# INTERNSPATIAL: A COMPREHENSIVE DATASET FOR SPATIAL REASONING IN VISION-LANGUAGE MODELS

#### **Anonymous authors**

Paper under double-blind review

#### **ABSTRACT**

Recent benchmarks and datasets have been proposed to improve spatial reasoning in vision-language models (VLMs), yet existing open resources remain limited in scale, visual diversity, and instruction expressiveness. In this work, we introduce InternSpatial, the largest open-source dataset for spatial reasoning in VLMs, along with InternSpatial-Bench, a corresponding evaluation benchmark designed to assess spatial understanding under diverse instruction formats. InternSpatial comprises 12 million QA pairs spanning both single-view and multi-view settings, drawn from diverse visual environments and supporting 19 instruction formats that reflect varied query styles. For evaluation, we propose InternSpatial-Bench for single-view tasks and expand multi-view reasoning by introducing a novel rotation angle prediction task that has not been explored in prior work. Experimental results show that models trained on InternSpatial achieve 12.1% improvement on InternSpatial-Bench and 10.7% on VSI-Bench, while maintaining strong performance on general-purpose benchmarks. We hope these resources will support the development of spatially capable VLMs in practical applications such as robotics and embodied AI.

### 1 Introduction

Vision-language models (VLMs) have achieved remarkable progress across a range of multimodal tasks such as visual question-answering (VQA), image captioning, and grounding, demonstrating their ability to align and reason over visual and textual inputs. Nonetheless, they still struggle with spatial reasoning, both in single-view settings (*e.g.*, identifying object position or size from a static image) and in multi-view scenarios (*e.g.*, estimating distances or tracking appearance order across dynamic video frames). Enhancing spatial reasoning capabilities in VLMs is crucial for real-world applications, including robotics, autonomous navigation, and augmented reality, where accurate spatial understanding is essential for interaction with complex environments.

Recent efforts have introduced spatially-relevant VQA datasets and corresponding evaluation benchmarks to enhance and assess VLMs' spatial reasoning capabilities (Cai et al., 2025; Cheng et al., 2024; Chen et al., 2024a; Yang et al., 2024). While these works have advanced the field, they still exhibit several notable limitations. (1) *Limited scene diversity*: existing datasets are typically drawn from narrow sources, primarily indoor or outdoor scenes, and fail to capture a broader spectrum of scenarios. (2) *Restricted instruction formats*: SpatialVLM (Chen et al., 2024a) and SpatialQA (Cai et al., 2025) rely exclusively on natural language, and OSD (Cheng et al., 2024) uses region masks. These limited formats fail to reflect the diversity of instruction types required for practical spatial reasoning tasks. (3) *Narrow training scope*: existing spatial training data primarily focus on single-view settings and cover only basic spatial concepts from a single static image, such as object position or existence, without providing multi-view supervision that captures spatial relationships across different viewpoints or temporal sequences. These limitations underscore the need for a more comprehensive dataset paired with a corresponding evaluation benchmark to advance spatial reasoning in VLMs.

To address these limitations, we propose the largest open-source spatial reasoning dataset, *InternSpatial*, and a corresponding evaluation benchmark, *InternSpatial-Bench*, specifically designed to enhance spatial reasoning capabilities in VLMs. InternSpatial comprises 9.5M single-view and 2.5M multi-view question-answer pairs, sourced from a broad spectrum of visual environments, in-

Table 1: Comparison of our InternSpatial with existing spatial reasoning datasets. W: in-the-wild, I: indoor, D: drive, E: embodied, O: object-centric

Dataset	# of QA	Scenario	Open-source	View Type	Instruction format
SpatialVLM (Chen et al., 2024a)	2B	W	X	Single-view	Single-format
SpatialQA (Cai et al., 2025)	0.9M	W,E	✓	Single-view	Single-format
OSD (Cheng et al., 2024)	8.7M	W	✓	Single-view	Single-format
InternSpatial	12M	W,I,D,E,O	✓	Single-view, Multi-view	Multiple-format

cluding in-the-wild scenes (Lin et al., 2014; Wang et al., 2024c; Krishna et al., 2017), structured indoor spaces (Wald et al., 2019; Dai et al., 2017; Mao et al., 2022), urban streetscapes (Cordts et al., 2016), object-centric scenes (Deitke et al., 2022), and embodied navigation contexts (Anderson et al., 2018). To enrich instruction formats, we incorporate a diverse set of query representations, including masks, bounding boxes, and numerical indicators embedded in images, as well as coordinate-based references and spatial cues expressed through textual instructions. In total, our dataset supports 19 distinct instruction formats, enabling broader coverage of spatial reasoning query types. We further introduce a novel multi-view task, rotation angle prediction, with 2.46M newly collected training question-answer pairs, which has not been addressed in prior spatial reasoning benchmarks. To facilitate evaluation, we construct InternSpatial-Bench with 6,008 question-answer pairs, serving as a comprehensive diagnostic benchmark for single-view spatial reasoning tasks. For multi-view evaluation, we extend the existing VSI benchmark by adding 1,000 additional question-answer pairs for the rotation angle prediction task. As shown in Table 1, our InternSpatial significantly expands scene coverage, instruction format diversity, and multi-view supervision compared to existing benchmarks.

In summary, our contributions are threefold:

- (1) We present InternSpatial, the largest open-source spatial reasoning dataset for VLMs, designed for supervised fine-tuning. It contains single-view and multi-view samples across diverse scenes and supports 19 instruction formats to support varied spatial query forms.
- (2) To support evaluation, we introduce InternSpatial-Bench for single-view tasks and extend the VSI benchmark for multi-view evaluation, incorporating a novel rotation angle prediction task not addressed in existing datasets.
- (3) Extensive experimental results demonstrate the effectiveness of InternSpatial, showing that it substantially improves spatial reasoning in VLMs, achieving a 12.1% improvement on InternSpatial-Bench and 10.7% on VSI-Bench while preserving general multimodal performance.

#### 2 Related Work

### 2.1 SPATIAL REASONING VIA VISION LANGUAGE MODELS

Recently, numerous large language models (LLMs) (Brown et al., 2020; Achiam et al., 2023; Touvron et al., 2023) and vision-language models (VLMs) (Zhu et al., 2022; Li et al., 2023a; Zhu et al., 2023a; Wang et al., 2023; Liu et al., 2023; Li et al., 2023b; Wang et al., 2024c; Chen et al., 2024c). have been developed. However, growing evidence indicates that VLMs still struggle with spatial reasoning tasks. (Cai et al., 2025; Chen et al., 2024a; Cheng et al., 2024; Yang et al., 2024) To address this limitation, several approaches have attempted to enhance spatial reasoning capabilities by incorporating additional information. For example, 3D-LLM (Hong et al., 2023b) and 3D-CLR (Hong et al., 2023a) introduce 3D representations and dense features; SpatialRGPT (Cheng et al., 2024) incorporates mask-based supervision; and SpatialBot (Cai et al., 2025) leverages depth information. Despite these efforts, current methods have not succeeded in enabling VLMs to perform end-to-end spatial reasoning effectively.

#### 2.2 Spatial Reasoning Datasets

To evaluate and improve the spatial reasoning capabilities of VLMs, several datasets and benchmarks have been proposed to cover a range of tasks and scenarios. One such benchmark, Spatial-Eval (Wang et al., 2024a), targets 2D spatial reasoning across tasks such as relation understanding,

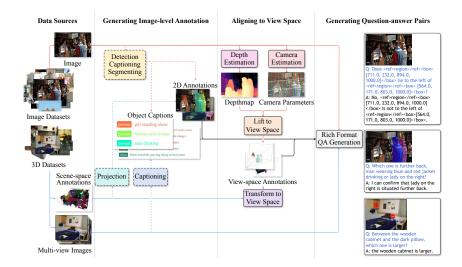


Figure 1: Generation pipeline for InternSpatial. The optional flows (represented by dashed lines and boxes) are only performed when the relevant annotations does not exist in the data source.

navigation, and counting. Another line of work explores spatial reasoning from a top-down perspective, emphasizing the need to enhance VLM performance in top-view settings (Li et al., 2024c). To enable VLMs to understand 3D spatial relationships from images, several datasets have been introduced that focus on answering 3D spatial reasoning questions (Cheng et al., 2024; Cai et al., 2025; Li et al., 2024d). However, these datasets are primarily tailored to specific models and often rely on additional inputs, such as segmentation masks or depth maps. An automatic data generation framework has also been developed to construct a large-scale 3D spatial VQA dataset using Internet images (Chen et al., 2024a), demonstrating that with appropriate training data, VLMs can infer spatial relationships without relying on auxiliary inputs. Nevertheless, the dataset is not publicly available. Spatial reasoning over image sequences or videos presents additional challenges. To assess such capabilities, the VSI benchmark (Yang et al., 2024) was proposed, evaluating a range of open-source and proprietary VLMs. Results show that current models still struggle with multi-frame spatial reasoning tasks. Our work addresses these limitations by introducing a dataset that integrates both single-view and multi-view tasks, significantly enhancing the spatial reasoning ability of VLMs across diverse contexts and highlighting their potential for deeper spatial understanding.

# 3 Dataset

#### 3.1 Data Engine for InternSpatial

We construct InternSpatial, a large-scale dataset comprising nearly 12 million Question-Answer(QA) pairs, to enable VLMs to perform 3D spatial reasoning through supervised fine-tuning. InternSpatial aggregates data from a wide range of sources, including in-the-wild scenes (Lin et al., 2014; Wang et al., 2024c; Krishna et al., 2017), structured indoor spaces (Wald et al., 2019; Dai et al., 2017; Mao et al., 2022), urban streetscapes (Cordts et al., 2016), object-centric scenes (Deitke et al., 2022), and embodied navigation contexts (Anderson et al., 2018).

To handle the heterogeneity of source data and support large-scale QA generation, we develop a fully automated and modular data engine that consolidates intermediate annotation extraction and QA synthesis into a unified pipeline applicable across diverse data sources. As illustrated in Figure 1, the pipeline begins by generating necessary annotations at the image level, followed by transforming the annotations into a canonical view space. Finally, QA pairs are constructed using a template-based approach that supports a wide variety of task types and instruction formats.

**Generating Image-level Annotation.** To generate 3D spatial reasoning QAs grounded in objects, we first obtain the necessary image-level annotations, including 2D bounding boxes, region descriptions, segmentation masks, etc. For image datasets that already provide such annotations,

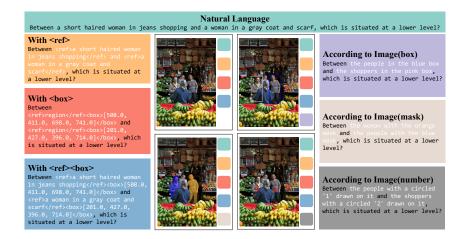


Figure 2: Examples of diverse instruction formats in text and image. The four images illustrate different visual formats: original (top-left), bounding boxes (top-right), segmentation masks (bottom-left), and numbered regions (bottom-right). Surrounding the images are seven corresponding text instruction formats. The color blocks beside each image indicate whether the corresponding image-text pair is included in InternSpatial and InternSpatial-Bench. Best viewed in color.

we directly utilize the existing labels. When annotations are missing, we employ pretrained models to generate them automatically. Specifically, we use open-source VLMs to extract object-level 2D boxes and associated textual descriptions, and apply the SAM2 model (Ravi et al., 2024) to generate segmentation masks within these boxes. These masks are subsequently lifted into 3D space to facilitate the construction of 3D bounding boxes. The prompts we used in this step can be found in Appendix C. In the case of 3D datasets, which typically include global 3D annotations and per-view camera parameters, we project the 3D information onto the image plane to obtain the corresponding 2D annotations. Although this projection is not strictly required for generating QAs, as the underlying 3D annotations are already available, it is necessary for supporting visual reference forms in prompts, such as bounding boxes and segmentation masks.

Aligning to View Space. To determine spatial relationships between objects, it is essential to obtain their positions and dimensions within a well-defined 3D coordinate system. We adopt a canonical view space as the reference frame, defined as a 3D Cartesian coordinate system centered at the camera's optical center. In this space, the y-axis aligns with the viewing direction, and the z-axis is perpendicular to the scene's horizontal plane, pointing upward. For 3D datasets, which provide global annotations and per-view camera parameters, transforming annotations into the canonical view space is straightforward. In contrast, image-only datasets contain only 2D visual information, requiring estimation of both camera parameters and depth maps. To address this, we follow the pipeline of SpatialRGPT (Cheng et al., 2024), leveraging WildCamera (Zhu et al., 2023b) for intrinsic parameter estimation, PerspectiveFields (Jin et al., 2023) for extrinsic parameter inference, and Metric3Dv2 (Hu et al., 2024) to predict dense depth maps. By combining the outputs of these models, we lift 2D annotations into the canonical 3D space, enabling accurate reasoning over object-level spatial relationships.

