



CityCube: Benchmarking Cross-view Spatial Reasoning on Vision-Language Models in Urban Environments

Anonymous ACL submission

Abstract

Cross-view spatial reasoning is essential for embodied AI, underpinning spatial understanding, mental simulation and planning in complex environments. Existing benchmarks primarily emphasize indoor or street settings, overlooking the unique challenges of open-ended urban spaces characterized by rich semantics, complex geometries, and view variations. To address this, we introduce **CityCube**, a systematic benchmark designed to probe cross-view reasoning capabilities of current VLMs in urban settings. CityCube integrates four viewpoint dynamics to mimic camera movements and spans a wide spectrum of perspectives from multiple platforms, e.g., vehicles, drones and satellites. For a comprehensive assessment, it features 5,022 meticulously annotated multi-view QA pairs categorized into five cognitive dimensions and three spatial relation expressions. A comprehensive evaluation of 33 VLMs reveals a significant performance disparity with humans: even large-scale models struggle to exceed 54.1% accuracy, remaining 34.2% below human performance. By contrast, small-scale fine-tuned VLMs achieve over 60.0% accuracy, highlighting the necessity of our benchmark. Further analyses indicate the task correlations and fundamental cognitive disparity between VLMs and human-like reasoning.

1 Introduction

Spatial reasoning across viewpoints and scales is fundamental to spatial intelligence. It goes beyond geometric measures (Yang et al., 2024; Cai et al., 2025a), also involves abilities such as relational reasoning (Chen et al., 2024; Song et al., 2025), perspective taking (Piaget, 2013; Li et al., 2025), mental simulation (Eslami et al., 2018), dynamic perception (Tversky, 2019; Ding et al., 2025) and world knowledge recalling (Jia et al., 2025). While humans naturally perform tasks like reasoning 3D scenes from streaming 2D views, replicating this in embodied AI remains challenging.

Recently, Vision-language Models (VLMs) have been increasingly adopted as the cognitive backbone for embodied agents, such as drones, mobile robots and autonomous vehicles (Majumdar et al., 2024). These agents dynamically operate and interact with the physical world, naturally requiring VLMs to perceive and understand multiple views (Hong et al., 2023; Zhang et al., 2024; Zhu et al., 2024; Qi et al., 2024). However, whether current VLMs possess the requisite capabilities for expansive urban spaces remains an open question, as comprehensive evaluations in this domain are still lacking.

To address this, we argue that benchmarks must extend beyond existing indoor-focused settings (Yang et al., 2025a; Du et al., 2024) to the broader context of urban open spaces. As shown in Fig. 1(a), urban environments present unique challenges for cross-view spatial intelligence (CvSI):

- **Richer semantics:** Dense and repetitive instances (e.g., signage and vehicles) pose strict demands on VLMs in spatial grounding and disambiguation.
- **Complex geometries:** Intricate urban 3D structures and road networks demand robust spatial mental reconstruction capabilities.
- **Cross-scale viewpoint variations:** Diverse perspectives, from egocentric to top-down (e.g., drones), necessitate rigorous cross-scale reasoning to maintain spatial consistency.

To the best of our knowledge, CvSI in urban embodied tasks remains underexplored. To address this gap, we introduce **CityCube** (Fig. 1), a systematic benchmark for cross-view urban spatial reasoning. CityCube is constructed by integrating urban views from real-world datasets and realistic simulators, forming a large-scale dataset with 18K images. It covers more than 70 representative

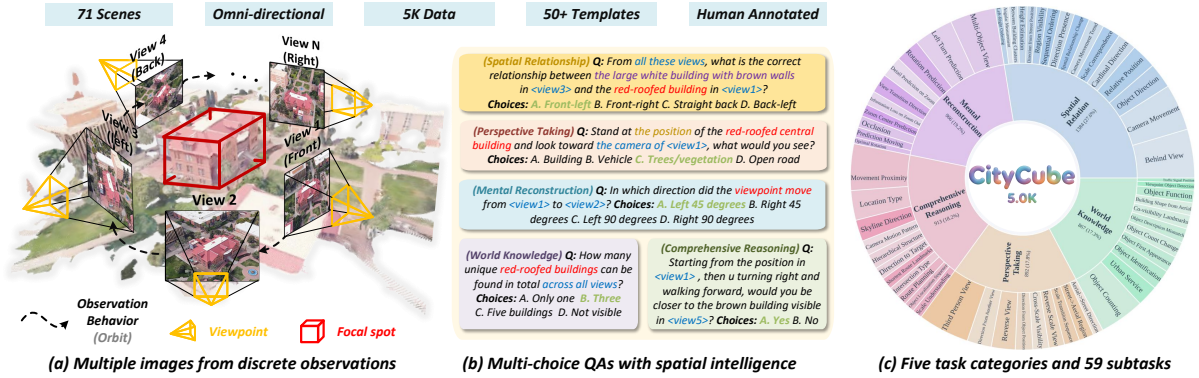


Figure 1: Illustration of the CityCube benchmark. **Left:** An illustration of an embodied orbiting observation, where an agent captures multi-view images by circling a focal object (highlighted in red). **Middle:** Examples of multi-choice QA designed to evaluate five dimensions of CvSI. **Right:** Task distributions on CityCube Benchmark.

cities, and includes multi-view imagery captured from heterogeneous platforms, providing a wide spectrum of first-person and aerial perspectives.

Based on this foundation, we design a comprehensive and challenging suite of spatial reasoning tasks to systematically probe CvSI of urban embodiments. Specifically, CityCube evaluates 5,022 QA pairs with five fundamental spatial intelligences in urban embodied scenes (Cai et al., 2025b), as illustrated in Fig. 1(b): Spatial Relations (SR), Perspective Taking (PT), Mental Reconstruction (MR), World Knowledge (WK), and Comprehensive Reasoning (CR). Each ability requires cross-scale and multi-centric spatial reasoning over diverse scenes, reflecting the strengths and weaknesses of an agent in geometry, semantics, viewpoint transformation, and contextual knowledge.

Upon these tasks, we conduct a large-scale evaluation of 33 mainstream VLMs. As shown in Fig. 1(c), it reveals substantial performance gaps between models, persistent discrepancies with human reasoning, and limitations of existing spatial benchmarks in urban cross-view reasoning. Beyond evaluation, CityCube is further split into training and testing sets. We fine-tune Qwen3-VL of variant scales on the training set using parameter-efficient LoRA, resulting in CityBot-2B, 4B and 8B. Experimental results show the potential of fine-tuning on CityCube in enhancing spatial reasoning, even for relatively small model scales.

In summary, our contributions are threefold:

- **Dataset:** We introduce CityCube, a comprehensive dataset dedicated to cross-view spatial reasoning in urban embodied environments, covering 18.1K images from diverse perspectives across a wide spectrum of urban scenes.

And this work further rearranges the imagery through four observation behavior primitives to mimic camera movements.

- **Benchmark:** We build a challenging CvSI benchmark including 5.0K QA pairs across 59 tasks under five fundamental cognition categories. The queries cover three kinds of relation expressions. Based on these problems, this work conducts a comprehensive evaluation over a diverse set of VLMs.
- **Findings:** We uncover several key findings, not limited to: (i) significance of the benchmark, on which leading proprietary and open-source VLMs exhibit lower than 54.1% accuracy, substantiating the challenge of CvSI, (ii) and disparity between current VLMs and human-like reasoning from correlation analyses and case studies.

2 Related Work

VLM Spatial Intelligence Benchmark VLMs have demonstrated significant potential in depth estimation and spatial cognition on spatial intelligence benchmarks like VSI-Bench (Yang et al., 2025a; Cai et al., 2025a; Cheng et al., 2024; Gholami et al., 2025). However, existing benchmarks (Yin et al., 2025) often overlook the inherent cross-view nature of urban environments, limiting their evaluation scope to a broader applications. While prior works such as ViewSpatial (Li et al., 2025), UrbanVideo (Zhao et al., 2025) and Urbench (Zhou et al., 2025) have explored embodied urban spaces to some extent, they primarily focus on single or restricted perspectives, such as bird’s-eye views.

Table 1: Comparison of the proposed and popular benchmarks for multi-view spatial intelligence. “/” indicates information not mentioned or included, “+” represents sequential operation. “Total Tasks” represents the amount of the task type. *Abbreviations*– *Ped.*: Pedestrian, *Veh.*: Vehicle, *Sat.*: Satellite, *Tem.*: Templates, *Ego.*: Egocentric, *Allo.*: Allocentric, *Exo.*: Exocentric.

Benchmark	Platform	Camera Orientation	QA Num.	Reasoning Annotation	Annotation	Environment	Embodied Questions	Cross Scale	CvSI Categories	Total Tasks
All-Angles (Yeh et al., 2025)	Pedestrian	Ego	2.1K	✗	Human	Indoor	✓	✗	SR, PT	6
MMSI Bench (Yang et al., 2025b)	Pedestrian	Ego, Exo	1K	✓	Human	Indoor	✓	✗	SR, MR, PT	50
ViewSpatial Bench (Li et al., 2025)	Pedestrian	Ego, Exo	5.7K	✗	Rules+Tem.	Indoor & Web images	✓	✗	SR, PT	5
MindCube (Yin et al., 2025)	Pedestrian	Ego, Exo	21.1K	✓	Rules+Tem.	Indoor	✓	✗	MR, PT	5
Ego3D Bench (Gholami et al., 2025)	Vehicle	Ego	8.6K	✓	Rules+Tem.	Outdoor driving	✓	✗	WK	10
OmniSpatial (Jia et al., 2025)	Ped./Vehicle	Ego	1.5K	✓	Human	Web images	✗	✗	PT	50
UrbanFeel (He et al., 2025)	Pedestrian	Ego	14.3K	✗	VFM+Human	Street views	✗	✗	/	11
Urbench (Zhou et al., 2025)	Ped./Satellite	Ego, Allo	11.6K	✗	Rules+LLM+Human	Satellite+Street	✗	✓	/	14
CityCube	Multi-source (Ped., Veh., Drone, Sat.)	Multi-centric (Ego, Exo, Allo)	5.0K	✓	Tem.+LLM+Human	Urban aerial & street views	✓	✓	SR, PT, MR, CR, WK	59

As shown in Tab 1, *CityCube* provides a unique and comprehensive benchmark by integrating multi-source imagery with multi-centric perspective modeling, and supports the richest set of task types with well-annotated reasoning processes.

Urban Visual Question Answering Urban Visual Question Answering (VQA) serves as a critical bridge between multi-modal learning and urban informatics. Unlike traditional urban QA systems that rely on retrieving structured information from static databases (Feng et al., 2025a, 2024), urban VQA emphasizes active perception through multi-modal cues, including vision and language. Existing works focus on perception at different scales. As shown in Tab 1, tasks at the macro level encompass geo-spatial querying, urban governance (Zhou et al., 2025; Feng et al., 2025b), and socioeconomic analysis (He et al., 2025; Liu et al., 2025a; Hao et al., 2024). At the micro level, the focus shifts to agentic perception, including object localization (Zhang et al., 2025) and motion planning (Zhao et al., 2025). However, current benchmarks still exhibit limitations in evaluating multi-image understanding, particularly for tasks requiring advanced cognitive abilities like cross-view reconstruction and perspective transformation in embodied contexts.

3 CityCube Benchmark

3.1 Overview

As depicted in Fig. 1, *CityCube* targets multi-scale spatial reasoning under partial observability and dynamical urban viewpoints. This benchmark is built on extensive real-world urban images (collected from 74 cities across the globe, including Singapore and Boston) and virtual images from 2 high-fidelity urban simulators (i.e. EmbodiedCity, and MatrixCity). In total, it offers **18.1K observation points** spanning diverse viewpoints, scales,

and scene compositions. Building on this image pool, we carefully curate **5.0K QA pairs** to form the dataset, as shown in Fig. 1(c).

Specifically, the benchmark are structured into four well-designed dimensions, as shown in Fig. 2. The images are rearranged four different observation behavior primitives (**dim1**) with three camera perspectives (**dim2**). On this basis, we further identify five critical task categories (**dim3**) and define three spatial relations (**dim4**) for different queries to assess the VLMs from different CvSI abilities. Structurally, the first two determine the visual perspective and content while collecting the images, while the latter two dimensions are manifested in the textual QA pairs. Together, these variables configure the multi-image QA instances, effectively modeling the evaluation of high-level abilities as a superposition of targeted dimensions. The dataset statistics are illustrated in Appendix B.1.

3.2 Dimension 1: Viewpoint Dynamics

To simulate active perception, we classify collected images to a set of behavior primitives: **(1) Rotation**: hovers and observes at a fixed position, changing only the camera orientation (yaw/pitch). **(2) Orbit**: surrounds a landmark or interested region, including planar and volumetric motions. **(3) Ego-Allocentric View**: simulates cognitive alignment from first-person views (egocentric) to third-person perspectives (allocentric). **(4) Dynamic Translation**: mimics coarse-to-fine movements, such as zooming in from a skyscraper to a billboard. It involves multiple scales of urban space.

