 'Mixed residential and commercial buildings' </pre>	<section-header></section-header>
	<pre>brownet =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' cogvIm2-llama3-chat-19B: Mixed residential and commercial buildings llava_next_yi_34b: 'Transportation and utility' llava_onexision_qwen2_7b_ov: Transportation and utility gpt-40-2024-08-06: "Family buildings" When-VL-Plus: 'Mixed residential and commercial buildings' GeoChat: 'Family buildings' </pre>	<section-header></section-header>
	<pre>trigure 15. Submitter intege sector enablines et View Image-Scene Classification (SVI-SC) prompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings InternVL2-408: 'Transportation and utility' cogvlm2-llama3-chat-19B: Mixed residential and commercial buildings llava_next_yi_34b: 'Transportation and utility' llava_onevision_qwen2_7b_ov: Transportation and utility gpt-40-2024-08-06: "Family buildings" Wen-VL-Plus: 'Mixed residential and commercial buildings' GeoChat: 'Family buildings' JrbanMLLM: Transportation and utility</pre>	<section-header><image/></section-header>
	<pre>et View Image-Scene Classification (SVI-SC) et View Image-Scene Classification (SVI-SC) comment = ''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILA1.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' cogvln2-llama3-chat-19B: Mixed residential and commercial buildings llava_next_yi_34b: 'Transportation and utility' llava_onevision_qwen2_7b_ov: Transportation and utility gpt-40-2024-08-06: "Family buildings" Owen-VL-Plus: 'Mixed residential and commercial buildings' JrbanMLLM: Transportation and utility</pre>	<section-header><image/></section-header>
	<pre>trigure 15: Street view image scene classified et View Image-Scene Classification (SVI-SC) prompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility //LA1.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' cogvlm2-llama3-chat-19B: Mixed residential and commercial buildings llava_next_yi_34b: 'Transportation and utility' llava_onevision_gwen2_7b_ov: Transportation and utility gpt-40-2024-08-06: "Family buildings" Wen-VL-Plus: 'Mixed residential and commercial buildings' aeoChat: 'Family buildings' Eigure 14: Street view image scene classified</pre>	Groundtruth: Transportation and utility
	<pre>bright Fist Street view image scene classified bright fist Street view image scene classified bright figure 15: Street view image scene classified bright figure 15: Street view image scene classified bright figure 14: Street view image sce</pre>	froundtruth: Transportation and utility
	<pre>bright Fist Street view image scene classified brompt =''' which scene category does this image fit into? Choose just one from: 'Family buildings', 'Mixed residential and commercial buildings', 'Commercial and office buildings', 'Industrial and manufacturing', 'Transportation and utility', 'Public facilities and institutions', 'Open space and outdoor recreation', 'Vacant land', 'Unknown'. Reply with only one of the quoted options. ''' answers = llava_next_llama3: Transportation and utility /ILAL.5-40b: Mixed residential and commercial buildings InternVL2-40B: 'Transportation and utility' CogvIm2-llama3-chat-19B: Mixed residential and commercial buildings Llava_next_yi_34b: 'Transportation and utility' llava_onevision_qwen2_7b_ov: Transportation and utility gpt-40-2024-08-06: "Family buildings" Dwen-VL-Plus: 'Mixed residential and commercial buildings' GoChat: 'Family buildings' JrbanMLLM: Transportation and utility Figure 14: Street view image scene classified or of the state of the state</pre>	froundtruth: Transportation and utility

1242 Satellite Geo-Localization (SI-GL)

prompt =''' From the options below, which county or administrative region is depicted in this image? Submit only the letter of the correct choice. Groundtruth: A. Ventura County, California В B. Salt Lake County, Utah C. Baltimore County, Maryland D. Maricopa County, Arizona Answer only with A, B, C, or D, without any additional text. Example output: 'A' answers = llava_next_llama3: D VILA1.5-40b: A InternVL2-40B: B cogvlm2-llama3-chat-19B: A llava_next_yi_34b: D llava_onevision_qwen2_7b_ov: B gpt-4o-2024-08-06: B Qwen-VL-Plus: A GeoChat: D UrbanMLLM: B

Figure 15: Satellite image geo-localization results.

1265 Street View Geo-Localization (SVI-GL)

57 58 59 70 71 72 73 74	<pre>prompt =''' What is the correct county or equivalent administrative region for the place shown in this picture? Respond with the letter of the correct choice. A. Cook County, Illinois B. Allegheny County, Pennsylvania C. Baltimore City, Maryland D. Jefferson County, Alabama Answer only with A, B, C, or D, without any additional text. Example output: 'A'</pre>	Groundtruth: B
75	answers =	Mar L
76	llava_next_llama3: B	
77	VILA1.5-40b: A	
78	cogvlm2-llama3-chat-19B: B	
79	llava_next_yi_34b: A	
30	llava_onevision_qwen2_7b_ov: B	The second se
31	gpt-40-2024-08-06: B Owen-VL-Plus: B	A CONTRACT OF A
32	GeoChat: A	
3	UrbanMLLM: B	
2/1		
-		

Figure 16: Street view image geo-localization results.