**Template-based QA Generation.** While prompting a large language model (LLM) to generate QA pairs directly for each image can produce diverse instructions, this approach is prohibitively expensive at scale in terms of computation and time. Instead, we adopt a template-based generation strategy that avoids invoking the LLM during QA construction. This approach not only improves efficiency but also facilitates flexible expansion to multiple prompt styles, such as object references via bounding boxes or segmentation masks. To ensure sufficient instruction diversity, we first prompt an LLM to generate several question-answer templates for each task type and answer format. These templates contain placeholders for object references and other variable content. During generation, we randomly select a subset of tasks and object instances (or pairs) for each image, derive the corresponding answers using the previously constructed annotations, and instantiate the templates accordingly. We then filter out low-quality QA pairs, such as those involving ambiguous spatial rela-

tionships caused by occlusion, and balance the number of positive and negative examples to produce a well-structured dataset. We generate templates for 4 single-view tasks, covering the position/size relationship of two objects, as well as relationship-constrained count and existence tasks. The list of templates are shown in Appendix B.

**Extending Instruction Formats.** To enhance dataset diversity and better reflect real-world usage scenarios, we extend each QA pair into multiple instruction formats. Specifically, we generate up to five textual formats and up to four image formats per QA pair. The image formats include: (1) the original image, (2) the image annotated with bounding boxes, (3) the image with segmentation masks, and (4) the image annotated with numbers over key objects. The textual formats include: (1) natural language descriptions, (2) text with <ref>{caption}</ref> (3) text with <ref>region</ref> (box) {bbox}</box> (4) text with <ref>{caption}</ref> (5) text automatically generated based on image content. Representative examples of these visual and textual formats are shown in Figure 2. As a result, each QA pair can produce up to 19 training samples, from which only suitable ones are retained. Additionally, certain prompt types, such as images with numbers on key objects, may not directly indicate the correct object. Therefore, in these cases, we utilize the position information from the segmentation mask to correctly identify and reference the target object.

**Generating Multi-view QA Pairs.** To develop a comprehensive multi-view dataset for spatial understanding, we systematically collected and integrated multi-view data derived from the training splits of the ScanNet (Dai et al., 2017), MultiScan (Mao et al., 2022), R2R (Anderson et al., 2018), and Objaverse (Deitke et al., 2022), subsequently formulating temporally-agnostic training samples that encapsulate inter-object relational attributes such as relative properties, scale variations, and spatial distances, and cross-view relationships of objects such as rotation. Scene-level geometric priors were established by estimating room dimensions via the Alpha Shape algorithm (Akkiraju et al., 1995) applied to the point clouds, with the room centroid defined as the geometric center of the minimal axis-aligned bounding box enclosing the scene. We meticulously cataloged instance counts for each object semantic category. For unambiguous objects within the point clouds exhibiting a principal dimension exceeding 15cm, annotations were standardized to the OrientedBoundingBox format using Open3D (Zhou et al., 2018). For remaining objects or those with initial ambiguities, we leveraged existing annotations to reduce the risk of shortcut learning by language models. Plausible alternative options were constructed by extracting distractors from other items within the dataset, thereby forming a corresponding multiple-choice question training set.

# 3.2 INTERNSPATIAL-BENCH

To evaluate the performance of VLMs on 3D spatial reasoning tasks, particularly under diverse instruction formats, we propose InternSpatial-Bench, a novel multi-task benchmark that features a broad range of input types. Existing benchmarks such as SpatialRGPT-Bench (Cheng et al., 2024) and SpatialBench (Cai et al., 2025) present several limitations. First, the question formats are overly simplistic and do not reflect real-world application scenarios. Second, these benchmarks are tailored to specialized models and require auxiliary inputs such as region masks or depth maps. As a result, many tasks are incompatible with general-purpose VLMs that operate solely on images and text. Furthermore, SpatialBench suffers from a limited number of QA pairs, reducing its effectiveness as a comprehensive evaluation suite.

InternSpatial-Bench expands and refines both SpatialRGPT-Bench and SpatialBench to overcome these limitations. Specifically, we enrich instruction formats and introduce 3,000 carefully curated QA pairs, resulting in a total of 5,300 high-quality examples that span diverse task types and input modalities. Certain tasks from the original benchmarks, such as reachability prediction and quantitative estimation of spatial extent, are excluded because they are unsuitable for general-purpose VLMs when only a single-view image is provided. In the absence of additional information, such as depth or camera parameters, these tasks become severely under-constrained and often ambiguous, even for human annotators.

**Refining and Expanding SpatialRGPT-Bench and SpatialBench.** Since SpatialRGPT-Bench (Cheng et al., 2024) already provides a sufficient number of QA pairs, our focus is on expanding the diversity of question formats rather than increasing the dataset size. Specifically, we augment

the instruction styles of the original questions that do not involve numerical reasoning, following the format extension strategy described in subsection 3.1. However, to avoid ambiguity caused by duplicate object labels, we exclude formats that rely on natural language references or textual content containing <ref>caption</ref>. For each selected question, we randomly sample three different formats and leverage both object mask and bounding box annotations to construct the final benchmark entries. SpatialBench (Cai et al., 2025) contains QA pairs exclusively in natural language form. To diversify its instruction formats, we first manually extract reference phrases corresponding to the mentioned objects and convert the questions into templates with placeholders. Next, we prompt the VLM to ground the objects based on these phrases and apply SAM2 to segment the corresponding regions. Using the resulting question templates, along with object bounding boxes and masks, we apply the format extension method described in subsection 3.1 to generate diverse instruction variants for each QA. Finally, all generated QA pairs are manually verified to ensure quality, with erroneous answers corrected and ambiguous or ill-formed questions removed

**Extending the Benchmark with Curated QA Pairs.** Unlike the large-scale training dataset, the benchmark is relatively small but demands higher annotation quality. To this end, we implement a dedicated pipeline for generating high-quality QAs used in the benchmark. This pipeline operates without relying on any pre-annotated information, making it applicable to any image-only data source. To encourage diversity and expressiveness in question formulation, we prompt the VLM to generate questions directly. Finally, we introduce a manual verification step to review all automatically constructed questions and answers, ensuring the overall quality and correctness of the benchmark data. Details of the construction process are provided in Appendix D.

#### 3.3 Dataset Statistics

 Statistics of InternSpatial. Our proposed dataset, InternSpatial, encompasses a diverse set of tasks and instruction formats to comprehensively enhance spatial reasoning capabilities. It consists of a total of 12,035,415 question-answer pairs, covering both single-view and multi-view spatial reasoning tasks. Specifically, the single-view tasks include *Position Comparison*, *Size Comparison*, Existence Estimation, and Object Counting, while the multi-view tasks include Rotation Estimation, Object Counting, Room Size Estimation, Object Size Estimation, Route Planning, and Appearance Order. Detailed task descriptions and corresponding statistics are provided in Appendix A, and visual examples are shown in Appendix F. In addition, InternSpatial incorporates images from various sources to enhance the robustness of the model. As illustrated in Figure 3, the dataset includes COCO (Lin et al., 2014), AS-1B (Wang et al., 2024c), and Visual Genome (VG) (Krishna et al., 2017) for in-the-wild imagery; 3RScan (Wald et al., 2019), ScanNet (Dai et al., 2017), and MultiScan (Mao et al., 2022) for indoor scenes; Cityscapes (Cordts et al., 2016) for street scenes; Objaverse (Deitke et al., 2022) for single-object scenarios; and R2R (Anderson et al., 2018) for embodied navigation tasks. Moreover, InternSpatial emphasizes diversity in instruction formats. As shown in Figure 3, the number of samples across different formats is carefully balanced to avoid bias and ensure uniform coverage during training. In summary, InternSpatial provides a large-scale, diverse resource spanning task types, visual domains, and instruction formats, making it well-suited for training VLMs to handle real-world spatial reasoning tasks effectively.

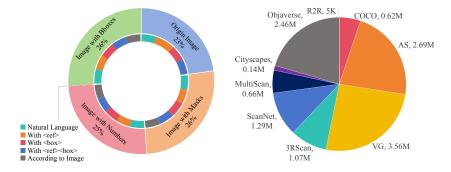


Figure 3: Distribution of instruction formats (**Left**) and data sources (**Right**) in InternSpatial.

Statistics of InternSpatial-Bench Following Spatial-Bench and Spatial-RGPT, our proposed benchmark, InternSpatial-Bench, includes five tasks—Position Estimation, Size Estimation, Rotation Estimation, Existence Estimation, and Object Counting—designed to systematically evaluate the spatial reasoning capabilities of VLMs. In total, InternSpatial-Bench consists of 6,008 QA pairs. Detailed task statistics are provided in Appendix A, and visual examples are shown in Appendix G. To ensure robustness and diversity, InternSpatial-Bench incorporates images from a broad range of domains. As shown in Fig 4, in addition to the sources used in Spatial-Bench and Spatial-RGPT, we include samples from the test sets of COCO, Flickr30K, Objaverse, ScanNet, and Cityscapes. This diverse image collection spans a wide range of real-world contexts, from indoor and outdoor environments to single-object scenarios and in-the-wild imagery. We apply the same instruction format expansion strategy as used in InternSpatial, with one exception: for the Rotation Estimation task, since each image contains only a single object, we only use the original image format and natural language instructions. Consequently, these formats have a higher proportion in this task compared to others. By combining diversity in task types, visual domains, and instruction formats, InternSpatial-Bench offers a comprehensive and realistic benchmark for evaluating the spatial reasoning abilities of VLMs across a wide range of practical scenarios.

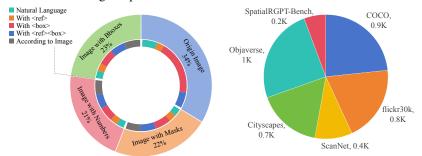


Figure 4: Distribution of instruction formats (Left) and data sources (Right) in InternSpatial-Bench.

### 4 EXPERIMENTS

We begin in Section 4.1 by introducing the baseline model and outlining the evaluation benchmarks used in our experiments. Section 4.2 then presents results on InternSpatial-Bench to assess the spatial reasoning capabilities of vision-language models. Section 4.3 reports performance on VSI-Bench (Yang et al., 2024), which further evaluates the models' multi-view spatial reasoning abilities. In Section 4.4, we conduct an ablation study to analyze the impact of different instruction formats on model performance. Finally, Section 4.5 evaluates whether training with InternSpatial affects general reasoning ability by benchmarking against a suite of standard vision-language tasks.

#### 4.1 EXPERIMENT SETUP

**Baseline.** We construct our baselines based on InternVL2.5-8B (Chen et al., 2024c), a representative traditional VLM. Following the training settings of InternVL2.5, we fine-tune our models from InternVL2.5-8B using a downsampled version of the general datasets employed in InternVL2.5, along with InternSpatial. Unless otherwise specified, we refer to the model trained on InternSpatial-Bench as **InternVL-Spatial-8B**. Detailed training configurations are provided in Appendix E.

**Evaluation.** We evaluate the models trained on InternSpatial using three types of benchmarks: our proposed InternSpatial-Bench, the multi-view spatial reasoning benchmark VSI-Bench (Yang et al., 2024), and several general-purpose benchmarks, including MathVision (Wang et al., 2024b), OCRBench (Liu et al., 2024), TextVQA (Singh et al., 2019), ChartQA (Masry et al., 2022), and MM-Star (Chen et al., 2024b). For InternSpatial-Bench, we follow the evaluation protocols of Spatial-Bench (Cai et al., 2025) and Spatial-RGPT (Cheng et al., 2024), reporting relative error for counting tasks, accuracy for multiple-choice questions, and GPT-4o-assigned (OpenAI, 2025) scores for quizstyle questions. For VSI-Bench, we adopt the official evaluation protocol, with the only modification being the use of 32 sampled frames per video during testing. For general benchmarks, we follow the evaluation procedures provided by OpenCompass (Contributors, 2023).

#### 4.2 EVALUATION ON INTERNSPATIAL-BENCH

Table 2: Results on InternSpatial-Bench. **Bold** indicates the best performance among all models, while <u>underline</u> denotes the second-best performance.