3.3 Dimension 2: View Perspectives

To ensure robustness across diverse embodiments, we standardize three acquisition protocols: **(1) Ground-level Panorama** represents views from ground vehicles (e.g., autonomous cars), covering discrete directions (front, rear, left, right, and diag-

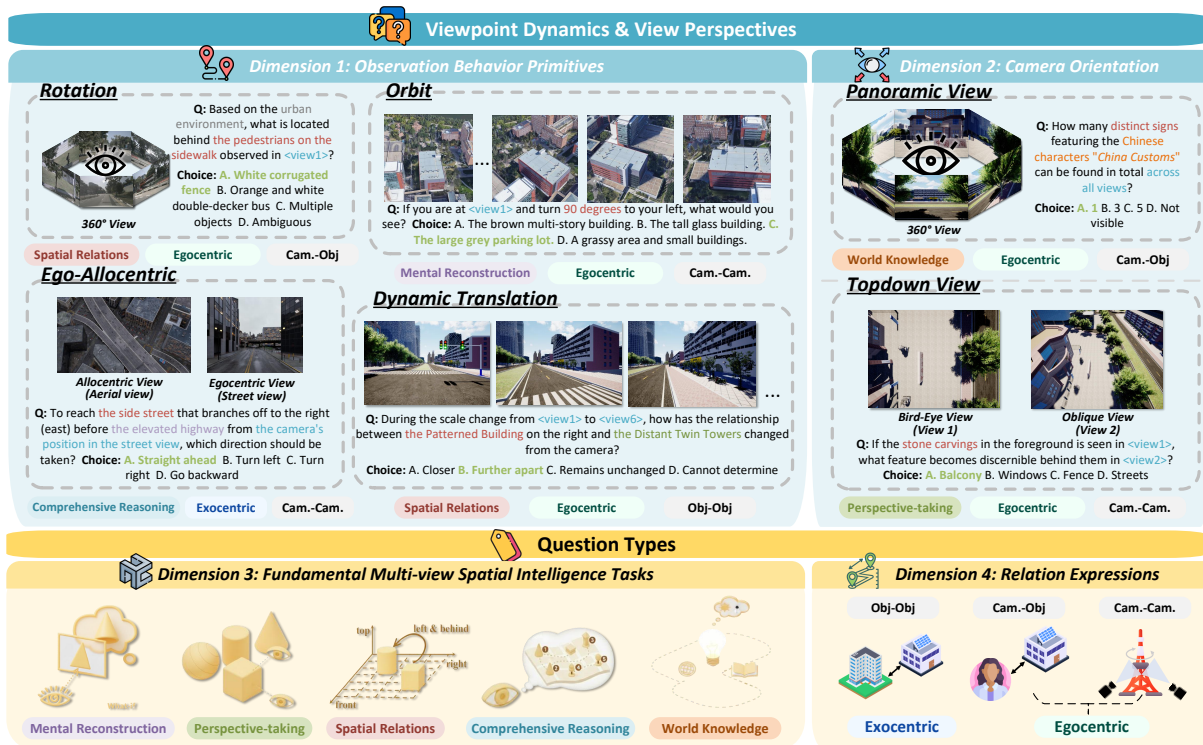


Figure 2: The systematic evaluation protocol of the CityCube benchmark. **Upper Left:** Dim 1 evaluates observation with four representative behavior; **Upper Right:** Dim 2 tests model across various camera orientations; **Bottom Left:** Dim 3 categorizes 59 tasks into 5 fundamental categories; and **Bottom Right:** Dim 4 labels the QA pairs with spatial reference frames.

onals) (Gholami et al., 2025). (2) **Low-altitude Oblique Imagery** represents views from aerial agents like drones, commonly used in 3D urban reconstruction to capture structural facades. (3) **Bird-Eye View** represents a global perspective with minimal occlusion, typical of satellite imagery or high-altitude mapping.

3.4 Dimension 3: Benchmark Tasks

CityCube establishes five fundamental CvSI task categories, built upon existing multi-view reasoning benchmarks in Tab 1 while extending coverage to underexplored dimensions (Cai et al., 2025c):

(1) **Mental Reconstruction** requires VLMs to infer spatial transformation between views by mentally simulating hypothetical movement through the environment.

(2) **Perspective Taking** assesses spatial consistency maintenance of VLMs, focusing on cross-view object grounding and inferring relationship shifts as viewpoints change.

(3) **Spatial Relation** targets the precise quantification of estimating the distance, direction, and topological relations between objects observed across different viewpoints.

(4) **Comprehensive Reasoning** demands **multi-step spatial inference** for a hypothesis testing. For example, hypothetical navigation requires VLMs to mentally execute actions starting from a specific view and verify the predicted targets referencing other observed views.

(5) **World Knowledge** probes urban common-sense of VLMs, such as object geometry, affordance, and visibility. Beyond that, we also design challenges for recalling co-visible landmarks, object counting and scene captions across views.

3.5 Dimension 4: Spatial Relation Expressions

To unify spatial reference frames in QA expressions, CityCube tasks utilize three kind of binary spatial relations: (1) **Object-to-Object (Exocentric)**: textual geometric relationships between external objects viewed from a spectator perspective. (2) **Camera-to-Object (Egocentric)**: spatial relations between the agent’s viewpoint and an observed target. (3) **Camera-to-Camera**: relative transformations between two distinct observation positions (e.g., Camera 1 vs. 2).

Rooted in these four critical dimensions, this design facilitates the construction of a systematic,

interpretable, and fine-grained diagnostic benchmark. This benchmark is specifically engineered to model and evaluate high-level spatial reasoning challenges, including multi-view observation integration, geometric consistency under motion, spatial reference frame alignment, and object identity preservation across scale transitions.

4 Benchmark Curation

This section presents how this work constructs the image dataset and benchmark as depicted in Sec. 3, including two main parts: image processing and QA generation.

4.1 Data Collection and Pre-processing

As shown in Fig. 5, we collect urban images that satisfy predefined patterns and explicitly organize them into multi-view sets. The detailed image source are given in the Appendix B.2. In pre-processing, we adopt two distinct pipelines for real-world datasets and 3D simulators, respectively. For real world, we design a multi-step image filtering and matching procedure to construct high-quality multi-view observations. For simulators, we manually collect viewpoint samples, record camera poses and motion trajectories. Across both real-world and simulated data, we ensure that multiple views correspond to the same underlying scene.

Real-world scenes. All real-world images are sourced from public academic datasets. To obtain ground-level panoramic views mentioned in dim 2, we collect images from nuScenes (Caesar et al., 2020), which covers urban environments in Singapore and Boston. Since driving recording contains a large number of duplicate frames, we apply temporal frame skipping and image-similarity deduplication. In addition, GeoText-1652 (Chu et al., 2024) provides satellite and drone imagery captured from arbitrary viewpoints across 72 cities worldwide. For this dataset, we manually design sampling intervals to obtain approximate orbital observations in dim 1 and further remove redundant images based on image similarity. Detailed data sources and sampling strategies are provided in Appendix B.3.

Simulated scenes. To further expand the scale and diversity of the dataset, we collect additional urban imagery from open-source simulators. EmbodiedCity (Gao et al., 2024) is a 3D urban simulator modeled after Beijing, containing multi-scale city elements ranging from large landmarks and

commercial buildings to fine-grained objects such as bicycles and billboards. Based on this environment, we manually record agent trajectories over large areas following two representative motion patterns mentioned in dim 1 to support dynamic translation views, such as “taking off” (vertical movement) and “approaching coffee shops” (horizontal movement), which are critical for cross-scale spatial reasoning. Besides, we also supplement the dataset with a set of panoramic views captured by piloting the drone. MatrixCity (Li et al., 2023) is a large-scale aerial-street view dataset from a virtual city. We compute geometric projections using precise camera poses to ensure view consistency between aerial and ground-level images, and manually filter image pairs that satisfy ego-allothetic views in dim 1 and 2. More detailed descriptions are provided in Appendix B.4.

4.2 Question-Answer Generation

To enrich multi-choice QA with CvSI tasks, we feed the processed image sets together with structured context into Gemini-2.5 Pro. The generation relies on three key context engineering strategies:

(1) Contextual Role-Playing: we prompt the VLM to serve as an urban embodiment, along with background knowledge of observation behaviors and perspectives. The specifications of prompts are listed in Sec B.5.1.

(2) Template Coverage: we also provide 59 distinct templates for formatting CvSI tasks, ensuring comprehensive coverage across all structural dimensions of the benchmark.

(3) Geometric Reference Injection: we explicitly supply ground-truth geometric information including camera position, orientation angle and image amount into the context. It is vital for model-based QA generation process, enhancing answer credibility, and mitigating hallucinations that violate physical constraints.

To mitigate potential construction biases arising from model, we implement a rigorous two-stage refinement pipeline as follows:

(1) Blind Filtering: we improve blind filtering (Zhao et al., 2025) for a text-only validation to mitigate textual biases. Specifically, an ensemble of models are used to assess questions without visual input. Questions are scored and stratified into difficulty levels (easy, moderate, hard etc) based on accuracy. By eliminating trivial samples, we obtain a QA dataset with a balanced and hierarchical difficulty distribution.

375 (2) **Human Verification:** our annotators verify 424
376 the rationality and authenticity of the QA pairs, 425
377 accepting or rejecting entries accordingly. Ac- 426
378 cepted questions undergo further proofreading to re- 427
379 solve ambiguities, invalid options, or erroneous rea- 428
380 soning, with human-authored reasoning processes 429
381 added to enrich the annotations. The details are 430
382 depicted in Sec B.6.2. 431

383 5 Experiments 432

384 5.1 Evaluation Setups 433

385 **Evaluated Models.** As shown in Tab 2, apart from 434
386 our fine-tuned models, we also evaluate the per- 435
387 formance of 25 models under the CityCube bench- 436
388 mark, including 6 state-of-the-art proprietary mod- 437
389 els, 16 mainstream open-source models and 3 spe- 438
390 cialized VLMs trained for spatial reasoning. No- 439
391 tably, the three spatial models are trained on re- 440
392 spective spatial intelligence benchmarks, i.e. Spa- 441
393 tialvlm (Chen et al., 2024), Omni-spatial (Jia et al., 442
394 2025) and SSRL (Liu et al., 2025b). 443

395 **Evaluation Protocol.** Leveraging the multiple- 444
396 choice format, we compute task-level and overall 445
397 average accuracy in a straightforward manner. For 446
398 the human baseline, we recruited two independent 447
399 groups of ten participants each, all with urban sci- 448
400 ence related academic backgrounds (master’s or 449
401 doctoral students). One group conducted verifi- 450
402 cations in Sec 4.2, while the other performed the 451
403 evaluations, ensuring no overlap between the two. 452
404 More implementation details of evaluation protocol 453
405 are described in Appendix E. 454

406 5.2 Main Results 455

407 Our main observations based on the results shown 456
408 in Table 2 are summarized as follows: 457

409 **Limited cross-view spatial intelligence across** 458
410 **current VLMs.** We find that *both* proprietary and 459
411 open-source VLMs perform poorly on CityCube, 460
412 indicating that CvSI remains largely unsolved by 461
413 existing model families. The best-performing pro- 462
414 prietary model achieves only 54.1% accuracy, ex- 463
415 ceeding the strongest open-source model by 9.2%, 464
416 yet still falling far short of human performance 465
417 (-34.2%). This consistent gap suggests that the 466
418 difficulty lies not in model architecture or scale 467
419 alone, but in the fundamental challenge of multi- 468
420 view spatial reasoning posed by CityCube. 469

421 **Effectiveness of fine-tuning on CityCube.** 470
422 Fine-tuning on CityCube consistently improves 471
423 model performance across all scales. Besides, train- 472

ing with human annotations (e.g., CityBot-4B and 424
-8 B w/ CoT) yields additional gains ($+0.6\%$ to 425
 $+3.6\%$), indicating the benefit of reasoning guid- 426
ance. Notably, the fine-tuned 2B model already 427
surpasses strong proprietary baselines, highlight- 428
ing the advantage of the benchmark. 429

430 **Reasoning-oriented VLMs struggle with spa-** 431
432 **tial tasks.** Despite their success in math and cod- 433
434 ing, reasoning-oriented models (e.g., Qwen3-VL- 434
8B-Thinking, Kimi-VL-A3B-Thinking) show no 435
consistent advantage over non-reasoning models on 436
multi-view spatial tasks such as *PT*. This suggests 437
that generic reasoning supervision alone is insuf- 438
ficient to induce strong spatial reasoning in urban 439
environments. We hypothesize that such reasoning 440
requires explicit modeling of view-dependent ge- 441
ometry and spatial transformations, which is not 442
encouraged by current post-training strategies. 443

444 **CityCube reveals limitations in existing** 444
445 **benchmarks.** General-purpose spatial models like 445
SpaceOM perform sub-optimally on our urban em- 446
bodiment scenarios. This finding validates the short- 447
comings of existing benchmarks and emphasizes 448
the unique value of CityCube in evaluating the dis- 449
tinct challenges of CvSI, which are not covered by 450
standard visual question answering evaluations. 451

452 **Different spatial dimensions exhibit uneven** 452
453 **difficulty.** Models perform relatively well on *WK* 453
454 and *CR* tasks. These tasks mainly rely on semantic 454
understanding and logical inference. In contrast, 455
performance drops significantly on *PT*, *SR*, and *MR*. 456
They require precise geometric reasoning and view- 457
point transformation across multiple views. The 458
results indicate that current VLMs favor semantic 459
priors over robust spatial representations. 460

461 5.3 Task Correlation Analysis 462

463 We posit that tasks requiring similar cognitive ca- 463
464 pabilities will elicit correlated model performance. 464
To analyze the structural relationships among CvSI 465
466 tasks, we compute the Pearson correlation ma- 466
467 trix over 59 tasks evaluated on 25 baseline VLMs. 467
As shown in Fig. 3, we report three key findings: 468

469 **Strong correlations across task categories.** At 469
470 the dimension level, we observe generally substan- 470
471 tial correlations across the five CvSI categories, 471
472 indicating that spatial intelligence is not naturally 472
473 decomposed into independent modules. Among 473
474 all pairs, *MR* and *PT* exhibit the highest inter- 474
dimension correlation ($r = 0.536$), suggesting 475
a shared reliance on underlying cognitive mech- 476
anism. 477

Table 2: Accuracy of 33 VLMs on overall QA pairs. **Only three selected tasks are displayed for each CvSI category besides the overall accuracy due to space limitations.** The best performing model in each category is highlighted **in-bold**, while the second-best is underlined.

