In this supplementary material, we present additional details and clarifications that are omitted in the main text due to space constraints.

- Appendix A Related Works.
- Appendix B Dataset Details.
- Appendix C Implementation Details.
- Appendix D Full Results.
- Appendix E More Qualitative Examples.
- 550 551 552

553

543

544

546

547

548 549

A RELATED WORKS

554 VLMs for Robotics. Vision-language models (VLMs) have emerged as pivotal tools in robotics, 555 enabling systems to interpret and act upon complex visual and textual information. By integrating 556 visual perception with language understanding, VLMs facilitate more intuitive human-robot interactions and enhance autonomous decision-making capabilities. Recent advancements have demon-558 strated the potential of VLMs in various robotic applications. For instance, vision-language-action 559 models (VLAs) Kim et al. (2024); Brohan et al. (2023); Octo Model Team et al. (2024) enable robots to interpret and execute complex instructions and output executable robot actions. Addition-561 ally, VLMs like GPT-4v OpenAI et al. (2024) have been utilized for high-level task planning Wake 562 et al. (2023), allowing robots to generate detailed action sequences from natural language instructions. Furthermore, VLMs have been used for keypoint/mask prediction Huang et al. (2024c); Wi 563 et al. (2023); Nasiriany et al. (2024), error analysis Duan et al. (2024), grasp pose prediction Huang et al. (2024a). Despite these advancements, integrating VLMs Cai et al. (2024); Cheng et al. (2024); 565 Yuan et al. (2024) into robotic systems presents challenges. One significant hurdle is the need for 566 precise spatial reasoning to navigate and manipulate objects effectively. While VLMs excel in un-567 derstanding and generating language, their ability to comprehend and reason about spatial relation-568 ships in dynamic environments remains limited Yamada et al. (2024); Xu et al. (2024); Wang et al. 569 (2024a). Therefore, ROBOSPATIAL aims to tackle this gap by presenting a large scale pretraining 570 and evaluation setup for teaching spatial understanding to VLM for robotics. 571

Spatial Understanding with VLMs. Spatial understanding has been implicitly and explicitly part 572 of various vision and question answering tasks Fu et al. (2024); Azuma et al. (2022); Jia et al. 573 (2024); Suhr & Artzi (2019); Salewski et al. (2022); Krishna et al. (2017); Johnson et al. (2017); 574 Hudson & Manning (2019). While many benchmarks and methods have been proposed, they often 575 come with limitations: some focus exclusively on simulations Szymanska et al. (2024) or generic 576 images Liu et al. (2023a); Rajabi & Kosecka (2024); Cheng et al. (2024); Chen et al. (2024); Fu 577 et al. (2024); Kamath et al. (2023); Shiri et al. (2024); Ranasinghe et al. (2024), others are difficult 578 to evaluate Szymanska et al. (2024); Du et al. (2024); Linghu et al. (2024), rely on complete 3D scans Zhang et al. (2024); Man et al. (2024); Ma et al. (2023); Linghu et al. (2024), or do not consider 579 reference frames Zhang et al. (2024); Man et al. (2024); Ma et al. (2023); Linghu et al. (2024); Chen 580 et al. (2024); Cheng et al. (2024); Fu et al. (2024); Ranasinghe et al. (2024). Furthermore, they 581 often fail to address actionable, robotics-relevant spatial relationships such as spatial compatibility 582 and context Du et al. (2024); Fu et al. (2024); Wang et al. (2024b); Shiri et al. (2024); Kamath et al. 583 (2023); Linghu et al. (2024); Ranasinghe et al. (2024). 584

- Inspired by prior works on spatial reasoning Liu et al. (2023a); Kamath et al. (2023)—where the impact of reference frames and spatial configurations was explored in generic images Lin et al. (2015); Hudson & Manning (2019)—we extend spatial understanding to a robotics-specific context with actionable spatial relationships such as spatial compatibility and spatial context. Our aim is to enable direct application to robotic workflows, such as task planning and verification.
- To achieve this, we have developed and are planning to open-source a large-scale 2D/3D ready
 pretraining dataset, an automated data annotation pipeline, and trained models. We further show how
 our dataset can be used to teach spatial reasoning to a suite of vision-language models (VLMs) in
 in-domain and out-of-domain spatial reasoning datasets. We hope these resources lower the barrier
 to entry for exploring spatial understanding tailored to robotics.