Model	Position Comparison	Size Comparison	Rotation Estimation	Object Counting	Existence Estimation	Average
GPT-4o-2024-11-20 (OpenAI, 2025)	71.2	71.5	26.7	63.5	74.9	61.6
Claude-3.7-Sonnet-20250219 (Anthropic, 2024)	73.2	72.3	25.9	59.2	70.5	60.2
Gemini-2.5-Flash(Comanici et al., 2025)	64.5	67.3	30.2	67.0	67.3	59.3
Llama-4-Scout(Meta Platforms, 2025)	42.2	45.0	20.8	44.0	25.7	35.5
Qwen2.5-VL-72B (Bai et al., 2025)	54.6	55.3	<u>30.6</u>	60.5	63.3	52.9
Pixtral-12B (Agrawal et al., 2024)	65.6	62.9	5.8	52.5	78.3	53.0
LLaVA-OneVision-72B(Li et al., 2024b)	<u>77.8</u>	<u>77.0</u>	25.8	64.5	77.6	<u>64.5</u>
InternVL2.5-8B (Chen et al., 2024c)	62.8	57.7	28.5	<u>67.8</u>	77.9	58.9
InternVL-Spatial-8B	87.8(+25.0)	78.6(+20.9)	33.6(+5.1)	71.3(+3.5)	83.9(+6.0)	71.0(+12.1)

To evaluate model performance in spatial reasoning, we conducted experiments on InternSpatial-Bench. The accuracy computation follows the methodology of Spatial-Bench (Cai et al., 2025) and Spatial-RPGT (Cheng et al., 2024), with a modification for the Object Counting task: since some VLMs struggle to follow instructions precisely, we extract the last number mentioned in the response as the predicted count and compute the relative error accordingly.

As shown in Table 2, our model, InternVL-Spatial-8B, outperforms the baseline InternVL2.5-8B (Chen et al., 2024c) by 12% in average accuracy. Notably, it achieves a 25% improvement in the Position Comparison task and a 20.9% gain in the Size Comparison task. Furthermore, InternVL-Spatial-8B surpasses advanced proprietary models such as GPT-4o (OpenAI, 2025) and Claude 3.5 Sonnet (Anthropic, 2024) across all tasks, demonstrating the effectiveness of InternSpatial in enhancing the spatial reasoning capabilities of VLMs.

#### 4.3 EVALUATION ON VSI-BENCH

Table 3: Results on VSI-Bench. **Bold** indicates the best performance among all models, while underline denotes the second-best performance.

Model	Obj.Count	Abs.Dist.	Obj.size	Room Size	Rel.Dist.	Route Plan	Appr.Order	Average
GPT-40 (OpenAI, 2025)	46.2	5.3	43.8	38.2	37.0	31.5	28.5	32.9
Gemini-1.5 Flash (Reid et al., 2024)	49.8	30.8	53.5	54.4	37.7	31.5	37.8	42.3
Gemini-1.5 Pro (Reid et al., 2024)	56.2	30.9	64.1	43.6	51.3	36.0	34.6	45.3
VILA-1.5-40B (Lin et al., 2024)	22.4	24.8	48.7	22.7	40.5	31.5	32.9	32.0
LLaVA-NeXT-Video-72B (Zhang et al., 2024)	48.9	22.8	57.4	35.3	42.4	<u>35.0</u>	48.6	41.5
LLaVA-OneVision-72B (Li et al., 2024a)	43.5	23.9	57.6	37.5	42.5	32.5	44.6	40.2
InternVL2.5-8B (Chen et al., 2024c)	51.7	32.9	45.1	42.3	40.8	27.8	50.5	41.6
InternVL-Spatial-8B	<b>68.7</b> (+17.0)	40.9(+8.0)	63.1(+18.0)	54.3(+12.0)	47.7(+6.9)	29.9(+2.1)	60.5(+10.0)	52.3(+10.7)

To evaluate the additional multi-view spatial reasoning capabilities of InternVL-Spatial-8B trained on InternSpatial, we conducted experiments on VSI-Bench (Yang et al., 2024). As shown in Table 3, InternVL-Spatial-8B achieves notable improvements over the baseline InternVL2.5-8B (Chen et al., 2024c) across all tasks in the benchmark. In particular, it surpasses the baseline by more than 10% in Object Counting, Object Size Estimation, and Appearance Order tasks.

When compared against both open-source and proprietary models, InternVL-Spatial-8B delivers top-tier performance: it ranks first in Object Counting, Absolute Distance Estimation, Object Size Estimation and Appearance Order, and second in the remaining tasks. Overall, it achieves the highest average score among all evaluated models, including GPT-40 (OpenAI, 2025) and Gemini-1.5 Pro (Reid et al., 2024). These results demonstrate that InternSpatial substantially enhances the spatial reasoning capabilities of vision-language models in multi-image scenarios.

#### 4.4 EFFECT OF THE VARIOUS QUESTION FORMATS

We conduct an ablation study on InternSpatial-Bench to evaluate the impact of different instruction formats in both the training step and evaluation step. Since the Rotation Estimation task does not include instruction format expansion, we exclude it from this analysis. Additionally, we train a

variant of InternVL2.5-8B using InternSpatial-Bench without instruction format expansion, referred to as **InternVL-Spatial-Raw-8B**.

As shown in Figure 5, the baseline model, InternVL2.5-8B (Chen et al., 2024c), performs best on original images and natural language instructions, which are prevalent in general-purpose training datasets. However, it performs significantly worse on formats involving elements such as <code>/box¿</code>, which are rare in typical datasets. In contrast, InternVL-Spatial-8B, trained on InternSpatial with diverse instruction format expansions, substantially narrows this performance gap across different instruction styles. Furthermore, comparing InternVL2.5-8B with InternVL-Spatial-Raw-8B reveals that even without instruction format expansion, InternVL-Spatial-Raw-8B consistently outperforms the baseline across all instruction styles. This indicates that the model gains a degree of generalization and cross-format transfer ability, even without being explicitly trained on diverse instruction forms. Finally, InternVL-Spatial-8B achieves the best performance across all instruction formats, including natural language and original image styles. This demonstrates that instruction format expansion not only improves the model's robustness to diverse input styles but also enhances its overall spatial reasoning capability.

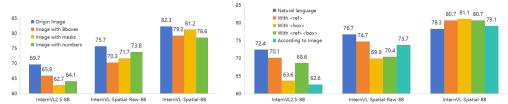


Figure 5: The results of the different image (Left) and text (Right) formats in the ablation study.

### 4.5 GENERAL VQA

For fairness, we re-evaluated InternVL2.5-8B (Chen et al., 2024c) under our experimental setup instead of directly using the results reported in its technical report. As shown in Table 4, InternVL-Spatial-8B achieves comparable performance to the baseline InternVL2.5-8B on general reasoning benchmarks. Specifically, InternVL-Spatial-8B shows a performance gain of +1.8% on Math-Vista (Wang et al., 2024b), -0.1% on OCRBench (Liu et al., 2024), +0.9% on TextVQA (Singh et al., 2019), -1.6% on ChartQA (Masry et al., 2022), and +0.2% on MMStar (Chen et al., 2024b). These results indicate that training with InternSpatial does not compromise the model's general reasoning capabilities, including mathematical reasoning, optical character recognition, visual question answering, and chart understanding.

Table 4: General benchmark results for InternVL2.5-8B vs. InternVL-Spatial-8B.

Model	MathVision (Wang et al., 2024b)	OCRBench (Liu et al., 2024)	TextVQA (Singh et al., 2019)	ChartQA (Masry et al., 2022)	MMStar (Chen et al., 2024b)
InternVL2.5-8B	19.0	82.3	79.0	83.0	62.9
InternVL-Spatial-8B	20.8(+1.8)	82.2(-0.1)	79.9(+0.9)	81.4(-1.6)	63.1(+0.2)

#### 5 CONCLUSIONS

We introduce InternSpatial, the largest open-source spatial reasoning dataset, and the benchmark InternSpatial-Bench, which together advance spatial understanding in VLMs through diverse scene coverage, rich instruction formats, and multi-view supervision. InternSpatial provides 12M high-quality QA pairs covering both single-view and multi-view settings, with broad scene diversity and 19 instruction formats that reflect the varied ways users express spatial queries. InternSpatial-Bench complements this with a diagnostic single-view benchmark and an extended multi-view evaluation via rotation angle prediction, a task not addressed in prior work. Extensive experiments show that training on InternSpatial yields substantial improvements on spatial reasoning benchmarks while maintaining strong performance on general multimodal tasks. Despite its scale and diversity, our template-based generation pipeline may underrepresent the full richness of natural language in real-world scenarios. Future work will explore more expressive QA generation and open-ended spatial reasoning in interactive environments. We anticipate that our dataset will support downstream applications such as robotics, embodied AI, and AR/VR, where spatial understanding is essential.

# REPRODUCIBILITY STATEMENT

All results reported in this paper are fully reproducible using the provided resources. The training configurations are detailed in Section 4 and Appendix E, while the dataset pipeline is described in Section 3, Appendix A, and Appendix B.

### ETHICS STATEMENT

Our work does not involve sensitive personal data. All dataset components were collected from open-source and publicly available sources, with careful filtering to exclude content that may be discriminatory or infringe copyright. We do not foresee any negative societal impacts arising from our methods or datasets.

### REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, Albert Q. Jiang, Kartik Khandelwal, Timothée Lacroix, Guillaume Lample, Diego Las Casas, Thibaut Lavril, Teven Le Scao, Andy Lo, William Marshall, Louis Martin, Arthur Mensch, Pavankumar Muddireddy, Valera Nemychnikova, Marie Pellat, Patrick Von Platen, Nikhil Raghuraman, Baptiste Rozière, Alexandre Sablayrolles, Lucile Saulnier, Romain Sauvestre, Wendy Shang, Roman Soletskyi, Lawrence Stewart, Pierre Stock, Joachim Studnia, Sandeep Subramanian, Sagar Vaze, Thomas Wang, and Sophia Yang. Pixtral 12b, 2024. URL https://arxiv.org/abs/2410.07073.
- N. Akkiraju, H. Edelsbrunner, M. Facello, P. Fu, and C. Varela. Alpha shapes: Definition and software. In GCG: International Computational Geometry Software Workshop, 1995.
- Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- Anthropic. The claude 3 model family: Opus, sonnet, haiku. https://www.anthropic.com, 2024. URL https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model\_Card\_Claude\_3.pdf.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL https://arxiv.org/abs/2502.13923.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 2020.
- Wenxiao Cai, Iaroslav Ponomarenko, Jianhao Yuan, Xiaoqi Li, Wankou Yang, Hao Dong, and Bo Zhao. Spatialbot: Precise spatial understanding with vision language models. In 2025 IEEE International Conference on Robotics and Automation (ICRA), 2025.
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14455–14465, June 2024a.

- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024b.
  - Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024c.
  - An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 135062–135093. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper\_files/paper/2024/file/f38cb4cf9a5eaa92b3cfa481832719c6-Paper-Conference.pdf.
  - Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv* preprint arXiv:2507.06261, 2025.
  - OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models. https://github.com/open-compass/opencompass, 2023.
  - Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
  - Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Niessner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
  - Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. *arXiv preprint arXiv:2212.08051*, 2022.
  - Yining Hong, Chunru Lin, Yilun Du, Zhenfang Chen, Joshua B. Tenenbaum, and Chuang Gan. 3d concept learning and reasoning from multi-view images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023a.
  - Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. *NeurIPS*, 2023b.
  - Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.
  - Mu Hu, Wei Yin, Chi Zhang, Zhipeng Cai, Xiaoxiao Long, Hao Chen, Kaixuan Wang, Gang Yu, Chunhua Shen, and Shaojie Shen. Metric3d v2: A versatile monocular geometric foundation model for zero-shot metric depth and surface normal estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10579–10596, 2024. doi: 10.1109/TPAMI.2024. 3444912.
  - Linyi Jin, Jianming Zhang, Yannick Hold-Geoffroy, Oliver Wang, Kevin Blackburn-Matzen, Matthew Sticha, and David F. Fouhey. Perspective fields for single image camera calibration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 17307–17316, June 2023.
  - Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017.

- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024a. URL https://arxiv.org/abs/2408.03326.
  - Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024b.
  - Chengzu Li, Caiqi Zhang, Han Zhou, Nigel Collier, Anna Korhonen, and Ivan Vulić. Topviewrs: Vision-language models as top-view spatial reasoners, 2024c. URL https://arxiv.org/abs/2406.02537.
  - Jianing Li, Xi Nan, Ming Lu, Li Du, and Shanghang Zhang. Proximity qa: Unleashing the power of multi-modal large language models for spatial proximity analysis, 2024d. URL https://arxiv.org/abs/2401.17862.
  - Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv* preprint arXiv:2301.12597, 2023a.
  - KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023b.
  - Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On pretraining for visual language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26689–26699, 2024.
  - Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
  - Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: On the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12), 2024. ISSN 1869-1919. doi: 10.1007/s11432-024-4235-6. URL http://dx.doi.org/10.1007/s11432-024-4235-6.
  - Zhaoyang Liu, Yinan He, Wenhai Wang, Weiyun Wang, Yi Wang, Shoufa Chen, Qinglong Zhang, Yang Yang, Qingyun Li, Jiashuo Yu, et al. Interngpt: Solving vision-centric tasks by interacting with chatbots beyond language. *arXiv preprint arXiv:2305.05662*, 2023.
  - Yongsen Mao, Yiming Zhang, Hanxiao Jiang, Angel Chang, and Manolis Savva. Multiscan: Scalable rgbd scanning for 3d environments with articulated objects. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 9058–9071. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper\_files/paper/2022/file/3b3a83a5d86e1d424daefed43d998079-Paper-Conference.pdf.
  - Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Proceedings of the annual meeting of the Association for Computational Linguistics*, pp. 2263–2279, 2022.
  - Inc. Meta Platforms. Llama-4-scout-17b-16e-instruct, 2025. URL https://huggingface.co/meta-llama/Llama-4-scout-17B-16E-Instruct. Accessed: 2024-05-15.
  - OpenAI. Gpt-4o system card. https://openai.com/index/gpt-4o-system-card/, 2025.
  - Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. URL https://arxiv.org/abs/2408.00714.

- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530, 2024.
  - Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8317–8326, 2019.
  - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
  - Johanna Wald, Armen Avetisyan, Nassir Navab, Federico Tombari, and Matthias Niessner. Rio: 3d object instance re-localization in changing indoor environments. In *Proceedings IEEE International Conference on Computer Vision (ICCV)*, 2019.
  - Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, Yixuan Li, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 75392–75421. Curran Associates, Inc., 2024a. URL https://proceedings.neurips.cc/paper\_files/paper/2024/file/89cc5e613d34f90de90c2le996e60b30-Paper-Conference.pdf.
  - Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multi-modal mathematical reasoning with math-vision dataset, 2024b. URL https://arxiv.org/abs/2402.14804.
  - Weiyun Wang, Min Shi, Qingyun Li, Wenhai Wang, Zhenhang Huang, Linjie Xing, Zhe Chen, Hao Li, Xizhou Zhu, Zhiguo Cao, et al. The all-seeing project: Towards panoptic visual recognition and understanding of the open world. In *ICLR*, 2024c.
  - Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *NeurIPS*, 2023.
  - Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces, 2024. URL https://arxiv.org/abs/2412.14171.
  - Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-next: A strong zero-shot video understanding model, April 2024. URL https://llava-vl.github.io/blog/2024-04-30-llava-next-video/.
  - Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. arXiv:1801.09847, 2018.
  - Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint* arXiv:2304.10592, 2023a.
  - Shengjie Zhu, Abhinav Kumar, Masa Hu, and Xiaoming Liu. Tame a wild camera: In-the-wild monocular camera calibration. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 45137–45149. Curran Associates, Inc., 2023b. URL https://proceedings.neurips.cc/paper\_files/paper/2023/file/8db9279f593652ee9bb2223b4a2c43fa-Paper-Conference.pdf.
  - Xizhou Zhu, Jinguo Zhu, Hao Li, Xiaoshi Wu, Hongsheng Li, Xiaohua Wang, and Jifeng Dai. Uni-perceiver: Pre-training unified architecture for generic perception for zero-shot and few-shot tasks. In *CVPR*, 2022.

### **APPENDIX**

### A EXPLANATION AND STATISTICS OF TASKS

InternSpatial and InternSpatial-Bench covers a total of 10 spatial reasoning tasks. The explanations of each task are shown in Table 5. We also count the number of QAs for each task in InternSpatial and InternSpatial-Bench, which are shown in Table 6 and Table 7 respectively.

Table 5: Explanation of tasks

Task	Description
Position Comparison	Compare the position of two objects in an image, involving three pairs of positional relationship: left/right, above/below, near/far.
Size Comparison	Compare the size of two objects in an image, involving three pairs of size relationship: wider/thinner, taller/shorter, larger/smaller.
Existence Estimation	Determine whether there are objects in the image whose positional/size relationships with the specified object meet the constraint conditions.
Object Counting	Estimate how many objects that meet the constraint conditions there are in a single image or multiple images.
<b>Rotation Estimation</b>	Estimate the rotation angle of an object between two images.
Absolute Distance	Estimate the closest distance between two objects given a serial of images.
Room Size	Estimate the volume of the room(s) given a serial of images.
Object Size	Estimate the longest dimension of an object given a serial of images.
Route Plan	Given a serial of images, choose what action should be performed be- tween a sequence of actions in order to route to from a start point to a target.
Appearance Order	Given a serial of images, determine the first-time appearance order of several objects.

Table 6: Statistics of tasks in InternSpatial

Table 6. Statistics of tasks in Internopatian						
Task	Related Views	# of QAs				
Position Comparison	Single	6,214,628				
Size Comparison	Single	3,227,124				
<b>Existence Estimation</b>	Single	50,845				
Object Counting	Single/Multiple	53,866				
<b>Rotation Estimation</b>	Multiple	2,464,500				
Absolute Distance	Multiple	14,596				
Room Size	Multiple	1,181				
Object Size	Multiple	3,709				
Route Plan	Multiple	4,966				
Appearance Order	Multiple	8,562				

# B TEMPLATES FOR GENERATING QAS IN INTERNSPATIAL

The QAs in InternSpatial are generated by template-based generation method. Here we provide the full list of templates. The "[...]" in templates are placeholders which will be replaced by object references in different formats, values, choices, and so on. Several candidates are provided to be randomly selected in generation process to enrich the structure of sentences.

Table 7: Statistics of tasks in InternSpatial-Bench

Task	Position	Size	Rotation	Object	Existence
	Comparison	Comparison	Estimation	Counting	Estimation
# of QAs	1845	1855	409	899	1000

#### Listing 1: Templates for task *Position Comparison*

```
764
       above_predict_templates =
765
         "question_templates": [
766
           "[A] is placed higher than [B], isn't it?",
767
           "Can we say that [A] is positioned above [B]?",
768
           "Is it correct to assume that [A] is located at a higher level than [
          B]?",
769
           "Is [A] placed higher than [B]?"
770
771
         "positive_answer_templates": [
772
           "Absolutely, [A] is clearly positioned above [B].",
773
           "Without a doubt, [A] is situated at a higher elevation than [B].",
774
           "Indeed, [A] is placed higher than [B].",
           "Certainly, [A] is located above [B]."
775
         1,
776
         "negative_answer_templates": [
777
           "Not at all, [A] is actually below [B].",
778
           "Definitely not, [A] is positioned lower than [B].",
779
           "Sorry, but [A] is not higher than [B].",
           "Unfortunately, [A] is not placed above [B]."
780
         ]
782
       below_predict_templates = {
783
         "question_templates": [
784
           "[A] is placed lower than [B], right?",
           "Can we say that [A] is positioned below [B]?",
785
           "Is it correct to assume that [A] is situated lower than [B]?",
786
           "Is [A] placed lower than [B]?"
787
788
         "positive_answer_templates": [
789
           "Absolutely, [A] is clearly positioned below [B].",
           "Without a doubt, [A] is located lower than [B].",
790
           "Indeed, [A] is situated beneath [B].",
791
           "Certainly, [A] is found at a lower level than [B]."
792
793
         "negative_answer_templates": [
794
           "Not at all, [A] is actually higher than [B].",
           "Definitely not, [A] is positioned above [B].",
795
           "In fact, [A] is situated higher than [B].",
796
           "Quite the opposite, [A] is at a higher level than [B]."
797
         ]
798
       left_predict_templates = {
         "question_templates": [
800
           "[A] is more to the left of [B], isn't it?",
801
           "Can we say that [A] is positioned more to the left than [B]?",
802
           "Is it correct to assume that [A] is situated to the left of [B]?",
803
           "Is [A] more to the left of [B]?"
804
         1.
         "positive_answer_templates": [
805
           "Absolutely, [A] is clearly positioned to the left of [B].",
806
           "Without a doubt, [A] is more to the left compared to [B].",
807
           "Indeed, [A] is located to the left of [B].",
808
           "Certainly, [A] is on the left side when compared to [B]."
809
         "negative_answer_templates": [
```

```
810
           "Not at all, [A] is not more to the left of [B].",
811
           "Actually, [A] is not positioned to the left of [B].",
812
           "Contrary to that, [A] is not situated to the left of [B].",
813
           "In fact, [A] is not on the left side when compared to [B]."
        1
814
815
       right_predict_templates = {
816
         "question_templates": [
817
           "[A] is more to the right of [B], isn't it?",
818
           "Can we say that [A] is positioned further to the right than [B]?",
           "Is it correct to assume that [A] is located to the right side of [B
819
          1?",
820
           "Is [A] more to the right of [B]?"
821
         ],
822
         "positive_answer_templates": [
           "Absolutely, [A] is clearly positioned to the right of [B].",
823
           "Indeed, [A] is noticeably more to the right compared to [B].",
824
           "Without a doubt, [A] is situated further to the right than [B].",
825
           "Certainly, [A] is distinctly to the right of [B]."
826
827
         "negative_answer_templates": [
828
           "Not at all, [A] is actually to the left of [B].",
           "Definitely not, [A] is not positioned to the right of [B].",
829
           "In fact, [A] is on the left side of [B].",
830
           "Contrary to that, [A] is not further to the right than [B]."
831
        1
832
833
       near_predict_templates = {
         "question_templates": [
834
           "Is [A] positioned in front of [B]?",
           "Does [A] precede [B] in this arrangement?",
836
           "Is [A] in front of [B]?",
837
           "Is [A] closer to the observer than [B]?"
838
         "positive_answer_templates": [
839
           "Without a doubt, [A] stands nearer to the viewer than [B].",
840
           "Definitely, [A] is more proximate to the observer than [B].",
841
           "Indeed, [A] is in front of [B].",
842
           "Absolutely, [A] is before [B].",
843
         "negative_answer_templates": [
844
           "Not at all, [A] is not closer to the observer than [B].",
845
           "No, [A] is not in front of [B].",
846
           "Unfortunately, [A] is not ahead of [B].",
847
           "Definitely not, [A] is not closer to the observer than [B]."
848
        ]
849
       far_predict_templates = {
850
         "question_templates": [
851
           "Is [A] situated behind [B]?",
852
           "Does [A] lie behind [B]?",
           "Is [A] to the rear of [B]?",
853
           "Is [A] farther from the observer than [B]?"
854
855
         "positive_answer_templates": [
856
           "Indeed, [A] is behind [B].",
857
           "Yes, [A] is behind [B].",
           "Without a doubt, [A] maintains a greater distance from the observer
858
859
           "Certainly, [A] is positioned further away from the observer than [B
860
          ]."
861
         1,
862
         "negative_answer_templates": [
           "No, [A] is not behind [B].",
863
           "Incorrect, [A] is not behind [B].",
```

```
864
           "That's wrong. [A] is not positioned behind [B].",
865
           "Unfortunately, [A] is not to the rear of [B]."
866
         1
867
       above_choice_templates = {
868
         "question_templates": [
869
           "Which one is positioned at a higher elevation, [A] or [B]?",
870
           "In terms of altitude, which comes first, [A] or [B]?",
871
           "Who stands taller, [A] or [B]?",
872
           "Which is placed higher, [A] or [B]?"
873
         "answer_templates": [
874
           "[0] is the one that is placed higher.",
875
           "The higher position belongs to [0].",
876
           "[0] occupies the superior location.",
           "It is [O] that is situated at a greater height."
877
         1
878
879
       below_choice_templates = {
880
         "question_templates": [
881
           "Which is positioned closer to the ground, [A] or [B]?",
882
           "Which one is situated at a lower elevation, [A] or [B]?",
           "Which of these is nearer to the base level, [A] or [B]?",
883
           "Which is placed lower, [A] or [B]?"
884
885
         "answer_templates": [
886
           "[0] is placed lower.",
           "The lower position belongs to [0].",
887
           "[0] occupies the lower spot.",
888
           "Lower down, you'll find [0]."
         1
890
891
       left_choice_templates = {
892
         "question_templates": [
           "Which is positioned further to the left, [A] or [B]?",
893
           "In terms of leftward placement, which comes first, [A] or [B]?", "When considering the left side, which one is closer, [A] or [B]?",
894
895
           "Which is more to the left, [A] or [B]?"
896
         1,
897
         "answer_templates": [
           "[0] is located more to the left.",
898
           "The position of [0] is further to the left.",
899
           "In comparison, [0] stands out as being more on the left.",
900
           "It is evident that [O] is situated more towards the left."
901
         ]
902
       right_choice_templates = {
903
         "question_templates": [
904
           "Which is positioned further to the right, [A] or [B]?",
905
           "In terms of horizontal alignment, which one is more to the right, [A
906
           ] or [B]?",
907
           "When comparing their positions, which one is situated more to the
           right, [A] or [B]?",
908
           "Which is more to the right, [A] or [B]?"
909
910
         "answer_templates": [
911
           "[0] is clearly more to the right.",
           "The position of [O] is further to the right.",
912
           "Comparing the two, [0] is definitively more to the right.",
913
           "It is evident that [O] is positioned more to the right."
914
         ]
915
916
       near_choice_templates = {
917
         "question_templates": [
           "Which one is positioned further forward, [A] or [B]?",
```

```
918
           "Between [A] and [B], which object is closer to the observer?",
919
           "Can you identify which of the two, [A] or [B], is in the foremost
920
           position?",
921
           "Of the two, [A] and [B], which is closer to the front?"
         ],
922
         "answer_templates": [
923
           "[0] is in front.",
924
           "The frontmost object is [0].",
925
           "[0] is situated at the foremost position.",
926
           "Among the options, [0] is the one that is most ahead."
         ]
927
928
       far_choice_templates = {
929
         "question_templates": [
930
           "Which one is further back, [A] or [B]?",
           "Can you tell me which is positioned more towards the back, [\mathtt{A}] or [\mathtt{B}
931
932
           "Between [A] and [B], which is more distant in the rear aspect?",
933
           "Comparing [A] and [B], which is more behind?"
934
935
         "answer_templates": [
936
           "[0] is definitely more behind.",
           "I can confirm that [0] is situated further back.",
937
           "[0] is clearly more behind than the other.",
938
           "There is no question that [O] is more behind."
939
        ]
940
941
       above_below_choice_templates = {
         "question_templates": [
942
           "Is [A] positioned higher or lower than [B]?",
943
           "Does [A] lie above or beneath [B]?",
944
           "Is [A] situated over or under [B]?",
945
           "Is [A] above or below [B]?"
946
         "above_answer_templates": [
947
           "[A] is above [B].",
948
           "[A] is positioned higher than [B].",
949
           "[A] lies over [B].",
950
           "[A] is situated above [B]."
951
         ],
         "below_answer_templates": [
952
           "[A] is below [B].",
953
           "[A] is positioned lower than [B].",
954
           "[A] lies under [B].",
955
           "[A] is situated below [B]."
         ]
956
957
       left_right_choice_templates = {
958
         "question_templates": [
959
           "Is [A] relatively farther to the left or right than [B]?",
           "Does [A] lie on the left or right side of [B]?",
960
           "Is [A] to the left or right of [B]?"
961
962
         "left_answer_templates": [
963
           "[A] is to the left of [B].",
964
           "[A] occupies the left side relative to [B].",
965
           "[A] lies on the left side of [B]."
966
         ],
         "right_answer_templates": [
967
           "[A] is to the right of [B].",
968
           "[A] occupies the right side relative to [B].",
969
           "[A] lies on the right side of [B]."
970
971
       near_far_choice_templates = {
```

```
972
         "question_templates": [
973
           "Is [A] relatively nearer or farther from the observer than [B]?",
974
           "Can you determine if [A] is closer or farther from the observer
975
          compared to [B]?",
           "Is [A] in front of or behind [B]?"
976
977
         "near_answer_templates": [
978
           "[A] is closer to the observer than [B].",
979
           "[A] is more proximate to the observer than [B].",
           "[A] comes before [B].",
980
           "[A] is in front of [B]."
981
         1,
982
         "far_answer_templates": [
983
           "[A] is farther from the observer than [B].",
984
           "[A] is less proximate to the observer than [B].",
           "[A] is behind [B]."
985
986
987
```