Method	Rank	Avg.	World Knowledge				Perspective Taking				Spatial Relation				Mental Recon.				Comp. Reasoning			
			Overall Acc.	Urban Service	Object Ident.	Object Counting	Overall Acc.	Another-view Dir.	Third Person View	Reverse View	Overall Acc.	Object Direction	Relative Pos.	Camera Move.	Overall Acc.	Multi-Obj View	Rotation Pred.	Left-turn Pred.	Overall Acc.	Route Planning	Target Direction	Location Type
<i>Baseline</i>																						
Random	-	22.8	19.2	24.0	3.19	24.4	20.5	18.3	21.9	11.5	25.5	25.0	25.2	25.0	22.0	25.1	15.0	24.2	25.2	16.0	28.1	28.6
Human Level	-	88.3	78.6	85.0	84.0	73.2	87.4	93.0	94.2	96.5	90.2	79.2	84.9	86.4	92.4	84.2	89.4	96.8	93.1	100.0	91.5	86.4
<i>Proprietary Models</i>																						
GPT-5.1-251113	3	53.4	58.3	47.0	<u>70.2</u>	53.1	46.9	32.2	27.7	44.3	51.6	56.7	44.5	42.7	<u>53.7</u>	42.9	52.2	<u>38.2</u>	57.8	14.0	59.3	48.6
Gemini-2.5-Pro	2	<u>53.8</u>	<u>57.9</u>	55.0	57.5	47.0	<u>50.9</u>	33.9	41.3	51.3	<u>50.9</u>	60.8	39.5	30.1	52.0	41.4	43.4	43.4	<u>59.3</u>	<u>24.0</u>	49.2	<u>55.7</u>
Qwen-3-VL-Plus	5	45.2	40.8	42.0	63.8	37.8	44.6	47.0	<u>41.7</u>	52.2	46.5	<u>60.0</u>	35.3	31.8	37.1	44.2	45.1	37.1	56.7	18.0	47.5	52.1
Step-1o-turbo-vision	4	51.8	55.9	42.0	41.3	46.0	48.2	<u>39.1</u>	35.5	47.8	45.8	56.7	39.5	34.6	52.5	41.9	45.1	37.1	59.5	10.0	<u>55.9</u>	57.1
Doubao-seed1.6-251015	1	54.1	57.7	43.0	73.4	<u>48.4</u>	56.3	<u>39.1</u>	51.7	52.2	46.9	55.0	<u>38.7</u>	35.9	54.8	<u>43.3</u>	<u>47.8</u>	36.6	58.8	28.0	<u>55.9</u>	<u>55.7</u>
Skywork-R1V4-Lite	6	40.1	38.6	<u>49.0</u>	43.6	34.3	35.8	27.0	18.2	46.9	34.6	35.8	29.4	19.6	42.9	31.6	43.4	29.6	51.0	8.0	39.0	48.6
<i>Open-source Models</i>																						
Qwen3-VL-8B-Instruct	3	43.1	36.1	20.0	28.7	14.1	37.6	24.4	27.3	25.7	45.8	<u>40.0</u>	30.3	37.3	<u>44.2</u>	37.2	38.1	34.4	49.6	22.0	44.1	30.0
Qwen3-VL-8B-Thinking	9	39.7	<u>41.8</u>	22.0	41.5	<u>35.7</u>	36.3	28.7	22.7	20.4	39.0	35.8	21.9	36.4	39.0	20.9	36.3	30.7	42.9	18.0	35.6	22.1
GLM-4.1V-9B-Base	5	42.6	39.7	19.0	36.2	24.4	36.2	32.2	28.5	31.0	43.0	42.5	<u>33.6</u>	38.6	42.4	28.4	36.3	37.6	51.3	<u>32.0</u>	44.1	<u>32.1</u>
GLM-4.1V-9B-Thinking	1	44.9	45.8	25.0	36.2	40.9	<u>42.8</u>	43.5	24.8	54.0	<u>44.8</u>	36.7	28.6	39.1	40.2	33.5	39.8	28.5	51.5	26.0	39.0	27.1
Kimi-VL-A3B-Instruct	10	39.7	36.1	25.0	12.8	23.5	33.9	27.8	24.0	45.1	36.4	35.8	24.4	35.0	43.8	28.4	<u>40.7</u>	35.5	49.3	10.0	39.0	23.6
Kimi-VL-A3B-Thinking	13	36.0	32.6	18.0	34.0	27.2	31.6	22.6	19.8	27.4	34.6	25.0	26.1	21.8	39.0	25.1	43.4	29.0	42.1	12.0	42.4	27.1
MiMo-VL-7B-SFT	8	40.2	36.9	23.0	30.9	13.2	39.8	32.2	28.9	37.2	38.1	21.7	21.0	37.3	38.9	32.1	31.9	26.3	48.4	16.0	44.1	23.6
MiMo-VL-7B-RL	7	40.9	38.2	21.0	33.0	16.9	39.9	35.7	30.2	30.1	38.4	29.2	26.9	37.3	41.2	34.0	31.9	30.1	48.2	18.0	<u>47.5</u>	23.6
MiniCPM-V-4.5	2	<u>43.9</u>	37.1	20.0	19.2	19.3	44.3	<u>40.0</u>	34.3	33.6	42.6	36.7	26.9	35.0	43.5	<u>35.8</u>	38.1	31.2	<u>52.5</u>	26.0	37.3	27.1
Ovis2.5-9B	4	42.7	40.7	20.0	30.9	24.9	41.1	36.5	<u>33.1</u>	38.1	39.3	20.8	22.7	<u>45.5</u>	40.9	24.7	39.8	31.2	53.2	40.0	<u>47.5</u>	28.6
LLaVA-NeXT-Video-7B	16	28.3	32.6	25.0	<u>42.6</u>	23.0	21.5	25.2	25.2	8.0	25.9	32.5	18.5	23.2	26.3	23.7	26.6	23.1	36.6	20.0	35.6	29.3
LLaVA-Onevision-7B	6	42.3	37.4	22.0	28.7	27.2	34.8	28.7	25.2	<u>48.7</u>	43.2	30.8	28.6	50.5	45.6	24.7	37.2	37.1	49.6	20.0	30.5	29.3
InternVL2.5-8B	12	38.7	36.1	16.0	<u>45.7</u>	14.1	31.7	17.4	18.6	44.3	37.2	34.2	<u>33.6</u>	33.2	40.4	24.2	<u>40.7</u>	37.6	48.2	10.0	50.9	25.0
Skywork-VL-Reward-7B	14	33.2	36.8	5.0	41.5	28.2	34.9	27.0	24.4	36.3	22.8	2.5	21.0	21.8	32.7	21.9	23.9	19.4	44.6	28.0	35.6	19.3
Molmo-7B-D-0924	11	38.7	33.3	17.0	25.5	<u>30.5</u>	33.5	31.3	20.7	30.1	36.8	28.3	34.5	32.3	41.8	26.5	33.6	28.0	48.2	26.0	44.1	33.6
Phi-4-multimodal-instruct	15	32.0	31.8	23.0	46.8	16.4	29.0	33.9	17.8	40.7	27.3	11.7	16.8	25.0	32.0	28.4	15.0	24.7	42.0	18.0	32.2	17.1
<i>Spatial Models</i>																						
Spatial-SSRL-4B	1	39.8	39.2	<u>22.0</u>	35.1	15.0	31.2	<u>31.3</u>	<u>26.5</u>	21.2	41.6	<u>30.0</u>	31.1	47.3	37.8	30.2	<u>35.4</u>	24.7	<u>48.0</u>	34.0	49.2	27.1
SpaceOm-4B	2	<u>38.9</u>	<u>38.6</u>	25.0	40.4	23.5	34.3	33.9	23.1	<u>48.7</u>	<u>37.1</u>	35.0	19.3	<u>33.2</u>	<u>37.9</u>	26.5	36.3	<u>34.9</u>	47.2	8.0	<u>45.8</u>	35.0
SpaceThinker-3B	3	38.7	35.8	21.0	<u>38.3</u>	<u>19.3</u>	<u>34.2</u>	25.2	29.8	55.8	35.6	25.0	<u>29.4</u>	30.9	40.5	<u>29.8</u>	33.6	36.0	48.5	<u>18.0</u>	33.9	<u>33.6</u>
<i>Fine-Tuning: Test set</i>																						
Qwen3-VL-2B (before)	9	30.4	26.4	30.0	10.0	22.7	23.4	25.0	12.0	33.3	37.1	50.0	<u>41.7</u>	31.8	25.4	13.6	16.7	26.3	36.7	20.0	50.0	28.6
CityBot-2B (CoT)	4	60.2	61.5	40.0	70.0	59.1	62.8	48.3	56.0	<u>83.3</u>	60.1	75.0	<u>41.7</u>	45.5	50.9	<u>31.8</u>	38.3	31.6	<u>67.4</u>	<u>60.0</u>	50.0	57.1
CityBot-2B (w/o CoT)	6	55.8	58.2	40.0	50.0	68.2	50.0	41.7	28.0	<u>75.0</u>	57.3	66.7	<u>41.7</u>	45.5	46.1	18.2	33.3	42.1	<u>67.4</u>	40.0	50.0	60.0
Qwen3-VL-4B (before)	8	36.6	38.5	20.0	60.0	40.9	27.7	41.7	28.0	12.5	37.8	50.0	25.0	36.4	39.2	27.3	<u>41.7</u>	15.8	38.8	40.0	33.3	28.6
CityBot-4B (CoT)	2	<u>61.0</u>	67.0	50.0	70.0	<u>72.7</u>	<u>59.6</u>	50.0	48.0	66.7	58.1	66.7	<u>41.7</u>	<u>50.0</u>	<u>54.9</u>	27.3	<u>41.7</u>	68.4	<u>67.4</u>	80.0	50.0	71.4
CityBot-4B (w/o CoT)	3	60.4	62.6	34.3	50.0	68.2	53.2	33.3	28.0	75.0	<u>58.0</u>	66.7	<u>41.7</u>	45.5	59.8	31.8	<u>41.7</u>	<u>52.6</u>	69.4	40.0	83.3	64.3
Qwen3-VL-8B (before)	7	37.1	40.7	20.0	60.0	45.5	26.6	8.3	16.0	58.3	42.0	66.7	33.3	31.8	38.2	<u>22.7</u>	<u>41.7</u>	26.3	35.7	40.0	50.0	57.1
CityBot-8B (CoT)	1	61.4	<u>64.8</u>	50.0	70.0	77.3	62.8	41.7	<u>52.0</u>	91.7	<u>58.0</u>	75.0	50.0	54.6	52.9	40.9	50.0	36.8	64.8	<u>60.0</u>	<u>66.7</u>	71.4
CityBot-8B (w/o CoT)	5	57.8	58.2	40.0	50.0	68.2	54.3	50.0	32.0	<u>83.3</u>	57.3	75.0	33.3	45.5	53.9	<u>31.8</u>	25.0	<u>52.6</u>	65.3	20.0	<u>66.7</u>	50.0

Dense correlations across different tasks. At a finer granularity, task-level analysis reveals dense cross-category correlations, with 81.4% of high-correlation pairs spanning different CvSI dimensions. For example, *Behind View* (T33, SR) and *Left-turn Prediction* (T15, CR) exhibit near-perfect correlation ($r = 0.954$), indicating that panoramic spatial perception is tightly coupled with dynamic viewpoint reasoning. Similarly, strong associations are observed between cross-scale semantic tasks such as *Hierarchical Structure* (T3, WK) and *Object Localization Sequence* (T8, CR) ($r = 0.965$), reflecting synchronized semantic-geometric reasoning under dynamic scale changes.

Metric estimation as a weakly correlated ca-

pability. Metric estimation shows weaker coupling with other tasks. For example, *Height Estimation* (T40) shows negligible correlation with most other tasks, suggesting that precise metric reasoning constitutes a distinct capability, weakly linked to view-dependent or semantic spatial reasoning processes.

5.4 The VLMs Error Analysis

Through a case-by-case qualitative investigation (as depicted in Appendix D), we identify four primary failure modes of current VLMs in urban spatial reasoning:

Limited sensitivity to small-scale urban entities. In complex urban environments, models frequently fail to recognize small-scale urban en-

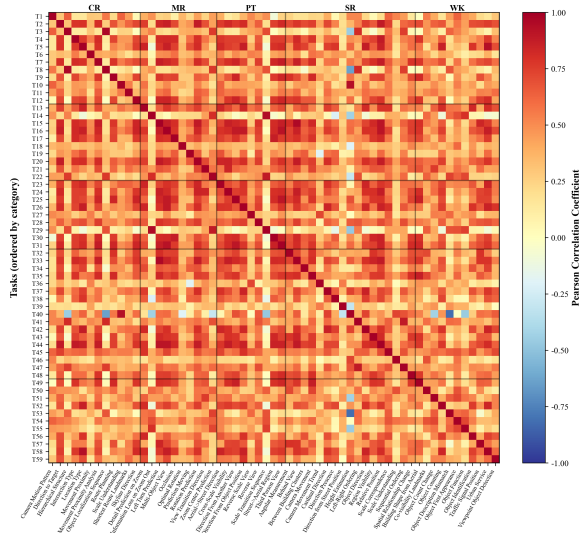


Figure 3: Task correlation matrix. Each axis corresponds to the 59 tasks after classification, and color intensity indicates the strength of correlations.

504 titles or entirely overlook critical landmarks. This
 505 issue is particularly pronounced in scenes with high
 506 semantic density, where numerous visually similar
 507 objects compete for attention.

508 **Failures in egocentric spatial reasoning.** Mod-
 509 els struggle to correctly predict changes in relative
 510 position, orientation, and visibility when the ego-
 511 centric viewpoint shifts. This failure in “mental
 512 rotation” indicates a lack of a reliable internal spa-
 513 tial representation.