Figure 3: Overview of the ROBOSPATIAL dataset. We automatically generate spatial relationship annotations from existing datasets with 3D point clouds, egocentric images, and 3D bounding box annotations. We create question/answer pairs covering three classes of spatial relationships, three spatial reference frames, and both binary (yes/no) and numeric (e.g., 2D image points) answers. From 1M images and 5K scans, we generate over 3M spatial question/answer pairs.

B DATASET DETAILS

B.1 DATASET STATISTICS

We provide the full dataset statistics in Table 3 For all training, we use only 900,000 spatial relationships, equally divided across all datasets, due to computational constraints. We further experiment on the effect of data scaling on Table 7 and explain the results. Notably, HOPE Tyree et al. (2022) and GraspNet-1B Fang et al. (2020) contain similar tabletop images captured from different perspectives, resulting in lower dataset diversity for the tabletop environment. We plan to enhance the diversity of our dataset by incorporating additional tabletop datasets.

B.2 DATASET GENERATION DETAILS

Frame Generation. We explain how answers are generated for each frame Figure 3 Each question type has three answers from ego-centric, object-centric, and world-centric perspectives. Each frame may share the same answer since not all frames lead to unique answers. For our answer generation, we used the following logic:

- Ego-centric: The default perspective from the ego-centric view.
- **Object-centric:** Using the oriented bounding box directions, we determine the front of the object. With this information, we assign front, behind, left, and right. Above and below remain the same as in the ego-centric perspective.
- World-centric: Using the z-coordinate of the oriented bounding box, we modify the above and below relationships to reflect whether an object is above or below another object with respect to elevation.

645 Compatibility Generation. For compatibility, we construct a top-down map as shown in Figure 4.
646 Using the top-down map and the top-down 2D bounding box of the object to be placed, we determine
647 if there exists an empty space within a threshold distance (i.e., 1% of the object's longest width or length). If the object meets this condition, it is deemed compatible to be fitted.

648

662 663

677 678 679

680

682

Category	Dataset	Split	Scans	Images	Configuration Q	Context Q	Compatibility Q
		Train	1859 scans	236243	298439	298439	298439
	Matterport3D Chang et al. (2017)	Validation	5 scans	100	100	100	100
		Test	J Scalls	100	200420	200.420	200.420
Indoor	ScanNet Dai et al (2017)	I rain Validation	12 scans	278402	298439	298439	298439
	Scanter Dar et al. 2017	Test	12 scans	1200	500	500	500
		Train	1543 scans	365355	298439	298439	298439
	3RScan Wald et al. (2019)	Validation	9 scans	900	400	400	400
		Test	9 scans	900	400	400	400
		Train	50 scenes	50000	36317	36317	36317
	HOPE Tyree et al. (2022)	Validation	10 scenes	50	500	500	500
Tabletop		Test	47 scenes	235	500	500	500
		Train	100 scenes	25500	36317	36317	36317
	GraspNet-1B Fang et al. (2020)	Validation	30 scenes	120	500	500	500
		Test	30 scenes	120	500	500	500

Table 3: Full dataset statistics for indoor and tabletop datasets.



Figure 4: An example of generated top-down map of the image from 3D bounding boxes.

C IMPLEMENTATION DETAILS

681 C.1 MODEL TRAINING

We further explain the training details for all 2D and 3D VLMs trained on ROBOSPATIAL. For all models, we perform instruction tuning using the model weights from public repositories. All training is done using 8 Nvidia H100 GPUs, with the training time between 20 and 40 hours.

VILA Lin et al. (2023) We initialize our model from Efficient-Large-Model/Llama-3-VILA1.5-8B
on Hugging Face. We use the fine-tuning script from the VILA GitHub repository to train our model
using the default hyperparameters.