### Listing 2: Templates for task Size Comparison

```
989
990
       wide_predict_templates = {
         "question_templates": [
991
           "Is [A] broader than [B]?",
992
           "Does [A] have a larger width compared to [B]?",
993
           "Can we say that [A] spans more horizontally than [B]?",
994
           "Is [A] wider than [B]?"
995
         "positive_answer_templates": [
996
           "Yes, [A] is noticeably broader than [B].",
           "Indeed, [A] has a significantly larger width than [B].",
998
           "Absolutely, [A] spans more horizontally than [B].",
999
           "Certainly, [A] is wider than [B]."
1000
         "negative_answer_templates": [
1001
           "No, [A] is not broader than [B].",
1002
           "In fact, [A] does not have a larger width compared to [B].",
1003
           "Sorry, but [A] does not span more horizontally than [B].",
1004
           "Unfortunately, [A] is not wider than [B]."
         1
1005
1006
       narrow_predict_templates = {
1007
         "question_templates": [
1008
           "Is [A] thinner than [B]?",
1009
           "Does [A] have a smaller width compared to [B]?",
1010
           "Is the width of [A] less than that of [B]?",
           "Is [A] narrower than [B]?"
1011
1012
         "positive_answer_templates": [
1013
           "Yes, [A] is noticeably narrower than [B].",
1014
           "Indeed, [A] has a significantly smaller width than [B].",
           "Absolutely, the width of [A] is less than that of [B].",
1015
           "Certainly, [A] is thinner than [B]."
1016
1017
         "negative_answer_templates": [
1018
           "No, [A] is not narrower than [B]; in fact, it's wider.",
1019
           "Definitely not; [A] has a larger width than [B].",
           "Not at all; the width of [A] exceeds that of [B].",
1020
           "No way; [A] is thicker than [B]."
1021
         ]
1022
1023
       tall_predict_templates = {
1024
         "question_templates": [
1025
           "[A] is taller than [B], isn't it?",
           "Can we say that [A] surpasses [B] in height?",
```

```
1026
           "Is it correct to assume that [A] is taller than [B]?",
1027
           "Is [A] taller than [B]?"
1028
1029
        "positive_answer_templates": [
           "Absolutely, [A] towers over [B].",
1030
           "Without a doubt, [A] is significantly taller than [B].",
1031
           "Indeed, [A] outshines [B] in terms of height.",
1032
           "Unquestionably, [A] is taller than [B]."
1033
1034
        "negative_answer_templates": [
           "Not at all, [B] is actually taller than [A].",
1035
           "Sorry, but [A] does not exceed [B] in height.",
1036
           "In fact, [B] surpasses [A] in height.",
1037
           "Regrettably, [A] falls short when compared to [B]'s height."
1038
        ]
1039
      vshort_predict_templates = {
1040
         "question_templates": [
1041
           "Is [A] shorter than [B] in vertical direction?",
1042
           "Does [A] have less height than [B]?",
1043
           "Is the vertical length of [A] smaller than that of [B]?",
           "Is [A] shorter than [B] in vertical direction?"
1044
1045
        "positive_answer_templates": [
1046
           "Yes, [A] is indeed shorter than [B] in the vertical direction.",
1047
           "Absolutely, [A] has less height compared to [B].",
1048
           "Certainly, the vertical length of [A] is smaller than that of [B].",
           "Without a doubt, [A] is shorter than [B] vertically."
1049
1050
        "negative_answer_templates": [
           "No, [A] is not shorter than [B] in the vertical direction.",
1052
           "Definitely not, [A] does not have less height than [B].",
1053
           "Not at all, the vertical length of [A] is not smaller than that of [A]
          B].",
1054
           "Certainly not, [A] is not shorter than [B] vertically."
1055
1056
1057
      large_predict_templates = {
1058
         "question_templates": [
           "[A] is larger than [B], isn't it?",
1059
           "Can we say that [A] has a bigger size compared to [B]?",
1060
           "Is it correct to assume that [A] surpasses [B] in size?",
1061
           "Is [A] larger than [B]?"
1062
1063
        "positive_answer_templates": [
           "Absolutely, [A] is noticeably larger than [B].",
1064
           "Without a doubt, [A] outsizes [B] significantly.",
1065
           "Indeed, [A] is clearly more expansive than [B].",
1066
           "Definitely, [A] dwarfs [B] in terms of size."
1067
1068
        "negative_answer_templates": [
           "Not at all, [B] is actually larger than [A].",
1069
           "Quite the opposite, [B] surpasses [A] in size.",
1070
           "In fact, [B] is the larger one when compared to [A].",
1071
           "Sorry, but [B] is bigger than [A]."
1072
        1
1073
      small_predict_templates = {
1074
         "question_templates": [
1075
           "[A] is smaller than [B], isn't it?",
1076
           "Can we say that [A] is smaller than [B]?",
1077
           "Is it true that [A] is smaller than [B]?",
1078
           "Is [A] smaller than [B]?"
1079
        "positive_answer_templates": [
```

```
1080
           "Absolutely, [A] is noticeably smaller than [B].",
1081
           "Yes, [A] is indeed smaller than [B].",
1082
           "Without a doubt, [A] is smaller than [B].",
1083
           "Definitely, [A] is smaller than [B]."
         ],
1084
         "negative_answer_templates": [
1085
           "Not at all, [A] is actually larger than [B].",
1086
           "No, [A] is not smaller than [B].",
1087
           "Quite the opposite, [A] is bigger than [B].",
1088
           "False, [A] is not smaller than [B]."
         ]
1089
1090
       wide_choice_templates = {
1091
         "question_templates": [
1092
           "Which has a greater width, [A] or [B]?",
           "In terms of width, which one is larger, [A] or [B]?",
1093
           "When comparing widths, which one comes out on top, [A] or [B]?",
1094
           "Which is wider, [A] or [B]?"
1095
1096
         "answer_templates": [
1097
           "[0] is wider."
           "The width of [O] is greater.",
1098
           "Comparing the two, [0] has the larger width.",
1099
           "In terms of width, [O] surpasses the other."
1100
1101
1102
       narrow_choice_templates = {
1103
         "question_templates": [
           "Which has a smaller width, [A] or [B]?",
1104
           "In terms of width, which one is less, [A] or [B]?",
1105
           "When comparing widths, which one comes out smaller, [A] or [B]?",
1106
           "Which is narrower, [A] or [B]?"
1107
         ],
         "answer_templates": [
1108
           "[0] is the narrower one.",
1109
           "The narrower object is [0].",
1110
           "[0] has the lesser width."
1111
           "Comparing the two, [O] is clearly narrower."
1112
        1
1113
       tall_choice_templates = {
1114
         "question_templates": [
1115
           "Which has a greater height, [A] or [B]?",
1116
           "In terms of height, which one is superior, [A] or [B]?",
1117
           "When comparing heights, which comes out on top, [A] or [B]?",
           "Which is taller, [A] or [B]?"
1118
1119
         "answer_templates": [
1120
           "[0] is the taller one.",
1121
           "The height of [O] surpasses the other.",
1122
           "[0] stands out as the taller between the two.",
           "Comparatively speaking, [O] is taller."
1123
1124
1125
       vshort_choice_templates = {
1126
         "question_templates": [
1127
           "Which has a shorter vertical length, [A] or [B]?",
           "In terms of vertical measurement, which one is shorter, [A] or [B]?"
1128
1129
           "When comparing the vertical dimensions, which is shorter, [A] or [B]
1130
          1?",
1131
           "Which is shorter in vertical direction, [A] or [B]?"
1132
1133
         "answer_templates": [
           "[0] is shorter in the vertical direction.",
```

```
1134
           "The vertical length of [O] is less than the other.",
1135
           "Comparing vertically, [0] comes out shorter.",
1136
           "In terms of height, [O] is the shorter one."
1137
         1
1138
       large_choice_templates = {
1139
         "question_templates": [
1140
           "Which has a greater size, [A] or [B]?",
1141
           "In terms of size, which one is bigger, [A] or [B]?",
1142
           "When comparing sizes, which one comes out on top, [A] or [B]?",
           "Which is larger, [A] or [B]?"
1143
         1,
1144
         "answer_templates": [
1145
           "[0] is the larger one.",
1146
           "The bigger size belongs to [0].",
           "[0] surpasses the other in size."
1147
           "Comparatively speaking, [0] is the larger."
1148
1149
1150
       small_choice_templates = {
1151
         "question_templates": [
           "Which has a smaller size, [A] or [B]?",
1152
           "In terms of size, which one is smaller, [A] or [B]?",
1153
           "When comparing sizes, which one comes out smaller, [A] or [B]?",
1154
           "Which is smaller, [A] or [B]?"
1155
1156
         "answer_templates": [
           "[0] is the smaller one.",
1157
           "The smaller of the two is [O].",
1158
           "[0] has the smaller size.",
           "Comparing the two, [0] is the smaller."
1160
1161
       wide_narrow_choice_templates = {
1162
         "question_templates": [
1163
           "Is [A] relatively wider or narrower than [B]?",
1164
           "How does the width of [A] compare to [B]?",
1165
           "Can you tell me if [A] has a greater or lesser width than [B]?",
1166
           "Is [A] wider or narrower than [B]?"
1167
         "wide_answer_templates": [
1168
           "[A] is wider than [B].",
1169
           "The width of [A] exceeds that of [B].",
1170
           "[A] has a larger width compared to [B].",
1171
           "In terms of width, [A] surpasses [B]."
1172
         "narrow_answer_templates": [
1173
           "[A] is narrower than [B].",
1174
           "The width of [B] is greater than that of [A].",
1175
           "[A] has a smaller width compared to [B].",
           "In terms of width, [B] surpasses [A]."
1176
1177
1178
       tall_short_choice_templates = {
1179
         "question_templates": [
1180
           "Is [A] relatively taller or shorter than [B]?",
1181
           "How does the height of [A] compare to [B]?",
1182
           "Can you determine if [A] is taller or shorter than [B]?",
           "Is [A] taller or shorter than [B]?"
1183
1184
         "tall_answer_templates": [
1185
           "[A] is taller than [B].",
1186
           "The height of [A] exceeds that of [B].",
           "[A] surpasses [B] in height."
1187
           "Compared to [B], [A] is definitely taller."
```

```
1188
1189
         "short_answer_templates": [
1190
           "[A] is shorter than [B].",
1191
           "In terms of height, [A] falls below [B].",
           "[B] is taller than [A].",
1192
           "[A]'s height is less than that of [B]."
1193
         ]
1194
1195
       large_small_choice_templates = {
1196
         "question_templates": [
           "Is [A] relatively larger or smaller than [B]?",
1197
           "How does the size of [A] compare to [B]?",
1198
           "Can you determine if [A] is bigger or smaller than [B]?",
1199
           "Is [A] larger or smaller than [B]?"
1200
         "large_answer_templates": [
1201
           "[A] is larger than [B].",
1202
           "The size of [A] exceeds that of [B].",
1203
           "[A] surpasses [B] in size.",
1204
           "Compared to [B], [A] is bigger."
1205
1206
         "small_answer_templates": [
           "[A] is smaller than [B].",
1207
           "The size of [A] is less than that of [B].",
1208
           "[B] is larger than [A].",
1209
           "In comparison to [B], [A] is smaller."
1210
1211
1212
```