514 **Insufficient cross-view consistency.** Models
 515 exhibit difficulty in establishing correct correspon-
 516 dences across perspectives. For instance, models
 517 often mismatch a building in street imagery with
 518 an unrelated structure in aerial views, leading to
 519 erroneous reasoning about spatial relationships and
 520 world knowledge.

521 **Misinterpretation of motion and scale dynam-
 522 ics.** Models struggle to interpret the camera move-
 523 ments (e.g., forward/backward translation or di-
 524 rectional turns) and the resulting dynamic scale
 525 changes. Such failures in motion understanding di-
 526 rectly impede the ability to perform complex scene
 527 reconstruction and navigation-related reasoning.

528 5.5 Human-AI Difficulty Bias Analysis

529 To examine whether the difficulty of spatial tasks
 530 for humans aligns with VLM performance, we an-
 531alyze the correlation between human baseline accu-
 532racy and the average performance of evaluated
 533 VLMs across 59 subtasks. Fig. 4 reveals a Pear-
 534son correlation coefficient of $r = 0.098$ with a

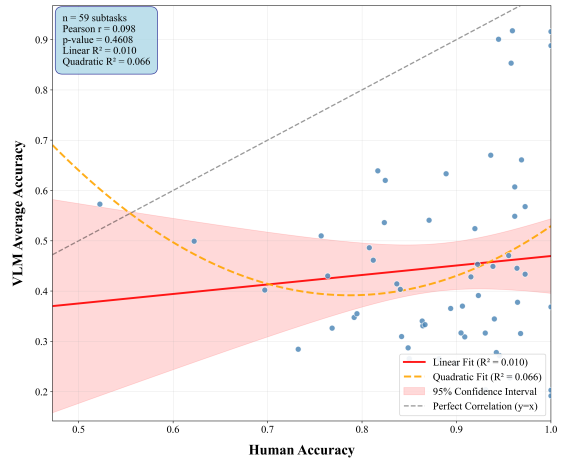


Figure 4: Human vs AI task correlation. The scatters illustrate the performance correlation across 59 tasks, where each individual point represents a specific task.

535 p-value of 0.4608. This extremely low correla-
 536 tion ($R^2 = 0.010$) indicates that tasks found dif-
 537 ficult by VLMs do not necessarily pose a chal-
 538 lenge for human, and vice versa. Notably, the
 539 AI performance distribution is highly dispersed
 540 (range: [19.2%, 91.8%]), suggesting that VLMs
 541 might merely sensitive to low-level visual features
 542 rather than spatial mental modeling. This diver-
 543 gence confirms that CityCube captures unique spa-
 544 tial challenges that are non-trivial for current model
 545 architectures, despite being intuitive for humans.

546 6 Conclusion

547 In this paper, we presented **CityCube**, a compre-
 548 hensive benchmark specifically designed to eval-
 549 uate CvSI of VLMs in urban environments. City-
 550 Cube encompasses 59 tasks across 5 cognitive cat-
 551 egories, supported by a large-scale collection of
 552 images from dynamic viewpoints and diverse ori-
 553 entations. The resulting dataset contains 5,022
 554 multiple-choice questions (MCQs), each rigorously
 555 annotated and verified by humans.

556 Our extensive evaluation of 33 VLMs reveals
 557 that CvSI remains extremely challenging, even for
 558 very large-scale models, “thinking” models and
 559 specialized spatial models. While our fine-tuned
 560 **CityBot** models (based on 2B, 4B and 8B back-
 561 bones) outperform leading proprietary models, a
 562 substantial gap to human-level spatial cognition
 563 persists. We hope that CityCube can serve as a
 564 foundation for future studies on spatially grounded
 565 learning paradigms and as a diagnostic tool for
 566 developing next-generation VLMs with stronger
 567 urban spatial intelligence.

7 Limitations

This work has several limitations. First, although CityCube covers diverse cross-view spatial tasks, we do not explicitly isolate perspective-induced cognitive biases, such as viewpoint asymmetry and reference-frame ambiguity, which may systematically influence spatial reasoning performance. Second, CityCube is constructed in a simulated urban environment; the impact of Sim-to-Real transfer and domain gaps is not evaluated in this study. Third, our analysis focuses on task-level performance and does not probe internal representations or multi-view fusion mechanisms, limiting interpretability of model failure modes. Finally, while we hypothesize the importance of explicit spatial supervision, we do not implement or validate dedicated post-training strategies (e.g., spatial Chain-of-Thought) in this work. Future work will address these issues by analyzing viewpoint biases, extending evaluation to Sim2real settings, and exploring spatially grounded architectural and post-training designs.

8 Ethics Statement

This research exclusively utilizes publicly available datasets, programs, and pre-trained models. All data annotation was conducted with informed consent from participants, and the datasets do not contain information that could compromise individual privacy or public safety. All procedures strictly adhere to the guidelines established by the ACL Code of Ethics. Therefore, this work does not raise any ethical concerns.

References

Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas Guibas. 2020. Referit3d: Neural listeners for fine-grained 3d object identification in real-world scenes. In *European conference on computer vision*, pages 422–440. Springer.

Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. 2022. Scanqa: 3d question answering for spatial scene understanding. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19129–19139.

Shuai Bai, Yuxuan Cai, and Ke Zhu. 2025. [Qwen3-vl technical report](#). *Preprint*, arXiv:2511.21631.

Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. 2020.

nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631.

Wenxiao Cai, Iaroslav Ponomarenko, Jianhao Yuan, Xiaoyi Li, Wankou Yang, Hao Dong, and Bo Zhao. 2025a. Spatialbot: Precise spatial understanding with vision language models. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9490–9498. IEEE.

Zhongang Cai, Ruisi Wang, and Lei Yang. 2025b. Scaling spatial intelligence with multimodal foundation models. *arXiv preprint arXiv:2511.13719*.

Zhongang Cai, Yubo Wang, Qingping Sun, Ruisi Wang, Chenyang Gu, Wanqi Yin, Zhiqian Lin, Zhitao Yang, Chen Wei, Xuanke Shi, and 1 others. 2025c. Has gpt-5 achieved spatial intelligence? an empirical study. *arXiv preprint arXiv:2508.13142*, 3.

Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. 2017. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*.

Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 2024. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465.

An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. 2024. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093.

Meng Chu, Zhedong Zheng, Wei Ji, Tingyu Wang, and Tat-Seng Chua. 2024. Towards natural language-guided drones: Geotext-1652 benchmark with spatial relation matching. In *European Conference on Computer Vision*, pages 213–231. Springer.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839.

Jingtao Ding, Yunke Zhang, Yu Shang, Yuheng Zhang, Zefang Zong, Jie Feng, Yuan Yuan, Hongyuan Su, Nian Li, Nicholas Sukiennik, and 1 others. 2025.

671	Understanding world or predicting future? a comprehensive survey of world models. <i>ACM Computing Surveys</i> , 58(3):1–38.	
672		
673		
674	Mengfei Du, Binhao Wu, Zejun Li, Xuan-Jing Huang, and Zhongyu Wei. 2024. Embspatial-bench: Benchmarking spatial understanding for embodied tasks with large vision-language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 346–355.	
675		
676		
677		
678		
679		
680		
681	SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, and 1 others. 2018. Neural scene representation and rendering. <i>Science</i> , 360(6394):1204–1210.	
682		
683		
684		
685		
686		
687	Jie Feng, Tianhui Liu, Yuwei Du, Siqi Guo, Yuming Lin, and Yong Li. 2025a. Citygpt: Empowering urban spatial cognition of large language models. In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2</i> , pages 591–602.	
688		
689		
690		
691		
692		
693	Jie Feng, Shengyuan Wang, Tianhui Liu, Yanxin Xi, and Yong Li. 2025b. Urbanllava: A multi-modal large language model for urban intelligence with spatial reasoning and understanding. <i>arXiv preprint arXiv:2506.23219</i> .	
694		
695		
696		
697		
698	Jie Feng, Jun Zhang, Junbo Yan, Xin Zhang, Tianjian Ouyang, Tianhui Liu, Yuwei Du, Siqi Guo, and Yong Li. 2024. Citybench: Evaluating the capabilities of large language model as world model. <i>arXiv e-prints</i> , pages arXiv–2406.	
699		
700		
701		
702		
703	Rao Fu, Jingyu Liu, Xilun Chen, Yixin Nie, and Wenhan Xiong. 2024. Scene-llm: Extending language model for 3d visual understanding and reasoning. <i>arXiv preprint arXiv:2403.11401</i> .	
704		
705		
706		
707	Chen Gao, Baining Zhao, Weichen Zhang, Jinzhu Mao, Jun Zhang, Zhiheng Zheng, Fanhang Man, Jianjie Fang, Zile Zhou, Jinqiang Cui, and 1 others. 2024. Embodiedcity: A benchmark platform for embodied agent in real-world city environment. <i>arXiv preprint arXiv:2410.09604</i> .	
708		
709		
710		
711		
712		
713	Mohsen Gholami, Ahmad Rezaei, Zhou Weimin, Sitong Mao, Shunbo Zhou, Yong Zhang, and Mohammad Akbari. 2025. Spatial reasoning with vision-language models in ego-centric multi-view scenes. <i>arXiv preprint arXiv:2509.06266</i> .	
714		
715		
716		
717		
718	Dong Guo, Faming Wu, and Zuquan Song. 2025. Seed1.5-vl technical report . <i>Preprint</i> , arXiv:2505.07062.	
719		
720		
721	Zoey Guo, Yiwen Tang, Ray Zhang, Dong Wang, Zhigang Wang, Bin Zhao, and Xuelong Li. 2023. Viewrefer: Grasp the multi-view knowledge for 3d visual grounding. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 15372–15383.	
722		
723		
724		
725		
726		
	Xixuan Hao, Wei Chen, Yibo Yan, Siru Zhong, Kun Wang, Qingsong Wen, and Yuxuan Liang. 2024. Urbanvlp: A multi-granularity vision-language pre-trained foundation model for urban indicator prediction. <i>CoRR</i> .	727
		728
		729
		730
		731
	Jun He, Yi Lin, Zilong Huang, Jiacong Yin, Junyan Ye, Yuchuan Zhou, Weijia Li, and Xiang Zhang. 2025. Urbanfeel: A comprehensive benchmark for temporal and perceptual understanding of city scenes through human perspective. <i>arXiv preprint arXiv:2509.22228</i> .	732
		733
		734
		735
		736
		737
	Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 2023. 3d-llm: Injecting the 3d world into large language models. <i>Advances in Neural Information Processing Systems</i> , 36:20482–20494.	738
		739
		740
		741
		742
	Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. 2024. An embodied generalist agent in 3d world. In <i>Proceedings of the 41st International Conference on Machine Learning</i> , pages 20413–20451.	743
		744
		745
		746
		747
		748
	Mengdi Jia, Zekun Qi, Shaochen Zhang, Wenyao Zhang, Xinqiang Yu, Jiawei He, He Wang, and Li Yi. 2025. Omnispatial: Towards comprehensive spatial reasoning benchmark for vision language models. <i>arXiv preprint arXiv:2506.03135</i> .	749
		750
		751
		752
		753
	Dingming Li, Hongxing Li, Zixuan Wang, Yuchen Yan, Hang Zhang, Siqi Chen, Guiyang Hou, Shengpei Jiang, Wenqi Zhang, Yongliang Shen, and 1 others. 2025. Viewspatial-bench: Evaluating multi-perspective spatial localization in vision-language models. <i>arXiv preprint arXiv:2505.21500</i> .	754
		755
		756
		757
		758
		759
	Yixuan Li, Lihan Jiang, Linning Xu, Yuanbo Xiangli, Zhenzhi Wang, Dahua Lin, and Bo Dai. 2023. Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 3205–3215.	760
		761
		762
		763
		764
		765
	Tianhui Liu, Jie Feng, Hetian Pang, Xin Zhang, Tianjian Ouyang, Zhiyuan Zhang, and Yong Li. 2025a. Citylens: Benchmarking large language-vision models for urban socioeconomic sensing. <i>arXiv preprint arXiv:2506.00530</i> .	766
		767
		768
		769
		770
	Yuhong Liu, Beichen Zhang, Yuhang Zang, Yuhang Cao, Long Xing, Xiaoyi Dong, Haodong Duan, Dahua Lin, and Jiaqi Wang. 2025b. Spatial-ssrl: Enhancing spatial understanding via self-supervised reinforcement learning. <i>arXiv preprint arXiv:2510.27606</i> .	771
		772
		773
		774
		775
		776
	Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav Putta, Sriram Yenamandra, Mikael Henaff, Sneha Silwal, Paul Mccvay, Oleksandr Maksymets, Sergio Arnaud, and 1 others. 2024. Openeqa: Embodied question answering in the era of foundation models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 16488–16498.	777
		778
		779
		780
		781
		782
		783
		784