- LLaVA-NeXT Liu et al. (2023b) We initialize our model from lmms-lab/llama3-llava-next-8b on
 Hugging Face. We use the LLaVA-Next fine-tuning script from the LLaVA-Next repository using
 the default hyperparameters.
- SpaceLLaVA Chen et al. (2024) As official code and weights for SpatialVLM Chen et al. (2024) is not released, we use a community implementation which is endorsed by SpatialVLM Chen et al. (2024) authors. We initialize our model from remyxai/SpaceLLaVA from Hugging Face. We use LLaVA-1.5 finetuning script from LLaVa Liu et al. (2023c) repository using the default hyperparameters.
- RoboPoint Yuan et al. (2024) We initialize our model from wentao-yuan/robopoint-v1-vicuna-v1.513b on Hugging Face. We use the fine-tuning script provided in the RoboPoint Yuan et al. (2024)
 GitHub repository to train our model using the default hyperparameters.
- **3D-LLM** Hong et al. (2023) We initialize our model using the pretrain_blip2_sam_flant5xl_v2.pth checkpoint downloaded from the official GitHub repository. Since the model requires preprocessing of multiview images, we follow the author's pipeline to process multiview images from our environments. Because the model does not accept image input, we append the following text in front of the

Model	Indoor			Tabletop			Average		
litituti	Configuration	Context	Compatibility	Configuration	Context	Compatibility	Indoor	Tabletop	Total
			Open-sourc	e VLMs					
VILA Lin et al. (2023) +ROBOSPATIAL	54.7 71.4 ↑	18.3 45.9 ↑	2D VL 56.3 77.2 ↑	Ms 45.1 71.8 ↑	13.2 43.7 ↑	53.8 73.3↑	43.1 64.8 ↑	37.4 62.9 ↑	40.2 63.9↑
LLaVA-NeXT Liu et al. (2023b) +ROBOSPATIAL	48.9 69.3 ↑	12.5 41.3 ↑	32.7 70.5 ↑	48.3 70.7 ↑	8.4 44.8 ↑	30.9 66.1 ↑	31.4 60.4 ↑	29.2 60.5 ↑	30.3 60.5 †
SpaceLLaVA Chen et al. (2024) +ROBOSPATIAL	52.6 76.0 ↑	15.3 50.7 ↑	49.0 76.6 ↑	66.5 74.9 ↑	12.2 46.4 ↑	60.1 70.5 ↑	38.9 67.8 ↑	46.2 63.6 ↑	43.6 65.7 ↑
RoboPoint Yuan et al. (2024) +ROBOSPATIAL	39.0 72.2 ↑	41.4 68.9 ↑	38.3 72.1 ↑	37.9 70.3 ↑	31.6 61.7 ↑	45.2 78.4 ↑	39.6 71.0 ↑	38.2 70.1 ↑	38.9 70.6↑
			3D VL	Ms					
3D-LLM Hong et al. (2023) +ROBOSPATIAL	54.5 76.3 ↑	8.1 35.4 ↑	53.6 77.5 ↑	59.2 76.2 ↑	10.6 46.8 ↑	57.4 75.0↑	37.6 63.1↑	42.4 66.0↑	40.0 64.6 ↑
LEO Huang et al. (2024b) +ROBOSPATIAL	56.1 80.2 ↑	11.3 56.7 ↑	58.3 82.5 ↑	60.8 78.1 ↑	11.1 55.2 ↑	59.3 78.9 ↑	41.9 73.1 ↑	43.7 70.7 ↑	42.8 71.9 ↑
			Not available fo 2D VL	r fine-tuning Ms					
Molmo Deitke et al. (2024) GPT-40 Open AL et al. (2024)	40.6	48.2	60.0 59.4	61.5 62.3	35.8 27.9	54.6 66.8	49.6 49.3	50.6 52.3	50.1 50.8

702 Table 4: Results of existing 2D/3D VLMs on a held-out test split of images and scans. All methods, for all 703 tasks, perform better ([†]) when fine-tuned on our ROBOSPATIAL dataset. The best result for each column is bolded. 704

722 question to ensure the model understands the perspective from which the question is being asked: "I 723 am facing ANCHOR OBJECT." We use the default hyperparameters and train the model for 20 epochs 724 per the author's guidelines. We choose the best model based on validation accuracy.