### Listing 3: Templates for task Existence Estimation

```
1214
       existence_left_templates = {
1215
         "question_templates": [
1216
           "Does [B] exist to the left of [A]?",
           "Is there [B] to the left of [A]?",
1217
           "Is there [B] more to the left than [A]?"
1218
1219
         "positive_answer_templates": [
1220
           "Yes."
1221
         "negative_answer_templates": [
1222
           "No."
1223
1224
1225
       existence_right_templates = {
1226
         "question_templates": [
           "Is there [B] positioned more to the right than [A]?",
1227
           "Does [B] exist to the right of [A]?",
1228
           "Is there [B] locating to the rightside of [A]?"
1229
1230
         "positive_answer_templates": [
           "Yes."
1231
1232
         "negative_answer_templates": [
1233
           "No."
1234
1235
1236
       existence_above_templates = {
         "question_templates": [
1237
           "Does [B] exist at higher elevation than [A]?",
1238
           "Can you find [B] above [A]?",
1239
           "Is there [B] that is located above [A]?"
1240
1241
         "positive_answer_templates": [
           "Yes."
```

```
1242
1243
         "negative_answer_templates": [
1244
           "No."
1245
1246
       existence_below_templates = {
1247
         "question_templates": [
1248
           "Is there [B] that is situated below [A]?",
1249
           "Does [B] exist below [A]?",
1250
           "Is there [B] positioned lower than [A]?"
1251
1252
         "positive_answer_templates": [
1253
           "Yes."
1254
         "negative_answer_templates": [
1255
           "No."
1256
1257
1258
       existence_near_templates = {
1259
         "question_templates": [
           "Does [B] exist near [A]?",
1260
           "Is there [B] that is in front of [A]?",
1261
           "Can you find [B] that is closer than observer than [A]?"
1262
1263
         "positive_answer_templates": [
1264
           "Yes."
1265
         "negative_answer_templates": [
1266
           "No."
1267
1268
1269
       existence_far_templates = {
         "question_templates": [
1270
           "Does [B] exist far from [A]?",
1271
           "Is there [B] that is behind [A]?",
1272
           "Does [B] exist behind [A]?"
1273
1274
         "positive_answer_templates": [
           "Yes."
1275
1276
         "negative_answer_templates": [
1277
           "No."
1278
1279
1280
       existence_wide_templates = {
         "question_templates": [
1281
           "Can you find [B] that is wider than [A]?",
1282
           "Is there [B] that is wider than [A]?",
1283
           "Is there [B] that has a larger extent in horizontal than [A]?"
1284
1285
         "positive_answer_templates": [
1286
           "Yes."
1287
1288
         "negative_answer_templates": [
1289
           "No."
1290
1291
       existence_narrow_templates = {
1292
         "question_templates": [
1293
           "Is there [B] that is narrower than [A]?",
1294
           "Can you find [B] that is narrower than [A]?",
           "Does [B] with smaller width than [A] exist?"
1295
        ],
```

```
1296
         "positive_answer_templates": [
1297
           "Yes."
1298
1299
         "negative_answer_templates": [
           "No."
1300
1301
1302
       existence_tall_templates = {
1303
         "question_templates": [
1304
           "Is there [B] that is taller than [A]?",
           "Can you find [B] that has a larger height than [A]?"
1305
           "Is there [B] that is larger in vertical than [A]?"
1306
1307
         "positive_answer_templates": [
1308
           "Yes."
1309
         "negative_answer_templates": [
1310
           "No."
1311
1312
1313
       existence_vshort_templates = {
1314
         "question_templates": [
           "Is there [B] that is shorter than [A] in vertical?",
1315
           "Does [B] shorter than [A] exists?",
1316
           "Is there [B] that has a smaller height than [A]?"
1317
1318
         "positive_answer_templates": [
           "Yes."
1319
1320
         "negative_answer_templates": [
1321
           "No."
1322
1323
1324
       existence_large_templates = {
         "question_templates": [
1325
           "Can you find [B] that is larger than [A]?",
1326
           "Is there [B] that is larger than [A]?",
1327
           "Does [B] exist that has a larger volume than [A]?"
1328
         "positive_answer_templates": [
1329
           "Yes."
1330
1331
         "negative_answer_templates": [
1332
           "No."
1333
1334
       existence_small_templates = {
1335
         "question_templates": [
1336
           "Is there [B] that is smaller in size than [A]?",
1337
           "Does [B] with smaller size than [A] exist?",
1338
           "Does [B] exist that is smaller than [A]?"
1339
         "positive_answer_templates": [
1340
           "Yes."
1341
1342
         "negative_answer_templates": [
1343
           "No."
         1
1344
1345
1346
```

Listing 4: Templates for task *Object Counting* 

```
count_above_templates = {
   "question_templates": [
        "How many [B] are located higher than [A]?",
```

```
1350
               "How many [B] are positioned higher than [A]?",
1351
               "How many [B] are above [A]?"
1352
1353
           "answer_templates": [
               "[V]."
1354
1355
1356
       count_below_templates = {
1357
           "question_templates": [
1358
               "How many [B] are lower than [A]?",
               "How many [B] are situated below [A]?",
1359
               "How many [B] are positioned lower than [A]?"
1360
1361
           "answer_templates": [
1362
               "[V]."
           1
1363
1364
       count_left_templates = {
1365
           "question_templates": [
1366
               "How many [B] are positioned to the left of [A]?",
1367
               "How many [B] are more to the left than [A]?",
               "How many [B] are on the leftside of [A]?"
1368
1369
           "answer_templates": [
1370
               "[V]."
1371
1372
1373
       count_right_templates = {
           "question_templates": [
1374
               "How many [B] are found to the right of [A]?",
               "How many [B] lie to the rightside of [A]?",
1376
               "How many [B] are more to the right than [A]?"
1377
           "answer_templates": [
1378
               "[V]."
1379
1380
1381
       count_near_templates = {
1382
           "question_templates": [
               "How many [B] are closer to the observer than [A]?",
1383
               "How many [B] are in front of [A]?",
1384
               "How many [B] are located nearer to the observer than [A]?"
1385
1386
           "answer_templates": [
1387
               "[V]."
1388
1389
       count_far_templates = {
1390
           "question_templates": [
1391
               "How many [B] are positioned farther from the observer than [A]?"
1392
               "How many [B] are located behind [A]?"
1393
               "How many [B] are farther from the observer than [A]?"
1394
1395
           "answer_templates": [
1396
               "[V]."
1397
1398
       count_wide_templates = {
1399
           "question_templates": [
1400
               "How many [B] have a larger width compared to [A]?",
1401
               "How many [B] are wider than [A]?",
1402
               "How many [B] have a larger extent in horizontal than [A]?"
1403
           "answer_templates": [
```

```
1404
               "[V]."
1405
           ]
1406
1407
       count_narrow_templates = {
           "question_templates": [
1408
               "How many [B] are narrower than [A]?",
1409
               "How many [B] have a less width than that of [A]?",
1410
                "How many [B] are thinner than [A]?"
1411
           1,
1412
           "answer_templates": [
               "[V]."
1413
1414
1415
       count_tall_templates = {
1416
           "question_templates": [
                "How many [B] are taller than [A]?",
1417
                "How many [B] surpass [A] in height?",
1418
                "How many [B] have a larger extent in vertical than [A]?"
1419
1420
           "answer_templates": [
1421
               "[V]."
1422
           ]
1423
       count_vshort_templates = {
1424
           "question_templates": [
1425
               "How many [B] have less height than [A]?",
1426
               "How many [B] are shorter than [A]?",
               "How many [B] have a smaller vertical length than that of [A]?"
1427
1428
           "answer_templates": [
               "[V]."
1430
1431
1432
       count_large_templates = {
           "question_templates": [
1433
                "How many [B] are larger than [A]?",
1434
               "How many [B] have a bigger size compared to [A]?",
1435
               "How many [B] surpass [A] in size?"
1436
           "answer_templates": [
1437
               "[V]."
1438
1439
1440
       count_small_templates = {
1441
           "question_templates": [
               "How many [B] have a smaller size compared to [A]?"
1442
               "How many [B] are smaller than [A]?",
1443
               "How many [B] are smaller in volume than [A]?"
1444
           ],
1445
           "answer_templates": [
1446
               "[V]."
1447
1448
1449
                              Listing 5: Templates for multi-view tasks
1450
```

```
object_rotation_predict_templates = {
  "question_templates": [
    "Here are two images of the same object:\nImage 1:\n<image>\nImage
    2:\n<image>\nPlease estimate how [A] in image 2 is rotated relative
    to image 1?",
    "Here are two images of the same object:\nImage 1:\n<image>\nImage
    2:\n<image>\nIn what direction and by what angle has [A] in image 2
    been rotated from its position in image 1?"
],
```