785	OpenAI. 2025. Gpt-5 system card . Accessed: 2026-01-03.	839
786		840
787	Jean Piaget. 2013. <i>Child’s conception of space: Selected works vol 4</i> . Routledge.	841
788		842
789	Zekun Qi, Runpei Dong, Shaochen Zhang, Haoran Geng, Chunrui Han, Zheng Ge, Li Yi, and Kaisheng Ma. 2024. Shapellm: Universal 3d object understanding for embodied interaction. In <i>European Conference on Computer Vision</i> , pages 214–238. Springer.	843
790		844
791		845
792		846
793		847
794	Chan Hee Song, Valts Blukis, Jonathan Tremblay, Stephen Tyree, Yu Su, and Stan Birchfield. 2025. Robospatial: Teaching spatial understanding to 2d and 3d vision-language models for robotics. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 15768–15780.	848
795		849
796		850
797		851
798		852
799		853
800	Barbara Tversky. 2019. <i>Mind in motion: How action shapes thought</i> . Basic Books.	854
801		855
802	Tianwen Wei, Liang Zhao, and Yahui Zhou. 2023. Skywork: A more open bilingual foundation model . Preprint, arXiv:2310.19341.	856
803		857
804		858
805	Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. 2025a. Thinking in space: How multimodal large language models see, remember, and recall spaces. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 10632–10643.	859
806		860
807		861
808		
809		
810		
811	Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. 2024. Depth anything: Unleashing the power of large-scale unlabeled data. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10371–10381.	
812		
813		
814		
815		
816		
817	Sihan Yang, Runsen Xu, Yiman Xie, Sizhe Yang, Mo Li, Jingli Lin, Chenming Zhu, Xiaochen Chen, Haodong Duan, Xiangyu Yue, and 1 others. 2025b. Mmsi-bench: A benchmark for multi-image spatial intelligence. <i>arXiv preprint arXiv:2505.23764</i> .	
818		
819		
820		
821		
822	Chun-Hsiao Yeh, Chenyu Wang, Shengbang Tong, Ta-Ying Cheng, Ruoyu Wang, Tianzhe Chu, Yuexiang Zhai, Yubei Chen, Shenghua Gao, and Yi Ma. 2025. Seeing from another perspective: Evaluating multi-view understanding in mllms. <i>arXiv preprint arXiv:2504.15280</i> .	
823		
824		
825		
826		
827		
828	Baiqiao Yin, Qineng Wang, Pingyue Zhang, Jianshu Zhang, Kangrui Wang, Zihan Wang, Jieyu Zhang, Keshige Yan Chandrasegaran, Han Liu, Ranjay Krishna, and 1 others. 2025. Spatial mental modeling from limited views. In <i>Structural Priors for Vision Workshop at ICCV’25</i> .	
829		
830		
831		
832		
833		
834	Sha Zhang, Di Huang, Jiajun Deng, Shixiang Tang, Wanli Ouyang, Tong He, and Yanyong Zhang. 2024. Agent3d-zero: An agent for zero-shot 3d understanding. In <i>European Conference on Computer Vision</i> , pages 186–202. Springer.	
835		
836		
837		
838		
	Weichen Zhang, Zile Zhou, Zhiheng Zheng, Chen Gao, Jinqiang Cui, Yong Li, Xinlei Chen, and Xiaoping Zhang. 2025. Open3dvqa: A benchmark for comprehensive spatial reasoning with multimodal large language model in open space. <i>arXiv preprint arXiv:2503.11094</i> .	
	Baining Zhao, Jianjie Fang, Zichao Dai, Ziyong Wang, Jirong Zha, Weichen Zhang, Chen Gao, Yue Wang, Jinqiang Cui, Xinlei Chen, and 1 others. 2025. Urbanvideo-bench: Benchmarking vision-language models on embodied intelligence with video data in urban spaces. <i>arXiv preprint arXiv:2503.06157</i> .	
	Baichuan Zhou, Haote Yang, Dairong Chen, Junyan Ye, Tianyi Bai, Jinhua Yu, Songyang Zhang, Dahua Lin, Conghui He, and Weijia Li. 2025. Urbanbench: A comprehensive benchmark for evaluating large multimodal models in multi-view urban scenarios. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 10707–10715.	
	Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang, and Xihui Liu. 2024. Llava-3d: A simple yet effective pathway to empowering llms with 3d-awareness. <i>arXiv preprint arXiv:2409.18125</i> .	

862 A Appendix

863 In the supplementary materials, we provide the
864 following:

- 865 • Details of data collection, processing, auto-
866 matic generation and refinement of **CityCube**
867 Benchmark (Sec B).
- 868 • Details of reproduction, including code for
869 QA generation, model training scripts, im-
870 age datasets, and the detail model information
871 (Sec C).
- 872 • Details of VLMs error and their correspond-
873 ing reason (Sec D).
- 874 • Details of experiment setup, including train-
875 ing hyperparameter, dataset settings and hu-
876 man evaluation UI interface (Sec E).
- 877 • Further discussion on relative research, espe-
878 cially in 3D QA (Sec F).

879 B Dataset Generation Pipeline

880 B.1 Dataset Visualization

881 We analyzed the distribution of spatial reference
882 frames and spatial task frequencies associated with
883 each observation behavior, with the results summa-
884 rized in the stacked bar chart in Fig. 6. The training
885 and test sets contain approximately 4.5k and 0.5k
886 questions, respectively. We construct each split
887 using a stratified sampling strategy, ensuring that
888 the proportion of each task type remains consis-
889 tent with the overall QA distribution of the dataset.
890 We further present word cloud visualizations of the
891 textual content in the training set, as illustrated in
892 Fig. 7.

893 B.2 Details of Image Acquisition

894 Our sources span both real-world environments
895 and photorealistic virtual scenes. We leverage real-
896 world sensor datasets—GeoText-1652 (Chu et al.,
897 2024) and nuScenes (Caesar et al., 2020), to ac-
898 quire raw data for perspectives involving camera
899 behavior primitives of Rotation and Orbit. We fur-
900 ther augment the benchmark with extensive first-
901 person to third-person camera motion imagery de-
902 rived from embodied vision, which are based on
903 virtual environments—EmbodiedCity (Gao et al.,
904 2024) and MatrixCity (Li et al., 2023).

The core datasets are selected for their com- 905
plementary strengths in providing diverse, high- 906
quality visual data while adhering to rigorous data 907
ethics standards. Key details are outlined below: 908

- 909 • **nuScenes (Caesar et al., 2020)**: This multi- 910
modal autonomous driving dataset provides 911
1,000 large-scale scenes from real-world ur- 912
ban environments (Boston and Singapore). Its 913
synchronized data from six cameras, LiDAR, 914
and radar, along with comprehensive 3D an- 915
notations and precise calibration, offers a rich, 916
real-world foundation for studying multi-view 917
geometry and complex camera motions.
- 918 • **GeoText-1652 (Chu et al., 2024)**: This bench- 919
mark extends real-world imagery with spa- 920
tial language annotations. Based on the es- 921
tablished University-1652 image set, it pro- 922
vides high-resolution ground-level and aerial 923
images. Its key contribution to our work is 924
the precise 360-degree image correspondence 925
the aerial scan of urban buildings, facilitating 926
tasks that require spatial and linguistic ground- 927
ing.
- 928 • **MatrixCity (Li et al., 2023)**: A large-scale 929
photorealistic synthetic dataset built with Un- 930
real Engine 5. It provides over 500,000 street- 931
view and aerial images with pixel-perfect 932
ground-truth information (e.g., camera poses, 933
depth, normal maps) and full control over en- 934
vironmental conditions (weather, lighting). It 935
serves as a primary, privacy-safe source for 936
generating diverse first-person camera motion 937
trajectories.
- 938 • **EmbodiedCity (Gao et al., 2024)**: This 939
benchmark platform supports embodied AI 940
agents in city-scale environments. It enables 941
the generation of extensive, realistic first- 942
person visual experience data for navigation 943
and interaction tasks, perfectly aligning with 944
our need for embodied, egocentric visual data 945
in complex virtual urban settings.

Privacy and Ethical Compliance Statement. We 946
strictly adhere to data privacy and ethical research 947
standards. All real-world data utilized in this study 948
(nuScenes and GeoText-1652) are sourced from es- 949
tablished, publicly released academic datasets that 950
have undergone formal anonymization and cura- 951
tion processes, involving no collection of personal 952
identifiers. Crucially, a significant portion of our 953

Algorithm 1: Multi-view Frame Temporal Sampling and Deduplication for nuScenes

Input: Driving dataset \mathcal{D} with scenes and synchronized multi-camera frames;
Temporal step size s ; similarity threshold τ ;
optional per-scene limit M

Output: A set of exported multi-view frame groups with metadata

```
foreach scene  $\mathcal{S}$  in  $\mathcal{D}$  do
  Retrieve ordered frame list
   $\{\mathbf{F}_1, \mathbf{F}_2, \dots, \mathbf{F}_N\}$ ;
  Initialize  $\mathbf{A}_{\text{prev}} \leftarrow \emptyset$ , counter  $k \leftarrow 0$ ;
  for  $i = 1, 1 + s, \dots, N$  do
    if  $M$  is set and  $k \geq M$  then
       $\perp$  break
    Let  $\mathbf{F}_i = \{I_i^{(1)}, \dots, I_i^{(6)}\}$  be six
      synchronized camera images;
    if any image in  $\mathbf{F}_i$  is missing then
       $\perp$  continue
    Construct concatenate view
     $\mathbf{A}_i \leftarrow \text{CONCATE}(\mathbf{F}_i)$ ;
    if  $\mathbf{A}_{\text{prev}} \neq \emptyset$  then
      Compute appearance difference
      
$$d \leftarrow \frac{1}{|\mathbf{A}_i|} \sum |\mathbf{A}_i - \mathbf{A}_{\text{prev}}|$$

      if  $d < \tau$  then
         $\perp$  continue
    Export all images in  $\mathbf{F}_i$  as one
      multi-view sample;
    Save associated metadata (scene ID,
      frame index, camera list);
     $\mathbf{A}_{\text{prev}} \leftarrow \mathbf{A}_i$ ;
     $k \leftarrow k + 1$ ;
```

983 preserves scene diversity while significantly reduc-
984 ing redundant observations. The retained frames
985 are stored as structured multi-view samples, each
986 associated with metadata describing scene identity
987 and camera configuration.

988 B.3.2 Geotext-1652

989 Unlike driving datasets with strong temporal con-
990 tinuity, images in GeoText-1652 are organized as
991 viewpoint sequences centered around prominent
992 landmarks, exhibiting systematic variations in cam-
993 era distance and altitude.

994 To obtain representative landmark-centric multi-
995 view observations while avoiding redundant sam-
996 ples, we adopt a rule-based sampling strategy with

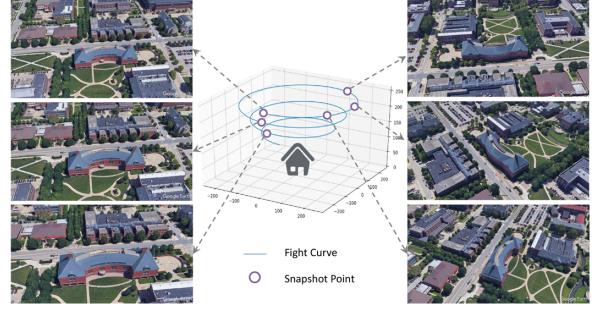


Figure 8: A data example of viewpoint dynamic in Geotext-1652.

manually designed intervals. Specifically, for each landmark sequence, we extract two complementary subsets: (i) *orbit views*, which capture appearance changes under small viewpoint variations at similar altitudes, and (ii) *spiral views*, which emphasize significant altitude and scale changes. These two subsets jointly approximate local and global perceptual variations around a landmark.

Given an ordered image sequence, orbit views are sampled with a short interval to retain fine-grained viewpoint diversity, while spiral views are sampled with a larger interval to highlight elevation changes. To further control redundancy and balance the dataset, we impose an upper bound on the number of retained images per subset. The selected images are exported into separate directories according to their view type, preserving the original filenames for traceability. This strategy yields compact yet diverse multi-view observations suitable for landmark-centric perception and reasoning tasks.

B.3.3 Cross-view Geometric Pairing in MatrixCity

MatrixCity is a large-scale virtual city dataset providing synchronized aerial and street-level imagery with precise camera poses. Unlike real-world datasets where camera geometry may be noisy or partially missing, MatrixCity offers accurate extrinsic parameters, enabling explicit geometric reasoning between ego-centric (street) and allocentric (aerial) views.

To construct reliable aerial-street image pairs, we perform geometry-aware cross-view matching based on camera poses and viewing configurations (as depicted in 3). Given a street-level image, candidate aerial views are filtered and scored through a sequence of geometric constraints. Only image pairs that satisfy both spatial overlap and view-

Algorithm 2: Landmark-centric Multi-view Sampling for GeoText-1652

Input: Image folders $\{\mathcal{F}_1, \dots, \mathcal{F}_K\}$;
orbit sampling step s_o ; spiral sampling step s_s ;
maximum samples per type M
Output: Two subsets per folder: orbit views and spiral views

```
foreach folder  $\mathcal{F}$  do
  Load ordered image sequence
   $\{I_1, I_2, \dots, I_N\}$ ;
  Initialize orbit index set  $\mathcal{I}_o \leftarrow \emptyset$ ;
  Initialize spiral index set  $\mathcal{I}_s \leftarrow \emptyset$ ;
  for  $i = 1, 1 + s_o, \dots, N$  do
    Append  $i$  to  $\mathcal{I}_o$ ;
  for  $i = 1, 1 + s_s, \dots, N$  do
    Append  $i$  to  $\mathcal{I}_s$ ;
  if  $N \notin \mathcal{I}_o$  then
    Append  $N$  to  $\mathcal{I}_o$ ;
  if  $N \notin \mathcal{I}_s$  then
    Append  $N$  to  $\mathcal{I}_s$ ;
  Truncate  $\mathcal{I}_o$  and  $\mathcal{I}_s$  to at most  $M$ 
  elements;
  foreach  $i \in \mathcal{I}_o$  do
    Export image  $I_i$  to orbit subset;
  foreach  $i \in \mathcal{I}_s$  do
    Export image  $I_i$  to spiral subset;
```

1035 point consistency are retained, while ambiguous or
1036 degenerate cases are manually filtered out. This
1037 process ensures strong correspondence between
1038 ego-centric observations and their allocentric coun-
1039 terparts.

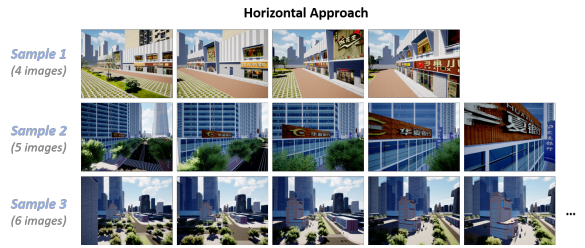
1040 Specifically, the pairing procedure relies on five
1041 core geometric factors: (1) horizontal Euclidean
1042 distance between camera positions in the ground
1043 plane; (2) height difference between aerial and
1044 street cameras; (3) overlap of projected viewing
1045 areas on the ground; (4) consistency of viewing
1046 orientation in the horizontal plane, and (5) ground-
1047 plane projection of camera viewing rays. These
1048 factors are applied sequentially to prune invalid
1049 candidates and to compute a final matching score
1050 for pair selection.