725 LEO Huang et al. (2024b) We initialize our model from the sft_noact.pth checkpoint downloaded from the official GitHub repository. 726

Since LEO supports dual image and 3D point cloud input, we input both of them and modify the 727 question as in 3D-LLM. We use the default hyperparameters and train the model for 10 epochs per 728 the author's guidelines, and choose the best model based on validation accuracy. 729

730 We could not fine-tune Molmo Deitke et al. (2024) from allenai/Molmo-7B-D-0924 or GPT-40 Ope-731 nAI et al. (2024) from the gpt-40-2024-08-06 API due to the unavailability of the fine-tuning script 732 at the time of this work, thus we use them as a zero-shot baselines.

733 734

735

743

744 745

746

747

748

C.2 ROBOT SETUP

736 For picking, we find which object the point maps to using SAM 2 Ravi et al. (2024) and execute 737 our picking behavior on that object. For placing, we simply compute the 3D coordinate based on 738 the depth value at that pixel and place the object at that coordinate. There were no failures due to 739 cuRobo Sundaralingam et al. (2023) failing. The experiments were purposely designed to consist of behaviors that our robot system can handle in order to avoid introducing irrelevant factors. The 740 picking behavior consists of computing a top-down grasp pose and reaching it with cuRobo Sundar-741 alingam et al. (2023). To compute the grasp pose: 742

- 1. We estimate the major axis of the object's point cloud in top-down view using PCA.
- 2. The grasp orientation is orthogonal to the major axis.
- 3. The grasp height is based on the highest point in the object's point cloud minus an offset of 3cm. This heuristic ensures the system can grip long objects.

749 The placing behavior is the same as picking, except that an area within 5cm of the placement coor-750 dinate is used as the point cloud for estimating orientation and height, and a vertical height offset is 751 added to account for the height at which the object was picked. 752

753 754

755

C.3 **OMITTED RESULTS IN THE MAIN TEXT**

We show the full results in held-out test split in Table 4 and out-of-domain splits in Table 5.

⁷²¹

Model	RO	BOSPATIAL-	Home	BLINK-Spatial	
Widder	Localization	Affordance	Compatibility	Accuracy	
	Open-sou	rce VLMs			
	2D V	LMs			
VILA Lin et al. (2023)	53.3	12.0	52.0	72.7	
+ROBOSPATIAL	$62.0\uparrow$	32.0 ↑	58.0 ↑	79.7 ↑	
LLaVA-NeXT Liu et al. (2023	3b) 48.0	9.3	37.3	71.3	
+ROBOSPATIAL	58.0 ↑	$24.0\uparrow$	44.0 ↑	79.0 ↑	
SpaceLLaVA Chen et al. (202-	4) 60.0	16.0	49.3	76.2	
+ROBOSPATIAL	68.7 ↑	38.0 ↑	56.0 ↑	81.8 ↑	
RoboPoint Yuan et al. (2024)	43.3	41.3	36.0	63.6	
+ROBOSPATIAL	50.0 ↑	54.0 ↑	48.0 ↑	$70.6\uparrow$	
	3D V	LMs	16.0	NT/A	
3D-LLM Hong et al. (2023)	40.0 48.0 *	8.0	46.0 52.7.↑	N/A N/A	
+ROBOSPATIAL	48.0	30.0	32.7	IN/A	
LEO Huang et al. (2024b)	50.7	10.0	48.0	N/A	
+ROBOSPATIAL	$64.0\uparrow$	$40.0\uparrow$	60.0 ↑	N/A	
	Not available	or fine-tuning	2		
Molmo Deitke et al. (2024)	44.7	38.0	58.0	67.1	
GPT-40 OpenAI et al. (2024)	64.0	25.3	56.7	76.2	

Table 5: Results on an out-of-domain test split comparing prior art VLMs. The results show improved ([†]) spatial understanding capabilities on similar domains. Bolded number is the best result for the column.

Table 6: Average accuracy for dataset generalization when training on indoor scenes and testing on tabletop scenes (indoor \rightarrow tabletop), and vice versa (tabletop \rightarrow indoor).