1452

1453

1454

1455

1456

```
1458
         "clockwise_answer_templates": [
1459
           "[A] rotates about [D] degrees clockwise.",
1460
           "[A] turns clockwise by about [D] degrees.",
1461
           "[A] undergoes approximately a [D] degree clockwise rotation."
        ],
1462
         "counterclockwise_answer_templates": [
1463
           "[A] rotates abount [D] degrees counterclockwise.",
1464
           "[A] turns counterclockwise by about [D] degrees.",
1465
           "[A] undergoes approximately a [D] degree counterclockwise rotation."
1466
         "rotate_180_answer_templates": [
1467
           "[A] rotates about [D] degrees.",
1468
           "[A] turns by about [D] degrees.",
1469
           "[A] undergoes approximately a [D] degree rotation."
1470
1471
      route_plan_templates = {
1472
         "question_templates": [
1473
           "Image-1: <image>\nImage-2: <image>\nImage-3: <image>\nImage-4: <
1474
          image>\nImage-5: <image>\nImage-6: <image>\nImage-7: <image>\nImage
1475
          -8: <image>\nImage-9: <image>\nImage-10: <image>\nImage-11: <image>\
1476
          nImage-12: <image>\nImage-13: <image>\nImage-14: <image>\nImage-15: <
          image>\nImage-16: <image>\nImage-17: <image>\nImage-18: <image>\
1477
          nImage-19: <image>\nImage-20: <image>\nImage-21: <image>\nImage-22: <
1478
          image>\nImage-23: <image>\nImage-24: <image>\nYou are a robot
1479
          beginning at the column and facing the staircase. You want to
1480
          navigate to the grand staircase. You will perform the following
1481
          actions (Note: for each [please fill in], choose either 'turn back,'
          'turn left,' or 'turn right.'): 1. Go forward until the columns. 2. [
1482
          please fill in]. 3. Go forward until the steps. 4. Stop on the
          landing.\nA. Turn Right\nB. Turn Left\nC. Turn Back\nAnswer with the
1484
          option's letter from the given choices directly."
1485
         "answer_templates": [
1486
          "[0]"
1487
1488
1489
      abs_dist_templates = {
1490
         "question_templates": [
           "Image-1: <image>\nImage-2: <image>\nImage-3: <image>\nImage-4: <
1491
          image>\nImage-5: <image>\nImage-6: <image>\nImage-7: <image>\nImage
1492
          -8: <image>\nImage-9: <image>\nImage-10: <image>\nImage-11: <image>\
1493
          nMeasuring from the closest point of each object, what is the
1494
          distance between [A] and [B] (in meters)?\nPlease answer the question
1495
           using a single word or phrase."
1496
         "answer_templates": [
1497
           "[V]"
1498
1499
1500
      obj_count_templates = {
         "question_templates": [
1501
           "Image-1: <image>\nImage-2: <image>\nImage-3: <image>\nImage-4: <
1502
          image>\nImage-5: <image>\nImage-6: <image>\nImage-7: <image>\nImage
1503
          -8: <image>\nImage-9: <image>\nImage-10: <image>\nImage-11: <image>\
1504
          nThese are frames of a video.\nHow many [A](s) are in this room?\
1505
          nPlease answer the question using a single word or phrase."
1506
        ],
        "answer_templates": [
1507
          "[V]"
1508
1509
1510
      room_size_templates = {
1511
        "question_templates": [
```

```
1512
           "Image-1: <image>\nImage-2: <image>\nImage-3: <image>\nImage-4: <
1513
          image>\nImage-5: <image>\nImage-6: <image>\nImage-7: <image>\nImage
1514
          -8: <image>\nImage-9: <image>\nImage-10: <image>\nImage-11: <image>\
1515
          nThese are frames of a video. \nWhat is the size of this room (in
          square meters)? \nIf multiple rooms are shown, estimate the size of
1516
          the combined space. \nPlease answer the question using a single word
1517
          or phrase."
1518
1519
        "answer_templates": [
1520
          "[V]"
1521
1522
      rel_dist_templates = {
1523
         "question_templates": [
1524
          "Image-1: <image>\nImage-2: <image>\nImage-3: <image>\nImage-4: <
          image>\nImage-5: <image>\nImage-6: <image>\nImage-7: <image>\nImage
1525
          -8: <image>\nImage-9: <image>\nImage-10: <image>\nThese are frames of
1526
           a video.\nMeasuring from the closest point of each object, which of
1527
          these objects ([B0],[B1],[B2],[B3]) is the closest to [A]?\nA. [B0]\
1528
          nB. [B1]\nC. [B2]\nD. [B3]\nAnswer with the option's letter from the
1529
          given choices directly."
        1,
1530
        "answer_templates": [
1531
           "[0]"
1532
1533
1534
      object_size_templates =
1535
         "question_templates": [
           "Image-1: <image>\nImage-2: <image>\nImage-3: <image>\nImage-4: <
1536
          image>\nImage-5: <image>\nImage-6: <image>\nImage-7: <image>\nImage
1537
          -8: <image>\nImage-9: <image>\nImage-10: <image>\nImage-11: <image>\
1538
          nThese are frames of a video.\nWhat is the length of the longest
1539
          dimension (length, width, or height) of [A], measured in centimeters
1540
          ?\nPlease answer the question using a single word or phrase."
1541
         "answer_templates": [
1542
           "[V]"
1543
1544
      appear_order_templates = {
1545
         "question_templates": [
1546
          "Image-1: <image>\nImage-2: <image>\nImage-3: <image>\nImage-4: <
1547
          image>\nImage-5: <image>\nImage-6: <image>\nImage-7: <image>\nImage
1548
          -8: <image>\nImage-9: <image>\nImage-10: <image>\nImage-11: <image>\
1549
          nThese are frames of a video.\nWhat will be the first-time appearance
1550
           order of the following categories in the video: [A0], [A1],
          A3]?\nA. [B0]\nB. [B1]\nC. [B2]\nD. [B3]\nAnswer with the option's
1551
          letter from the given choices directly."
1552
        1.
1553
        "answer_templates": [
           "[0]"
1555
1556
```

### C VLM-ASSISTED ANNOTATION FOR INTERNSPATIAL

1557 1558

1559 1560

1561

1562

1563

1564

1565

As described in Dataset section, we involved open-source VLM to do the object detection, captioning, and grounding in the pipeline of InternSpatial generation. We use QWen2.5-VL 72B(Bai et al., 2025) as the assistant. For each process, we design corresponding prompt to make the VLM understand what should do and what should output. Here we provide the prompts for these processes.

Listing 6: Prompts for detecting objects in images

```
messages = [{"role": "system", "content": f"""
```

```
1566
        You are an object detector. Given an image, you should find all objects
1567
           in image with grounding. The term "object" includes all living and
1568
          non-living things. For each detected object, you should assign a
          label, which represents what the object is. You should also describe
1569
          each detected object in detail with a phrase. The description can
1570
          contain appearance, function, action, etc.
1571
1572
        Output format: The response should be in json format, which contains a
1573
          list of dicts. Each dict is for an detected object and has three keys
          : "label" for the label, "caption" for the description and "box" for
1574
          the grounding. The description should be lowercases and no period at
1575
          the end. The grounding should be a list of four ints [x1, y1, x2, y2
1576
          ], where (x1, y1) is the top-left coord and (x2, y2) is the bottom-
1577
          right coord. Compact the responsed json in one line.
1578
      messages.append({"role": "user", "content": '\n'.join(query)})
1579
```

### Listing 7: Prompts for captioning objects given bounding boxes

```
messages = [{"role": "system", "content": f"""
  You are an language assistant. You will be given an array dict. Each
  dict contains a field "box" for grounding box and an optional field "
  label" for label of a object in image. Your task is to generate brief
  descriptions with less than ten words for these objects. Output one
  description per line.

Here is an example:

Input:
[{"box": [10, 20, 300, 400], "label": "bus"}, {"box": [42, 512, 64,
  890]}]

Output:
  a blue bus seat with a suitcase partially resting on it
  a red car on right side
"""}]
messages.append({"role": "user", "content": '\n'.join(query)})
```

### Listing 8: Prompts for grounding objects given captions

```
messages = [{"role": "system", "content": f"""
  You are a professional image annotator. I will give you an image and a
  phrase about one or more objects in the image. Please detect all
  objects matching the phrase. The response should be a JSON object,
  containing a field "boxes". "boxes" is a list of grounding boxes [x1,
    y1, x2, y2].

Example Output:
{
    "boxes": [[192, 29, 321, 49], [19, 65, 392, 569], [59, 102, 439,
    139]]
  }
"""}]
messages.append({"role": "user", "content": '\n'.join(query)})
```

### D More Details about InternSpatial-Bench Generation Pipeline

The generation pipeline of InternSpatial-Bench does not rely on existing annotations. Starting from the images, we carried out four steps including image filtering, image captioning, question design, and object grounding to obtain the necessary 2D annotations for the questions of the benchmark and the generation of answers. In this steps, we design prompts respectively to enable the Visual Language Model (VLM) to automatically generate intermediate results. These prompts are presented in Listing 9, Listing 10, Listing 11, Listing 12, Listing 13, and Listing 8. Subsequently, we reused

 the processes in the second and third stages of the training dataset pipeline to generate answers and expand the instruction formats. After generated the QA pairs, we invited experienced human annotators to conduct manual verification of all the pairs to ensure the quality of the benchmark.

### Listing 9: Prompts for filtering images

```
messages = [{"role": "system", "content": f"""
  You are a helpful visual assistant. Please determine whether the
    following conditions are met:
1. 5-or-more-objects: There are at least 5 objects in the image.
2. cartoon: This is a cartoon image.
3. group-of-images: This is a group of images.
4. screenshot: This is a screenshot.
5. realistic-image: This is a realistic image captured by camera.

For each condition, answer true of false. Response in JSON dict with
    five fields: "5-or-more-objects", "cartoon", "group-of-images", "
    screenshot", "realistic-image
"""}]
messages.append({"role": "user", "content": '\n'.join(query)})
```

### Listing 10: Prompts for captioning images

```
messages = [{"role": "system", "content": f"""
  You are a helpful visual assistant. Please describe the image as detail
   as possible. Then detect all top-level objects and return their
   detailed descriptions (top-level means it's not a part of another
   object).
"""}]
messages.append({"role": "user", "content": '\n'.join(query)})
```

### Listing 11: Prompts for design questions of task *Position Comparison*

```
1647
      messages = [{"role": "system", "content": f"""
1648
        I will give you an image and a description about the image. You should
          design 2 questions regarding position judgments around top-level
1649
          objects in the image (top-level means it's not a part of another
1650
          object). The question should involve two object (anchor, target) and
1651
          a type of relationship:
1652
        - *The anchor and target object* should be randomly chosen from the top
1653
          -level objects. You possibly need to add the attributions about
1654
          appearance, behavior, posture, position in the image, etc. to the
1655
          description of the anchor and target objects so that they can be
1656
          distinguished from others.
1657
         - *The relationship* should be randomly selected from: more to the left
1658
          , more to the right, closer (to the observer), farther (from the
          observer), higher, lower.
1659
        - Only design questions about top-level objects. Ignore those not in
1660
          top-level objects list. Ignore environment objects such as water, sky
1661
          , grass, cloud, etc.
        After that, generate 2 more questions based on the designed questions
1663
          by choose another relationship and keep other parts unchanged.
1664
1665
        Please respond in JSON format. All content should be in English. Here
1666
          is an example of output:
1667
           "questions": [
1668
1669
               "question": "Is the wooden chair positioned higher than the blue
1670
          table?",
1671
               "anchor": "blue table",
1672
               "target": "wooden chair"
               "relationship": "higher",
1673
               "task": "position"
```

```
1674
1675
1676
               "question": "Does the red bicycle locate more to the left than
1677
           the man in a floral shirt?",
               "anchor": "man in a floral shirt",
1678
               "target": "red bicycle",
1679
               "relationship": "more left",
1680
               "task": "position"
1681
1682
           ],
           "modified_questions": [
1683
1684
               "question": "Is the wooden chair farther from the observer than
1685
           the blue table?",
1686
               "anchor": "blue table",
               "target": "wooden chair"
1687
               "relationship": "farther",
1688
               "task": "position"
1689
             },
1690
1691
               "question": "Is the red bicycle located at a lower elevation than
1692
            the man in a floral shirt?",
               "anchor": "man in a floral shirt",
1693
               "target": "red bicycle",
1694
               "relationship": "lower",
1695
               "task": "position"
1696
1698
      messages.append({"role": "user", "content": '\n'.join(query)})
1700
```