1051 We emphasize that the image data does not as-
1052 sume perfect viewpoint control in real-world data.
1053 Instead, it focuses on constructing view sets with
1054 consistent scene identity and interpretable view-

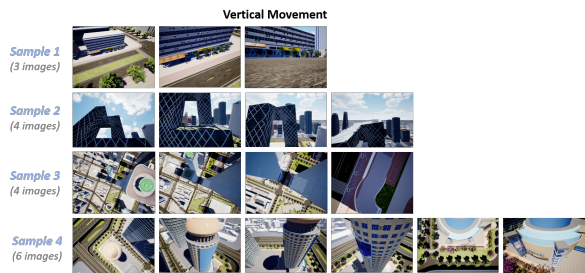
Algorithm 3: Geometry-aware Aerial-Street Pairing for MatrixCity

Input: Street images $\mathcal{S} = \{S_i\}$ with poses;
Aerial images $\mathcal{A} = \{A_j\}$ with poses;
Height range $[h_{\min}, h_{\max}]$; distance
threshold d_{\max} ;
Orientation threshold θ_{\max}
Output: Paired aerial-street image set \mathcal{P}
Initialize $\mathcal{P} \leftarrow \emptyset$;

```
foreach street image  $S_i \in \mathcal{S}$  do
  Extract street camera position  $\mathbf{p}_s$  and
  orientation  $\mathbf{o}_s$ ;
  if  $S_i$  violates diversity constraints then
    continue;
  Compute ground projection  $(\mathbf{g}_s, r_s)$ 
  from  $(\mathbf{p}_s, \mathbf{o}_s)$ ;
  Initialize best match  $A^* \leftarrow \emptyset$ , best
  score  $c^* \leftarrow \infty$ ;
  foreach aerial image  $A_j \in \mathcal{A}$  do
    Extract aerial camera position  $\mathbf{p}_a$ 
    and orientation  $\mathbf{o}_a$ ;
    // Height filtering
    if  $(\mathbf{p}_a^z - \mathbf{p}_s^z) \notin [h_{\min}, h_{\max}]$  then
      continue;
    // Horizontal distance
    filtering
    Compute  $d_h \leftarrow \|\mathbf{p}_a^{xy} - \mathbf{p}_s^{xy}\|_2$ ;
    if  $d_h > d_{\max}$  then
      continue;
    // Ground projection and
    viewing overlap
    Compute ground projection  $(\mathbf{g}_a, r_a)$ 
    from  $(\mathbf{p}_a, \mathbf{o}_a)$ ;
    Compute viewing center distance
     $d_g \leftarrow \|\mathbf{g}_a^{xy} - \mathbf{g}_s^{xy}\|_2$ ;
    if  $d_g > r_a + r_s$  then
      continue;
    // Cross-view orientation
    consistency
    Compute horizontal orientation
    angle  $\theta \leftarrow \angle(\mathbf{o}_a^{xy}, \mathbf{o}_s^{xy})$ ;
    if  $\theta > \theta_{\max}$  then
      continue;
    // Matching score
    Compute score
     $c \leftarrow \alpha d_h + \beta d_g + \gamma(\mathbf{p}_a^z - \mathbf{p}_s^z)$ ;
    if  $c < c^*$  then
       $c^* \leftarrow c, A^* \leftarrow A_j$ ;
  if  $A^* \neq \emptyset$  then
    Add pair  $(S_i, A^*)$  to  $\mathcal{P}$ ;
    Update diversity state;
```



(a) Multi-scale views in horizontal approach during dynamic translation.



(b) Multi-scale views in vertical movement during dynamic translation.

Figure 9: Examples of manually captured simulated images under dynamic translation with different motion patterns.

point variation, which is sufficient for evaluating cross-view spatial reasoning.

B.4 Details of Manually Collected Data

Cross-scale perception refers to partitioning the urban environment into multiple levels of spatial granularity, for example, progressively reducing the scale from a block-level viewpoint to that of an individual building (or facility), and further down to the object level, such as billboards, trees, or sculptures. With the assistance of experienced simulator operators in our team, we collected cross-scale image sequences under large-range urban motions, resulting in a total of 246 sequences.

As shown in Fig. 9(a), we select a visually distinctive target as an anchor point, observe it from a long distance, gradually approach it, and finally localize the target for close-range observation. During this process, the number of captured images is not fixed (typically 4–8), but is instead determined by the observability of the target.

As illustrated in Fig. 9(b), we perform takeoff or landing from an open plaza with salient landmarks. Throughout this process, we record the changes in perceived urban scale, with the number of images ranging from 3 to 6.

As shown in Fig. 10, we collected 42 sets of panoramic observations at multiple locations in

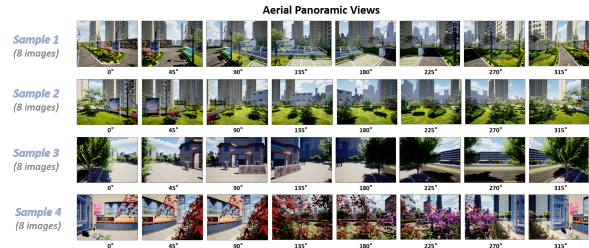


Figure 10: Examples of simulated images captured manually in rotation (panoramic views).

EmbodiedCity. Each set consists of eight camera images with a 90-degree field of view, which together form a complete 360-degree observation for the aerial agent.

B.5 Details of Multi-choice QA Generation

B.5.1 Role Templates

We adopt explicit role templates to frame the VLM as an embodied agent or an expert evaluator, which has been shown to be critical for inducing structured reasoning behaviors and consistent input–output formats in large language models.

Multi-scale Role (Fig. 11). The first role template frames the model as an observer reasoning across a hierarchy of spatial scales, ranging from regional to building and detail levels. By explicitly labeling scale transitions and enforcing scale-aware placeholders, this template encourages hierarchical spatial reasoning and object emergence analysis, which aligns with human spatial cognition theories emphasizing multi-level environmental representation.

Egocentric–Allocentric Role (Fig. 12). The second role template assigns the model the task of jointly reasoning over egocentric street-level views and allocentric aerial views of the same urban area. This design explicitly bridges self-centered perception and map-like global understanding, a distinction widely studied in cognitive science and embodied navigation. Such a role formulation enables the model to align local observations with global spatial context.

Surrounding Multi-view Role (Fig. 13). The third role template positions the VLM as an agent observing an urban object or location from multiple surrounding viewpoints, forming a 360-degree or multi-angle observation. This setting encourages cross-view consistency reasoning and spatial relation inference, which is essential for robust scene understanding under viewpoint variations.

```

MAIN INSTRUCTIONS - DYNAMIC SCALE

If mode == "horizontal movement":
    mode_desc="The camera is moving horizontally, simulating an approach towards the target."
Else:
    mode_desc = "The camera is moving vertically, simulating a descent or ascent towards the ground."

ROLE: You are given multiple images of the same urban scene captured at different scale levels. Each image will be clearly labeled with its scale identifier (e.g., Regional Scale, Building Scale, Detail Scale, etc.) immediately before the image. (mode_desc)
The images typically follow a progression from large scale (regional/area level) to medium scale (building level) to small scale (detail level), showing the same location with increasing zoom/magnification.

PLACEHOLDERS:
- <obj>: Categories or names. Named as <tree1>, <building1>, <window3>, <signboard1>, etc.
- <view>: View ID representing scale level in sequence (like large→medium→small), named as <view1>, <view2>, <view3>, etc.
- <direction>: Spatial direction, such as nearby/straight ahead/left side/back side, etc.
- <scale>: Scale hierarchy, such as "Regional level / Building level / Detail level"
- <action>: Camera movement actions, such as forward translation, backward translation, arise descent, etc.
- <option>: Options for multiple choice answers, named as <option1>, <option2>, <option3>, and <option4>

TASK REQUIREMENT: Your job is to generate ONE multiple-choice question strictly following the provided template below.

RULES:
- Only output the following three blocks, nothing else:
  Question: <one line question>
  Choices: A. ... B. ... C. ... D. ... (and E. ... if provided)
  Reasoning Process: ...
  Answer: <one of A/B/C/D/E>
- Output ONLY: Question, Choices, Reasoning Process, Answer
- Use exact viewpoint labels (view1, view2, ...)
- Use only objects visible in the images (no assumptions)
- Ensure the question is solvable by reasoning across the multiple surrounding viewpoints.
- Focus on Scale transition, Object emergence at different scales, Hierarchical relationships, Spatial detail changes, and Scale-aware reasoning.
- If this task is NOT APPLICABLE to the given images (e.g., required objects/relations are absent), OUTPUT EXACTLY: SKIP_THIS_SPEC

STRUCTURE INFORMATION:
Task Name: {task_name}
Task Instruction: {task_instruction}
View Counts: {image_count}
Few-shot Examples: {example_cases}

OUTPUT FORMAT EXAMPLE:
{
  Question: ...?
  Choices: A. ... B. ... C. ... D. ... E. ...
  Reasoning Process: ...
  Answer: B
}

```

Figure 11: Prompt template for role playing (Dynamic-scale views).

Rotation-based Viewpoint Role (Fig. 14). The fourth role template emphasizes rotational viewpoint changes around a fixed location, with each image annotated by precise viewing directions or camera poses. By enforcing strict viewpoint identifiers, this role promotes reasoning about self-position changes relative to a static environment, supporting rotation-aware spatial inference and mental viewpoint transformation.

B.5.2 Task Templates

As shown in Fig 15 to 18, we sequentially present the task templates for each camera dynamic, because these tasks are unique to this viewpoint mode.

B.6 Quality Control

B.6.1 Blind Filter

This filtering strategy aims to exclude questions that can be resolved solely through commonsense reasoning, without relying on explicit visual evidence. Specifically, we employ multiple VLMs (six open-source models in this study) to answer each question without providing any multi-view image inputs. If all VLMs correctly predict the answer

```

MAIN INSTRUCTIONS - EGO-ALLOCENTRIC

ROLE: You are given two images of the same urban area:
- Aerial View: Allocentric aerial view (high-altitude city map-like view)
- Street View: Egocentric street view (first-person street-level view)
This task focuses on understanding the relationship between individual perspective (Egocentric) and environmental perspective (Allocentric).

PLACEHOLDERS:
- <obj>: Categories or names. Named as <brown building1>, <black tower1>, etc.
- <direction>: Absolute direction, such as North, South, West, East, etc.
- <option>: Options for multiple choice answers, named as <option1>, <option2>, <option3>, and <option4>

TASK REQUIREMENT: Your job is to generate ONE multiple-choice question strictly following the provided template below.

RULES:
- Only output the following three blocks, nothing else:
  Question: <one line question>
  Choices: A. ... B. ... C. ... D. ... (and E. ... if provided)
  Reasoning Process: ...
  Answer: <one of A/B/C/D/E>
- Output ONLY: Question, Choices, Reasoning Process, Answer
- Use exact viewpoint labels (view1, view2, ...)
- Use only objects visible in the images (no assumptions)
- Ensure the question is solvable by reasoning across the multiple surrounding viewpoints.
- Focus on Appearance, Measurable indicators (height, size), Spatial relationships, Path planning.
- If this task is NOT APPLICABLE to the given images (e.g., required objects/relations are absent), OUTPUT EXACTLY: SKIP_THIS_SPEC

STRUCTURE INFORMATION:
Task Name: {task_name}
Task Instruction: {task_instruction}
Few-shot Examples: {example_cases}

OUTPUT FORMAT EXAMPLE:
{
  Question: ...?
  Choices: A. ... B. ... C. ... D. ... E. ...
  Reasoning Process: ...
  Answer: B
}

```

Figure 12: Prompt template for role playing (Ego-allocentric views).

```

MAIN INSTRUCTIONS - ORBIT

ROLE: You are given multiple images captured from different viewpoints (view1, view2, ...) surrounding the same urban object or location, forming a 360-degree or multi-angle observation. The task requires reasoning about urban scenes and spatial relationships across these viewpoints.

PLACEHOLDERS:
- <obj>: Categories or names. Named as <car1>, <car2>, <brown building1>, <black tower1>, etc.
- <view>: View ID in multi-view arrangements, named as <view1>, <view2>, <view3>, etc.
- <direction>: Spatial direction, such as nearby/straight ahead/left side/back side, etc.
- <count>: Counts of objects, such as <count1>, <count2>, <count3>, etc.
- <action>: Camera movement actions, such as forward translation, backward translation, clockwise rotation, etc.
- <option>: Options for multiple choice answers, named as <option1>, <option2>, <option3>, and <option4>

TASK REQUIREMENT: Your job is to generate ONE multiple-choice question strictly following the provided template below.

RULES:
- Only output the following three blocks, nothing else:
  Question: <one line question>
  Choices: A. ... B. ... C. ... D. ... (and E. ... if provided)
  Reasoning Process: ...
  Answer: <one of A/B/C/D/E>
- Output ONLY: Question, Choices, Reasoning Process, Answer
- Use exact viewpoint labels (view1, view2, ...)
- Use only objects visible in the images (no assumptions)
- Ensure the question is solvable by reasoning across the multiple surrounding viewpoints.
- Focus on Spatial relationships, Counting, Direction understanding, Mental rotation, and Perspective transformation.
- If this task is NOT APPLICABLE to the given images (e.g., required objects/relations are absent), OUTPUT EXACTLY: SKIP_THIS_SPEC

STRUCTURE INFORMATION:
Task Name: {task_name}
Task Instruction: {task_instruction}
View Counts: {image_count}
Few-shot Examples: {example_cases}

OUTPUT FORMAT EXAMPLE:
{
  Question: ...?
  Choices: A. ... B. ... C. ... D. ... E. ...
  Reasoning Process: ...
  Answer: B
}