	$\textbf{Indoor} \rightarrow \textbf{Tabletop}$	$\textbf{Tabletop} \rightarrow \textbf{Indoor}$
RoboPoint Yuan et al. (2024)	38.7	38.2
+ROBOSPATIAL	48.9 ↑	51.3 ↑
LEO Huang et al. (2024b)	41.9	43.7
+ROBOSPATIAL	47.2 ↑	54.5 ↑

C.4 CROSS-DATASET GENERALIZATION

We evaluate the generalization capability of our method by testing it across different scene types—specifically, both indoor and tabletop scenes—to control for any bias in the annotations of the underlying datasets that make up our benchmark. We train on data derived from subsets of the datasets corresponding to one scene type (either indoor or tabletop) and test on held-out datasets from the other scene type, representing unseen environments. We expect that even when training on a subset of datasets, the performance on unseen scene types will improve if our method generalizes well. The results of this cross-dataset evaluation are shown in Table 6

C.5 DATA SCALING

In Table 7, we experiment with scaling the number of annotations while keeping images fixed. We found that even though the number of images stays consistent, increasing the number of annotations can improve performance. For future work, we plan to apply our data generation pipeline to a diverse set of indoor and tabletop environments to further improve the performance of our models.

C.6 ACCURACY PER FRAME OF REFERENCE

We show the results per frame in Table 8 for our out-of-domain test set. From the results, we can see a distinct difference between 2D and 3D VLMs in understanding the world-centric frame before training with ROBOSPATIAL. Baseline 2D VLMs have trouble understanding the worldTable 7: Results of scaling experiment on LLaVa-Next Liu et al. (2023b) with varied spatial relationship annotations. Average accuracy on held-out test set is reported.

Annotation Size	100K	300K	900k (Default)	1.8M	3M (Full)
LLaVa-Next Liu et al. (2023b)	38.1	46.7	60.5	65.8	72.4

Table 8: Results of per frame accuracy of existing 2D/3D VLMs on a held-out test split of images and scans.
All methods, for all tasks, perform better ([†]) when fine-tuned on our ROBOSPATIAL dataset. The best result for each column is bolded.

Model		Indoor		Tabletop			Average		
Model	Ego-centric	Object-centric	World-centric	Ego-centric	Object-centric	World-centric	Indoor	Tabletop	Total
			Open-sourc	e VLMs					
		10 .	2D VL	Ms					
VILA Lin et al. (2023)	55.9	40.5	32.9	43.6	39.7	28.9	43.1	37.4	40.2
+ROBOSPATIAL	74.3↑	57.8↑	62.3 ↑	70.3↑	58.1↑	60.3 ↑	64.8 ↑	62.9↑	63.9↑
LLaVA-Next Liu et al. (2023b)	35.2	24.3	34.7	36.4	28.5	22.7	31.4	29.2	30.3
+ROBOSPATIAL	75.4 ↑	54.1 ↑	68.8 ↑	67.9 ↑	54.7 ↑	58.9 ↑	60.4 ↑	$60.5\uparrow$	$60.5\uparrow$
SpaceLLaVA Chen et al. (2024)	40.6	36.0	30.1	52.3	32.8	53.5	38.9	46.2	43.6
+ROBOSPATIAL	78.5 ↑	60.6 ↑	64.3 ↑	73.0 ↑	49.5 ↑	68.3 ↑	67.8 ↑	63.6 ↑	$65.7\uparrow$
RoboPoint Yuan et al. (2024)	41.9	36.2	40.7	46.2	30.5	37.9	39.6	38.2	38.9
+ROBOSPATIAL	76.4 ↑	58.3 ↑	78.3 ↑	76.7 ↑	62.6 ↑	$71.0\uparrow$	$71.0\uparrow$	$70.1\uparrow$	70.6 ↑
			3D VL	Ms					
3D-LLM Hong et al. (2023)	28.9	38.3	45.6	38.9	35.7	52.6	37.6	42.4	40.0
+ROBOSPATIAL	60.7 ↑	52.1 ↑	76.5 ↑	57.9 ↑	62.8 ↑	77.3 ↑	63.1 ↑	66.0 ↑	$64.6\uparrow$
LEO Huang et al. (2024b)	46.9	30.6	48.2	41.4	34.3	55.4	41.9	43.7	42.8
+ROBOSPATIAL	$68.1\uparrow$	71.6 ↑	79.6 ↑	71.4 ↑	$60.2\uparrow$	80.5 ↑	73.1 ↑	70.7 ↑	71.9 ↑
			Not available fo 2D VL	r fine-tuning Ms					
Molmo Deitke et al. (2024)	50.4	50.8	47.6	64.4	33.6	53.8	49.6	50.6	50.1
GPT-40 OpenAI et al. (2024)	52.9	38.7	56.3	62.5	30.7	63.7	49.3	52.3	50.8