#### Listing 12: Prompts for design questions of task Size Comparison

```
1702
      messages = [{"role": "system", "content": f"""
1703
        I will give you an image and a description about the image. You should
1704
          design 2 questions regarding size judgments around top-level objects
1705
          in the image (top-level means it's not a part of another object). The
1706
           question should involve two object (anchor, target) and a type of
1707
          relationship:
1708
        - *The anchor and target object* should be randomly chosen from the top
1709
          -level objects. You possibly need to add the attributions about
1710
          appearance, behavior, posture, position in the image, etc. to the
1711
          description of the anchor and target objects so that they can be
1712
          distinguished from others.
         - *The relationship* should be randomly selected from: larger, smaller,
1713
           taller, shorter, wider, narrower.
1714
        - Only design questions about top-level objects. Ignore those not in
1715
          top-level objects list. Ignore environment objects such as water, sky
1716
          , grass, cloud, etc.
1717
        After that, generate 2 more questions based on the designed questions
1718
          by choose another relationship and keep other parts unchanged.
1719
1720
        Please respond in JSON format. All content should be in English. Here
1721
          is an example of output:
1722
           "questions": [
1723
1724
               "question": "Is the green vase taller than the brown table?",
1725
               "anchor": "brown table",
1726
               "target": "green vase",
               "relationship": "taller",
1727
               "task": "size"
```

```
1728
             },
1729
1730
               "question": "Is the plate with food on it narrower than the white
1731
            box in the middle?",
               "anchor": "white box in the middle",
1732
               "target": "plate with food on it",
1733
               "relationship": "narrower",
1734
               "task": "size"
1735
1736
           ],
           "modified_questions": [
1737
1738
               "question": "Is the green vase smaller than the brown table?",
1739
               "anchor": "brown table",
1740
               "target": "green vase",
               "relationship": "smaller",
1741
               "task": "size"
1742
             },
1743
1744
               "question": "Is the plate with food on it larger than the white
1745
          box in the middle?",
               "anchor": "white box in the middle",
1746
               "target": "plate with food on it",
1747
               "relationship": "larger",
1748
               "task": "size"
1749
1750
1751
       """}]
1752
      messages.append({"role": "user", "content": '\n'.join(query)})
1753
```

Listing 13: Prompts for design questions of task Existence Estimation and Object Counting

```
1755
      messages = [{"role": "system", "content": f"""
1756
        I will give you an image and a description about the image. You should
1757
          design 2 questions regarding existence judgments and 2 questions
1758
          regarding counting around top-level objects in the image (top-level
1759
          means it's not a part of another object). The conditions in the
1760
          question need to involve an anchor object and a type of relationship:
1761
        - *The anchor object* should be randomly chosen from the top-level
1762
          objects. If multiple objects in the image are similar to anchor
1763
          object, you need to add the attributions about appearance, behavior,
1764
          posture, position in the image, etc. to the description of anchor
1765
          object so that it can be distinguished from others.
1766
         - *The relationship* should randomly selected from: more to the left,
          more to the right, closer (to the observer), farther (to the observer
1767
          ), higher, lower, larger, smaller, taller, shorter, wider, narrower.
1768
        - Only design questions about top-level objects. Ignore those not in
1769
          top-level objects list. Ignore environment objects such as water, sky
1770
          , grass, cloud, etc.
1771
        After that, generate 4 more questions based on the designed questions
1772
          by choose another type of relationship and keep other parts unchanged
1773
1774
1775
        Please respond in JSON format. All content should be in English. Here
1776
          is an example of output:
1777
           "questions": [
1778
1779
               "question": "Are there chairs wider than the blue table?",
1780
               "anchor": "blue table",
               "target": "chair",
1781
               "relationship": "wider",
```

```
1782
               "task": "existence"
1783
             },
1784
               "question": "Is there a bicycle located more to the left than the
1785
            man in a floral shirt?",
1786
               "anchor": "man in a floral shirt",
1787
               "target": "bicycle",
1788
                "relationship": "more left",
1789
               "task": "existence"
1790
             },
1791
                "question": "How many green vases are positioned higher than the
1792
           middle wooden table?",
1793
               "anchor": "middle wooden table",
               "target": "green vase",
1794
               "relationship": "higher",
1795
                "task": "count"
1796
             },
1797
1798
               "question": "How many plates are larger than the white box {\operatorname{in}} the
1799
            middle?",
               "anchor": "white box in the middle",
1800
               "target": "plate",
1801
               "relationship": "larger",
1802
               "task": "count"
1803
             }
1804
           "modified_questions": [
1805
1806
               "question": "Are there chairs which have lower elevation than the
1807
            blue table?",
1808
               "anchor": "blue table",
1809
               "target": "chair",
               "relationship": "lower",
1810
                "task": "existence"
1811
             },
1812
1813
               "question": "Is there a bicycle closer to the observer than the
1814
           man in a floral shirt?",
               "anchor": "man in a floral shirt",
1815
               "target": "bicycle",
1816
               "relationship": "closer",
1817
               "task": "existence"
1818
             },
1819
               "question": "How many green vases are wider than the middle
1820
           wooden table?",
1821
               "anchor": "middle wooden table",
1822
               "target": "green vase",
1823
               "relationship": "wider",
1824
               "task": "count"
             },
1825
1826
               "question": "How many plates are shorter than the white box in
1827
           the middle?",
1828
               "anchor": "white box in the middle",
1829
               "target": "plate",
               "relationship": "shorter",
1830
               "task": "count"
1831
1832
           ]
1833
       """}]
1834
       messages.append({"role": "user", "content": '\n'.join(query)})
1835
```

## E TRAINING DETAILS

We train InternVL-Spatial-8B using LoRA(Hu et al., 2022) with approximately 291K general training samples from InternVL2.5 (Chen et al., 2024c) and 2M samples from InternSpatial, counted with repetition. The training is conducted on 16 A100 GPUs for approximately 14 hours. We report the models and training hyperparameters of InternVL-Spatial-8B in Table 8.

Table 8: Training settings and hyperparameters for InternVL-Spatial-8B models. Key configurations for InternVL-Spatial-8B, including model architectures and training parameters.

#### InternVL-Spatial-8B ViT InternViT-300M LLM $Internlm2\_5\text{--}7b\text{--}chat$ Tile Resolution Lora Rank Training Hyperparameters Packed Batch Size Optimizer AdamW Learning Rate 2.00E-05 Warmup Ratio 0.03 LR Scheduler Cosine Weight Decay 0.05 ViT Drop Path 0.1 Image Tile Threshold Context Length 12.8K **Epochs**

# F VISUALIZATION OF INTERNSPATIAL











# Image Source: Visual Genome

\text{duestion:} \\
\text{image>} \text{Comparing \ref>region</ref>} \text{cox>}[367.0, 520.0, \\
441.0, 735.0]</ri>\text{box>} and \\
\text{ref-region</ref>} \text{cox}>(436.0, 490.0, 545.0, \\
734.0]</ri>\text{box>}, which is more behind?

Answer: <ref>region</ref><box>[367.0, 520.0, 441.0, 735.0]</box> is clearly more behind than the



Task type: Position Comparisor

orative black roof on the ref><box>[125.0, 5.0, 977.0, 145.0]</box>



Task type: Size Comparise

inswer: the height of the tapestry in the orange box surpasses the other.



Task type: Size Compari

Instruction Format: Image with numbers+ According to Image

Answer: The width of the island with a circled '1' drawn on it



Task type: Size Compari

Instruction Format: Image with masks+ With <box

Question:
<a href="mailto:question:decoration:geo-color:

Answer:
Sorry, but <ref>region</ref><box>{510.0, 282.0, 983.0, 495.0}</box> is bigger than <ref>region</ref><box>{295.0, 89.0, 475.0,



question: <image> Are there any people closer to the observer than <ref>region</re>/ref><box>[282.0, 331.0, 615.0, 558.0]</box>?

Answer:



Image Source: COCO

Task type: Object Counting



#### Image Source: Cityscapes

Task type: Size Comparison



Answer: the car in the pink box is the larger one.



#### Image Source: Cityscapes

Task type: Positional Comparison

Answer: the car with the red mask is in front of the car with the green mask.

# G VISUALIZATION OF RESULTS ON INTERNSPATIAL-BENCH



Question: <image> Where is <ref>the woman in the red coat</ref> positioned in relation to <ref>the sign</ref>? Answer with the option's letter from the given choices directly. (A) To the right side of <ref>the sign</ref> (B) In front of <ref>the sign</ref> (C) Behind <ref>the sign</ref> (D) To the left side of <ref>the sign</ref>

GT: A

InternVL-8B Pred: (C) Behind the sign

InternVL-Spatial-8B Pred: (A) To the right side of the sign

Task: positional



Question: <image> Are there any vehicles more to the right than <ref>the silver car in the left lane</ref>? Answer with a single word or option's letter.

GT: Yes

InternVL-8B Pred: No

InternVL-Spatial-8B Pred: Yes

Task: existence



Question: <image> Is <ref>region</ref><box>[[599, 390, 654, 472]]</box> positioned farther from the observer than <ref>region</ref><box>[[918, 378, 1000, 538]]</box>? Answer with a single word or or other settler.

GT: Yes

InternVL-8B Pred: No

InternVL-Spatial-8B Pred: Yes

Task: positional





Question: <image> Is <ref>the blue parking sign on the right side of the road</ref> shorter than <ref>the vhite building on the left side</ref>? Answer with a single word or option's letter.

GT: Yes

InternVL-8B Pred: No

InternVL-Spatial-8B Pred: Yes



Question: <image> Which is above, <ref>region</ref><box>[210, 179, 453, 566]</box> or <ref>region</ref><box>[209, 0, 506, 516]</box>?

GT: Positioned higher is <ref>region</ref><box>[209, 0, 506, 516]</box>.

InternVL-8B Pred: region[210, 179, 453, 566] is above region[209, 0, 506, 516].

InternVL-Spatial-8B Pred: region[209, 0, 506, 516] is the one that is

Task: above\_choice



Question: <image> Is the picture with the purple mask below the pillow with the pink mask?

 $\operatorname{GT:}$  Actually, the picture with the purple mask is above the pillow with the pink mask.

InternVL-8B Pred: Yes, the picture with the purple mask is below the pillow with the pink mask.

InternVL-Spatial-8B Pred: No, the picture with the purple mask is not situated beneath the pillow with the pink mask.

Task: below\_predicate



Question: <image> Does the counter with a circled '1' drawn on it come behind the picture with a circled '2' drawn on it?

GT: No, it is not.

InternVL-8B Pred: Yes, the counter with a circled '1' is behind the picture with a circled '2'.

InternVL-Spatial-8B Pred: No, the counter with a circled '1' drawn on i is not behind the picture with a circled '2' drawn on it.

Task: behind\_predicate





Question: <image> Does the books with the yellow mask come in front of the chair with the blue mask?

GT: No, it is not.

InternVL-8B Pred: Yes, the books with the yellow mask are placed in front of the chair with the blue mask.

InternVL-Spatial-8B Pred: No, the books with the yellow mask is not ir front of the chair with the blue mask.

Task: front\_predicate



Question: <image> Which is more to the left, the coffee maker in the red box or the box in the yellow box?

GT: From the viewer's perspective, the box in the yellow box appears more on the left side.

InternVL-8B Pred: The coffee maker in the red box is more to the left than the box in the yellow box.

InternVL-Spatial-8B Pred: the box in the yellow box is positioned further to the left.

Task: left\_choice



Question: <image> Can you confirm if the pillow with a circled '1' drawn on it is positioned to the right of the books with a circled '2' drawn on it.

GT. No, the pillow with a circled '1' drawn on it is not to the right of the books with a circled '2' drawn on it.

InternVL-8B Pred: Yes, the pillow with a circled '1' is positioned to the right of the books with a circled '2'.

InternVL-Spatial-8B Pred: the pillow with a circled '1' drawn on it is no to the right of the books with a circled '2' drawn on it, actually.

Task: right predicate



Question: <image> Does the toaster with a circled '1' drawn on it have a larger size compared to the counter with a circled '2' drawn on it?

GT: Incorrect, the toaster with a circled '1' drawn on it is not larger than the counter with a circled '2' drawn on it.

InternVL-8B Pred: Yes, the toaster with a circled '1' drawn on it is larger in size compared to the counter with a circled '2' drawn on it.

InternVL-Spatial-8B Pred: Not according to the information, the toaster with a circled '1' drawn on it is not larger than the counter with a circled '2' drawn on it.

Task: big\_predicate



Question: <image> Does the barrier in the purple box have a greater height compared to the truck in the pink box?

GT: In fact, the barrier in the purple box may be shorter than the truck in the pink box.

InternVL-8B Pred: Yes, the barrier in the purple box appears to have a greater height compared to the truck in the pink box.

InternVL-Spatial-8B Pred: Sorry, but the barrier in the purple box does not exceed the truck in the pink box in height.

Task: tall\_predicate



Question: <image> Can you confirm if the painting in the purple box is thinner than the potted plant in the blue box?

 $\mbox{GT:}$  In fact, the painting in the purple box might be wider than the potter plant in the blue box.

InternVL-8B Pred: Yes, the painting in the purple box is thinner than the potted plant in the blue box.

InternVL-Spatial-8B Pred: No, the painting in the purple box is not narrower than the potted plant in the blue box.

Task: thin\_predicate

# H THE USE OF LARGE LANGUAGE MODELS (LLMS)

We used large language models (LLMs) as assistive tools during the preparation of this work. Specifically, LLMs were employed for language polishing, LaTeX code editing, and debugging of prompts in the dataset construction process. The authors take full responsibility for the content of the paper.