```

Figure 13: Prompt template for role playing (Orbit views).

```

MAIN INSTRUCTIONS - ROTATION
ROLE: You are given multiple images of the same urban area taken from different viewpoints. Each image will be clearly labeled with its viewpoint identifier immediately before the image. {format_type}. The viewpoint labels provided with each image that indicate the exact viewing direction or camera position. This task focuses on understanding urban environments and the relationship between self-position and scene context through multi-view rotation.
PLACEHOLDERS:
- <obj>: Categories or names. Named as <brown building1>, <black tower1>, etc.
- <view>: View ID in multi-view arrangements, named as <view1>, <view2>, <view3>, etc.
- <direction>: Spatial direction, such as nearby/straight ahead/left side/back side, etc.
- <option>: Options for multiple choice answers, named as <option1>, <option2>, <option3>....
TASK REQUIREMENT: Your job is to generate ONE multiple-choice question strictly following the provided template below.
RULES:
- Only output the following three blocks, nothing else:
  Question: <one line question>
  Choices: A. ... B. ... C. ... D. ... (and E. ... if provided)
  Reasoning Process: ...
  Answer: <one of A/B/C/D/E>
- Output ONLY: Question, Choices, Reasoning Process, Answer
- Use exact viewpoint labels (view1, view2, ...)
- Use only objects visible in the images (no assumptions)
- Ensure the question is solvable by reasoning across the multiple surrounding viewpoints.
- Focus on Spatial relationships, Object counting, Direction understanding, Mental rotation, and Perspective transformation.
- If this task is NOT APPLICABLE to the given images (e.g., required objects/relations absent), OUTPUT EXACTLY: SKIP_THIS_SPEC
STRUCTURE INFORMATION:
  Task Name: {task_name}
  Task Instruction: {task_instruction}
  Few-shot Examples: {example_cases}
OUTPUT FORMAT EXAMPLE:
{
  Question: ...?
  Choices: A. ... B. ... C. ... D. ... E. ...
  Reasoning Process: ...
  Answer: B
}

```

Figure 14: Prompt template for role playing (Rotation views).

under this setting, the question is discarded, indicating that it can be solved using general perceptual priors alone and does not require detailed visual interpretation of the scene.

This procedure ensures that the retained questions genuinely demand complex image-grounded reasoning and the integration of information across multiple viewpoints, thereby sharpening the dataset’s emphasis on challenging visual perception tasks. Conversely, questions for which all VLMs fail to produce the correct answer are also removed, as such cases suggest that the question may be ill-posed, ambiguous, or outside the intended scope of the task.

Together, this bidirectional filtering mechanism ensures that the final dataset consists exclusively of questions that require authentic multi-view visual reasoning—namely, those that cannot be trivially answered by all models, yet remain solvable by some VLMs when visual input is absent.

B.6.2 Details of Human Refinement

This stage adopts a structured two-stage human refinement protocol to ensure annotation reliability and consistency. We recruit volunteers on campus and pay them reasonable compensation commensurate with the region. In detail, we recruit ten annotators with master’s or doctoral training and research experience related to urban spatial understanding. Annotators are divided into two

independent groups with complementary roles: a *refinement group*, responsible for revising automatically generated QA pairs, and a *verification group*, responsible for independent validation and consistency checking.

Each QA instance is reviewed by at least one annotator from each group, ensuring dual human coverage. Disagreements between the two groups are explicitly recorded and resolved through adjudication by the verification group, which serves as the final decision authority. This process implicitly enforces inter-annotator consistency by requiring agreement across independent reviewers before acceptance.

During refinement, annotators explicitly flag ambiguous or ill-posed questions, including cases with unclear spatial references, underspecified viewpoints, or multiple plausible answers. Such cases are either revised through rewording and answer option correction or removed entirely if ambiguity cannot be reliably resolved. This auditing process prevents semantically underspecified samples from entering the benchmark.

To further characterize refinement outcomes, we categorize common annotation issues encountered during this stage, including (i) incorrect spatial relations, (ii) misleading or visually unsupported answer options, (iii) inconsistent reasoning chains, and (iv) viewpoint-dependent ambiguities. Accepted QA pairs undergo final proofreading to eliminate residual ambiguities and formatting errors, and human-authored reasoning processes are added to improve clarity and interpretability.

The refinement interface used by annotators is illustrated in Fig. 19.

Overall, this refinement protocol prioritizes annotation reliability over scale, ensuring that subtle geometric relations are consistently grounded in the provided multi-view visual evidence.

C Experimental Details

All local model inference and fine-tuning is performed on 4×A100-SXM4-80GB. The code of our program is available at anonymous.4open.science/r/CityCube-Bench-9E72/.

C.1 Brief Introduction on Baselines

Our evaluation covers both proprietary and open-source MLLMs trained to receive multi-image inputs. The evaluated models, as well as random and

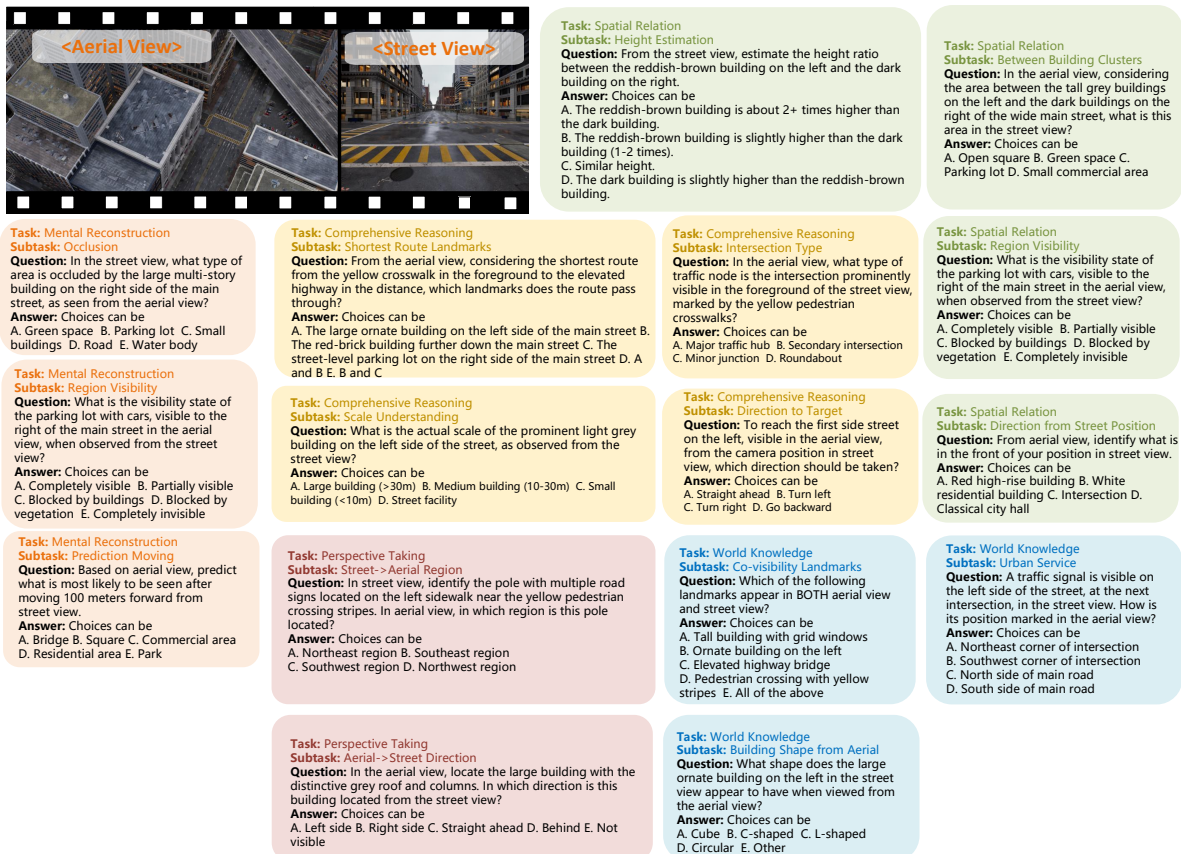


Figure 15: Examples of task templates and their corresponding images (Ego-allocentric views).



Figure 16: Examples of simulated images captured manually in rotation (Orbit views).



<p>Task: World Knowledge Subtask: Object Identification Question: Based on five views, identify what object is to the east of the large brown-roofed building in view1? Answer: Choices can be A. Car B. Truck C. Pedestrian D. Traffic sign E. Other</p>	<p>Task: World Knowledge Subtask: Urban Service Question: Based on view1, what kind of urban service is clearly identifiable in this location? Answer: Choices can be A. Parking B. Shopping C. Public transit D. Outdoor recreation E. Other</p>	<p>Task: World Knowledge Subtask: Object Counting Question: How many distinct parking areas can be found in total across all views? Answer: Choices can be A. 1 B. 3 C. 5 D. Not visible</p>	<p>Task: Spatial Relation Subtask: Direction Presence Question: Based on these images, if standing at the same position as view3, is there a grassy field with soccer goals to your right? Answer: Choices can be A. Yes B. No</p>	<p>Task: Spatial Relation Subtask: Camera Movement Question: What is the camera movement direction from view1 to view2? Answer: Choices can be A. Forward translation B. Backward translation C. Left turn D. Right turn</p>
<p>Task: Perspective Taking Subtask: Direction From Another View Question: If standing at view2, in which direction is the large yellow-roofed building from view1? Answer: Choices can be A. Left B. Right C. Straight ahead D. Straight back</p>	<p>Task: Mental Reconstruction Subtask: Multi-Object View Question: Currently facing the central brown-roofed building, which viewpoint should you rotate to in order to see both the central brown-roofed building and the black tower building simultaneously? Answer: Choices can be A. view1 B. view2 C. view3 D. view4 E. view5</p>	<p>Task: Mental Reconstruction Subtask: View Transition Direction Question: From view1 to view2, in which direction did the camera move around the central brown-roofed building? Answer: Choices can be A. Left -45 degrees B. Right -45 degrees C. Left -90 degrees D. Right -90 degrees</p>	<p>Task: Comprehensive Reasoning Subtask: Location Type Question: From view1 to view6 is depicted in these multiple viewpoints? Answer: Choices can be A. Intersection B. Mid-section of main road C. T-junction D. Roundabout</p>	<p>Task: Comprehensive Reasoning Subtask: Skyline Direction Question: From this location, in which viewpoint is the city skyline most open? Answer: Choices can be A. <view1> B. <view2> C. <view3> D. <view4></p>
<p>Task: Perspective Taking Subtask: Reverse View Question: If you stand at the position of the large brown-roofed building in the center and look toward the camera's position in view1, what would you primarily see? Answer: Choices can be A. Building B. Vehicle C. Trees/vegetation D. Open road</p>	<p>Task: Mental Reconstruction Subtask: Left Turn Prediction Question: If you are at viewpoint view1 and turn left 90 degrees, what would you see? Answer: Choices can be A. The large yellow building B. The green field C. The tall black-roofed building D. The Central Tan-Roofed Building</p>	<p>Task: Mental Reconstruction Subtask: Rotation Prediction Question: If currently facing the central brown-roofed building in view5, what would be seen directly ahead after rotating the view left 90 degrees? Answer: Choices can be A. The large yellow building B. The green field and parking lot C. The black tower building D. Uncertain</p>	<p>Task: Comprehensive Reasoning Subtask: Object Localization Sequence Question: From view1 to view6 observation, if you want to locate the 'Chessport' signboard, in what sequence do you need to lock on? Answer: Choices can be A. Region → Building → Signboard B. Region → Signboard → Building C. Building → Region → Signboard D. Signboard → Building → Region</p>	<p>Task: Comprehensive Reasoning Subtask: Movement Proximity Question: If standing at the same position as view1, facing the same direction, then turning right and walking forward, would you be closer to the small parking lot immediately south of the large brown-roofed building? Answer: Choices can be A. Yes B. No</p>
<p>Task: Perspective Taking Subtask: Direction From Object Position Question: If standing at the position of the large building with the brown roof and white/green walls and facing the tall black building, what is on its back side? Answer: Choices can be A. The large yellow building B. The large parking lot filled with cars C. The green sports field D. Not visible</p>	<p>Task: Perspective Taking Subtask: Third Person View Question: If you are a driver at the black multi-story building, can you see the scene primarily depicted in view2 from your perspective? Answer: Choices can be A. Yes, straight ahead B. Yes, to the side C. Yes, behind D. Cannot see</p>	<p>Task: Spatial Relation Subtask: Relative Position Question: Based on these images, what is the positional relationship of the yellow building relative to the brown-roofed building from view2? Answer: Choices can be A. From long shot to medium shot B. From medium shot to close-up C. From long shot directly to close-up D. Uncertain</p>	<p>Task: Spatial Relation Subtask: Cardinal Direction Question: Taking the top of view1 as North, what are the main landmarks to the North and East of the central brown-roof building? Answer: Choices can be A. North: The large yellow building, East: The large grassy field and adjacent parking lot B. North: A black tower building, East: The large grassy field and adjacent parking lot C. North: ... East: ... D. No obvious landmarks</p>	<p>Task: Spatial Relation Subtask: Behind View Question: If standing at the same position as view1 and facing the same direction, what is behind you? Answer: Choices can be A. The large yellow-roofed building B. A large parking lot filled with cars C. A green football field D. A multi-story black building and other grey-roofed buildings.</p>
<p>Task: World Knowledge Subtask: Object Description Mismatch Question: Based on these images, which option does NOT match the description of the main central building with the brown roof? Answer: Choices can be A. A large yellow building is located to the north of the main central building. B. A large green field is situated to the east of the main central building. C. ... D. ...</p>	<p>Task: Spatial Relation Subtask: Sequential Ordering Question: Looking across all views, what is the order of appearance from first to last for the large yellow building, the tall black rectangular building, and the large green athletic field? Answer: Choices can be A. Yellow building, Black building, Green field B. Yellow building, Green field, Black building C. ... D. ...</p>	<p>Task: Spatial Relation Subtask: Cardinal Direction Question: Taking the top of view1 as North, what are the main landmarks to the North and East of the central brown-roof building? Answer: Choices can be A. North: The large yellow building, East: The large grassy field and adjacent parking lot B. North: A black tower building, East: The large grassy field and adjacent parking lot C. North: ... East: ... D. No obvious landmarks</p>	<p>Task: Spatial Relation Subtask: Cardinal Direction Question: Taking the top of view1 as North, what are the main landmarks to the North and East of the central brown-roof building? Answer: Choices can be A. North: The large yellow building, East: The large grassy field and adjacent parking lot B. North: A black tower building, East: The large grassy field and adjacent parking lot C. North: ... East: ... D. No obvious landmarks</p>	<p>Task: Spatial Relation Subtask: Behind View Question: If standing at the same position as view1 and facing the same direction, what is behind you? Answer: Choices can be A. The large yellow-roofed building B. A large parking lot filled with cars C. A green football field D. A multi-story black building and other grey-roofed buildings.</p>

Figure 17: Examples of simulated images captured manually in rotation (panoramic views).