> centric frame, which involves understanding elevation, while 3D VLMs comparatively excel at it. Furthermore, we can see that since baseline 3D VLMs are trained on point clouds without information of perspective, their accuracy in ego-centric and object-centric frames is lower. However, with ROBOSPATIAL training, we were able to teach the 3D VLMs to think in a certain frame, thus considerably improving their performance on ego-centric and object-centric frames. However, we hypothesize that, due to their design—specifically, the lack of a means to visually inject perspective information since they require complete 3D point clouds—3D VLMs still lag behind 2D VLMs on ego-centric and object-centric frames.

D FULL RESULTS

D.1 ROBOT EXPERIMENTS

We present additional results from our robot experiments in Figure 5 and Figure 6. We observe that models trained with ROBOSPATIAL consistently outperform baseline models in most cases, even though the prompt is not optimized for ROBOSPATIAL-trained models. This demonstrates that the power of VLMs enables templated language to generalize to language unseen during training while maintaining spatial understanding capabilities. However, even with ROBOSPATIAL training, the models struggle with understanding stacked items, indicating a need for further data augmen-tation with diverse layouts. In a few cases, ROBOSPATIAL training adversely affects performance, especially with RoboPoint Yuan et al. (2024). We hypothesize that mixing the dataset with Robo-Point training data and ROBOSPATIAL training data may lead to unforeseen side effects, particularly in grounding objects. Nevertheless, we demonstrate that ROBOSPATIAL training enhances VLM's spatial understanding in real-life robotics experiments, even with freeform language.

E MORE QUALITATIVE EXAMPLES

Figure 7 and Figure 8 present additional qualitative comparisons between models trained on RO BOSPATIAL. Our findings demonstrate that models trained on ROBOSPATIAL consistently exhibit spatial understanding in the challenging ROBOSPATIAL-Home dataset, even outperforming closed



lava gpt-40 molmo robopoint spatial_robopoint spatint spatial_robopoint spatial_robopoint spatial_robo

Task: Place the object in a free space in front of the orange juice box.



Figure 5: Robotics experiments: the red dot shows the model output (if not present, the model failed to provide a valid point in the image); green dots are used to show when a model outputs multiple points. The robot motion generator, cuRobo Sundaralingam et al. (2023), is used to grasp the item referenced by the generated point. The *spatial*- prefix indicates model trained with ROBOSPATIAL.

models like GPT-40 OpenAI et al. (2024). However, we observed that object grounding is a crucial prerequisite for spatial understanding; the improvement is often hindered by the model's inability to ground objects in cluttered scenes, where GPT-40 performs more effectively. Additionally, in Figure 8 we show that the ROBOSPATIAL-trained model successfully generalizes to unseen spatial relationships in Blink-Spatial Fu et al. (2024), including those involving distance, such as "touching."