<p>Task: Spatial Relation Subtask: Left-Right Ordering Question: From viewpoint view_0, how are the Tall light-colored residential building, the Blue container structure, and the Brown two-story building ordered from left to right? Answer: Choices can be A. Tall light-colored residential building, Blue container structure, Brown two-story building B. Tall light-colored residential building, Brown two-story building, Blue container structure C. ... B. ...</p>	<p>Task: Spatial Relation Subtask: Angular Measurement Question: If standing at the same position as view_0, what is the approximate angle between the tall light-colored residential building on the left and the brown building on the right? Answer: Choices can be A. 45 degrees B. 90 degrees C. 180 degrees D. 270 degrees</p>	<p>Task: Spatial Relation Subtask: Object Direction Question: What object is to the right of the blue container observed in view_0? Answer: Choices can be A. The tall beige residential building B. The orange-red gate C. The brown commercial building D. Not visible</p>	<p>Task: World Knowledge Subtask: Viewpoint Object Detection Question: Which viewpoint can see a pink figure on the grass? Answer: Choices can be A. view_0 B. view_45 C. view_90 D. view_135</p>	<p>Task: Mental Reconstruction Subtask: Optimal Rotation Question: Which direction should one rotate to move from view_0 to view_45 most efficiently? Answer: Choices can be A. Left -45 degrees B. Right -45 degrees C. Left -90 degrees D. Right -90 degrees</p>	<p>Task: Comprehensive Reasoning Subtask: Route Planner Question: Which view can you directly move to the fence and exit the residential quarter across all views? Answer: Choices can be A. View 1 B. view2 C. View 3 D. View 6 E. View 8</p>
--	--	--	---	--	--

Figure 18: Examples of simulated images captured manually in rotation (Rotation views).

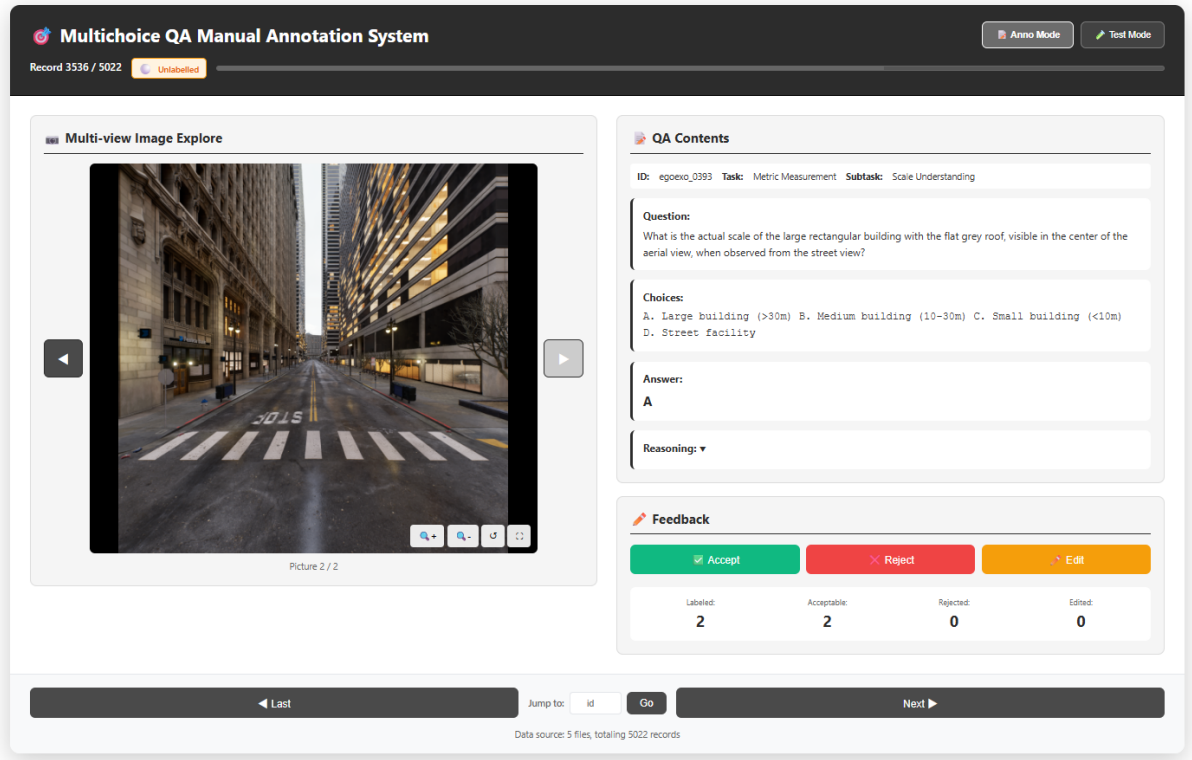


Figure 19: A demonstration Website of Multi-choice QA Annotation System on CityCube Dataset.

human baselines, are briefly introduced as follows:

Random. A random baseline model serving as the lowest performance benchmark for comparison.

Human. Human expert performance baseline representing the upper limit of human-level capability on this task.

GPT-5.1. A state-of-the-art proprietary multi-modal large language model from OpenAI (OpenAI, 2025), featuring advanced visual reasoning and long-context understanding capabilities.

Gemini-2.5-Pro. A high-performance multi-modal model developed by Google (Comanici et al., 2025), designed for complex reasoning over visual and textual inputs with strong generalization ability.

Qwen-3-VL-Plus. A large-scale vision-language model from Alibaba (Bai et al., 2025), optimized for multi-image understanding and instruction-following across diverse visual reasoning tasks.

Step-1o-turbo-vision. A reasoning-optimized vision-language model from StepFun, demonstrating competitive performance on multi-step and compositional visual reasoning benchmarks.

Doubao-seed1.6-251015. A proprietary multi-modal model from ByteDance (Guo et al., 2025), designed to support general-purpose vision-

language understanding and reasoning tasks.

Skywork-R1V4-Lite. The largest vision-language reasoning model from SkyworkAI (Wei et al., 2023), focusing on high-performance inference while maintaining strong visual comprehension capabilities.

Qwen3-VL. An open-source vision-language model family supporting multi-image inputs, widely adopted as a strong baseline for multimodal reasoning and perception tasks.

GLM-4.1V. A multimodal extension of the GLM series, featuring enhanced visual understanding and cross-modal reasoning abilities.

Kimi-VL-A3B. A compact vision-language model emphasizing efficient visual perception and instruction-following under limited parameter budgets.

MiMo-VL-7B. A 7B-parameter open-source vision-language model designed for general multimodal understanding and reasoning.

MiniCPM-V-4.5. A lightweight multimodal model from the MiniCPM series, targeting efficient deployment with competitive visual reasoning performance.

Ovis2.5. A medium-scale vision-language model from Alibaba, supporting multi-image inputs and fine-grained visual reasoning.



Figure 20: The four common failure cases for VLMs in multi-view spatial reasoning.

1374 reality, the camera motion or perspective shift was
 1375 misinterpreted, leading to a reversed understanding
 1376 of the perspective effect.

1377 E Implementation Details

1378 The dataset is split into training and test sets with a
 1379 ratio of 9:1, and data loading is parallelized using
 1380 four worker processes. The fine-tune experiment
 1381 results are produced on test set.

1382 E.1 Training Details

1383 We fine-tune **Qwen3-VL-2B**, **Qwen3-VL-4B** and
 1384 **Qwen3-VL-8B** using supervised fine-tuning with
 1385 Low-Rank Adaptation (LoRA). As described in
 1386 Tab 3, the models are trained for 5 epochs using
 1387 the Adam optimizer with an initial learning rate of
 1388 1×10^{-4} and a warmup ratio of 0.05. To reduce
 1389 memory consumption while maintaining training
 1390 stability, we employ bfloat16 precision, gradient ac-
 1391 cumulation with 4 steps, and FlashAttention. The
 1392 maximum input sequence length is set to 4096 to-
 1393 kens, and images are resized to ensure the total
 1394 pixel count does not exceed 1.6M.

1395 LoRA is applied to all linear layers with a rank
 1396 of 8 and a scaling factor of 32. We do not freeze the

Table 3: Training hyperparameters for **CityBot** series.

Category	Hyperparameter	Value
Training Setup	Training paradigm	Supervised fine-tuning (LoRA)
	Number of epochs	5
	Precision	bfloat16
Batching	Per-device train batch size	1
	Per-device eval batch size	1
	Gradient accumulation steps	4
	Effective batch size	4
Optimization	Learning rate	1×10^{-4}
	Warmup ratio	0.05
LoRA Configuration	LoRA rank (r)	8
	LoRA scaling factor (α)	32
	Target modules	All linear layers
	Training scope	Full model (no freezing)
Architecture & Memory	Attention implementation	FlashAttention
	Padding-free training	Enabled
	Sample packing	Enabled
	Gradient checkpointing	Enabled (ViT excluded)
Input	Maximum sequence length	4096
	Maximum image pixels	1,605,632
Evaluation & Logging	Evaluation interval	Every 100 steps
	Checkpoint saving interval	Every 100 steps
	Max checkpoints kept	3
	Logging interval	Every 10 steps
Data Loading	Test split ratio	0.1
	Dataset preprocessing workers	4
	Dataloader workers	4

1397 vision encoder or the vision-language alignment
 1398 module to allow full end-to-end adaptation. We
 1399 further enable padding-free training and sample
 1400 packing to improve computational efficiency.

1401 During training, evaluation and checkpoint sav-
 1402 ing are performed every 100 steps, and at most
 1403 three checkpoints are retained.

E.2 Human Evaluation Details

As shown in Figure 21, we collected a set of volunteer evaluation results through an interactive web interface. These volunteers are recruited independently and unrelated to the annotators who participated in refining the questions. They also did not receive any extra training in spatial knowledge.

F Further Discussion

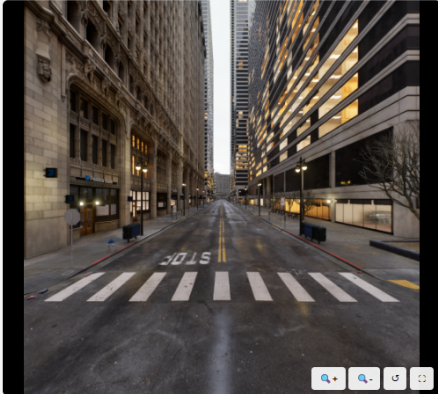
3D Question Answering A precise understanding of objects and their spatial relationships within 3D scenes is critical for applications in robotics and autonomous driving. Leveraging the holistic 3D scene datasets such as ScanNet (Dai et al., 2017) and Matterport3D (Chang et al., 2017), the research community has developed a variety of 3D Question Answering (3DQA) benchmarks like ScanQA (Azuma et al., 2022; Achlioptas et al., 2020; Hong et al., 2023). The core paradigm of these benchmarks involves inferring a comprehensive 3D understanding from a limited set of 2D views. In line with this, employing multiple views as input for the LLMs has become a primary methodology in 3DQA (Fu et al., 2024; Huang et al., 2024; Guo et al., 2023). Our work extends this research trajectory and systematically incorporates typical view combinations found in existing literature and, more importantly, pioneers the extension of the 3DQA application domain from indoor environments to the more expansive and challenging context of urban spaces.

Multi-choice QA Manual Annotation System

Anno Mode
Test Mode

Record 3536 / 5022
Unlabeled

Multi-view Image Explore



Picture 2 / 2

QA Contents

ID: egoexo_0393 Task: Metric Measurement Subtask: Scale Understanding

Question:
What is the actual scale of the large rectangular building with the flat grey roof, visible in the center of the aerial view, when observed from the street view?

Choices:
A. Large building (>30m) B. Medium building (10-30m) C. Small building (<10m)
D. Street facility

Please Select An Answer:

A Large building (>30m)

B Medium building (10-30m)

C Small building (<10m)

D Street facility

Tip: select the option will be automatically submitted, you can also press the number key 1-6 to quickly select A-F

Statistics

Total 5022	Answered 5022	Corrected 2590	Wrong 2432	Accuracy 52%
----------------------	-------------------------	--------------------------	----------------------	------------------------

[Export Results](#)

◀ Last
Jump to: [Go](#)
Next ▶

Data source: 5 files, totaling 5022 records

Figure 21: A demonstration Website of Multi-choice QA Test System of CityCube Dataset.