918				Question. Is there seem to slot the new
919	2	Question: pick lone object		cake mix in the middle of the row of
920		GPT-40 OpenAL et al. (2024)		boxes
921		Molmo Deitke et al. (2024)		GPT-40 OpenAI et al. (2024) √
922		LLaVa-Next Liu et al. (2023b) >		Molmo Deitke et al. (2024)
023		RoboPoint Yuan et al. (2023)	c contraction of the second	S-LLaVa-Next Liu et al. (2023b)
923		S-RoboPoint Yuan et al. (2024) √		RoboPoint Yuan et al. (2024) ×
924				S-RoboPoint Yuan et al. (2024) \checkmark
925			_	
920	1	Question: Is there space in the white container for the orange juice box		Question: alphabet soup fit in the pur-
927		LI aVa-Next Liu et al (2023b)		ple box
928		S-LLaVa-Next Liu et al. (2023b)		LLaVa-Next Liu et al. (2023b) \checkmark
929		RoboPoint Yuan et al. (2024) >		RoboPoint Yuan et al. (2023) \checkmark
930		Molmo Deitke et al. (2024)		S-RoboPoint Yuan et al. (2024)
931		GPT-40 <mark>OpenAI et al. (2024</mark>) √		GPT-40 OpenAI et al. (2024) ×
932				
933		Question: pick object behind the mid-		Ouestion: pick shortest object
934		dle container		LL aVa-Next Lin et al. (2023b)
935		LLaVa-Next Liu et al. (2023b)		S-LLaVa-Next Liu et al. (2023b)
936		RoboPoint Yuan et al. (2024)		RoboPoint Yuan et al. (2024)
937	PORCEN	S-RoboPoint Yuan et al. (2024)		Molmo Deitke et al. (2024) \checkmark
938		GPT-40 OpenAI et al. (2024)		GPT-40 OpenAI et al. (2024) √
939			_	
940		Ouestion: place object in container be-		Question: place the object inside the
941		hind popcorn		smallest box
942		LLaVa-Next Liu et al. (2023b)		LLaVa-Next Liu et al. (2023b) ×
943		S-LLaVa-Next Liu et al. (2023b)		RoboPoint Yuan et al. (2024) \checkmark
944		S-RoboPoint Yuan et al. (2024)		S-RoboPoint Yuan et al. (2024)
945		Molmo Deitke et al. (2024)		GPT-40 OpenAI et al. (2024) ×
946		GP1-40 OpenAI et al. (2024)	<u> </u>	
947		Oractions and the actual discretion with	<u> </u>	
948		the red orange peaches can without dis-		Question: is there an object that is not
949		turbing other objects?		in a stack?
950		LLaVa-Next Liu et al. (2023b) √		GPT-40 OpenAI et al. (2024)
951		S-LLaVa-Next Liu et al. (2023b) RoboPoint Yuan et al. (2024)		Molmo Deitke et al. (2024) \checkmark
952		S-RoboPoint Yuan et al. (2024)		S-LLaVa-Next Liu et al. (2023b) ✓
053		Molmo Deitke et al. (2024)		RoboPoint Yuan et al. (2024)
957			_	
055		Question: can the macaroni and cheese		Question: is there space to place one of
056		be placed on top of cheez-it without		the cans on the cheez-it box?
057		touching other objects?		GPT-40 OpenAI et al. (2024) ×
957		LLaVa-Next Liu et al. (2023b)		Molmo Deitke et al. (2024) ×
958		S-LLaVa-Next Liu et al. (2023b) > RoboPoint Yuan et al. (2024)		LLaVa-Next Liu et al. (2023b) × S-LLaVa-Next Liu et al. (2023b) ×
959		S-RoboPoint Yuan et al. (2024)		RoboPoint Yuan et al. (2024) ×
960		Molmo Deitke et al. (2024)		S-RoboPoint Yuan et al. (2024) ×
961				
962		Question: place on the chiest to the left		Question: pick the highest object on the stack of two objects
963		of macaroni and cheese		
964		GPT-40 OpenAI et al. (2024)		Molmo Deitke et al. (2024) ×
965		Molmo Deitke et al. (2024)		LLaVa-Next Liu et al. (2023b) ×
966		LLaVa-Next Liu et al. (2023b)		S-LLaVa-Next Liu et al. (2023b) × RoboPoint Yuan et al. (2024) ×
967		RoboPoint Yuan et al. (2024)		S-RoboPoint Yuan et al. (2024) ×
968		S-RoboPoint Yuan et al. (2024) √		
969				

Figure 6: Additional robot experiments. A green check mark indicates that the model answered correctly. The S- prefix denotes a model trained with ROBOSPATIAL. The questions are purposely not cleaned to reflect realistic language inputs.

970 971





Question: Is the dining table touching the donut?

Answer: Yes VILA: No S-VILA: Yes LLaVa: Yes SpaceLLaVa: Yes SpaceLLaVa: Yes RoboPoint: No S-RoboPoint: Yes Molmo: No GPT-4o: No



Question: Is the couch under the suitcase?

Answer: Yes VILA: No S-VILA: Yes LLaVa: No S-LLaVa: Yes SpaceLLaVa: No S-SpaceLLaVa: Yes RoboPoint: No S-RoboPoint: Yes Molmo: No GPT-40: Yes

Figure 8: Qualitative results on Blink-Spatial Fu et al. (2024). ROBOSPATIAL-trained model can generalize to unseen spatial relationships.

1023 1024 1025

1022

1013