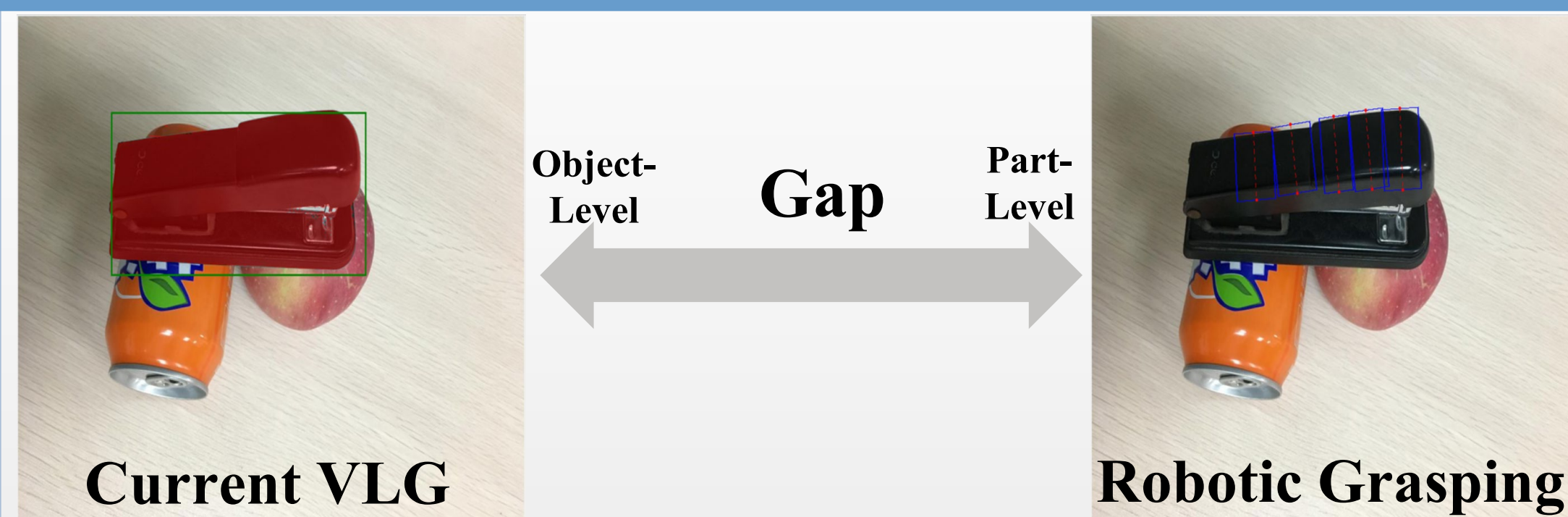
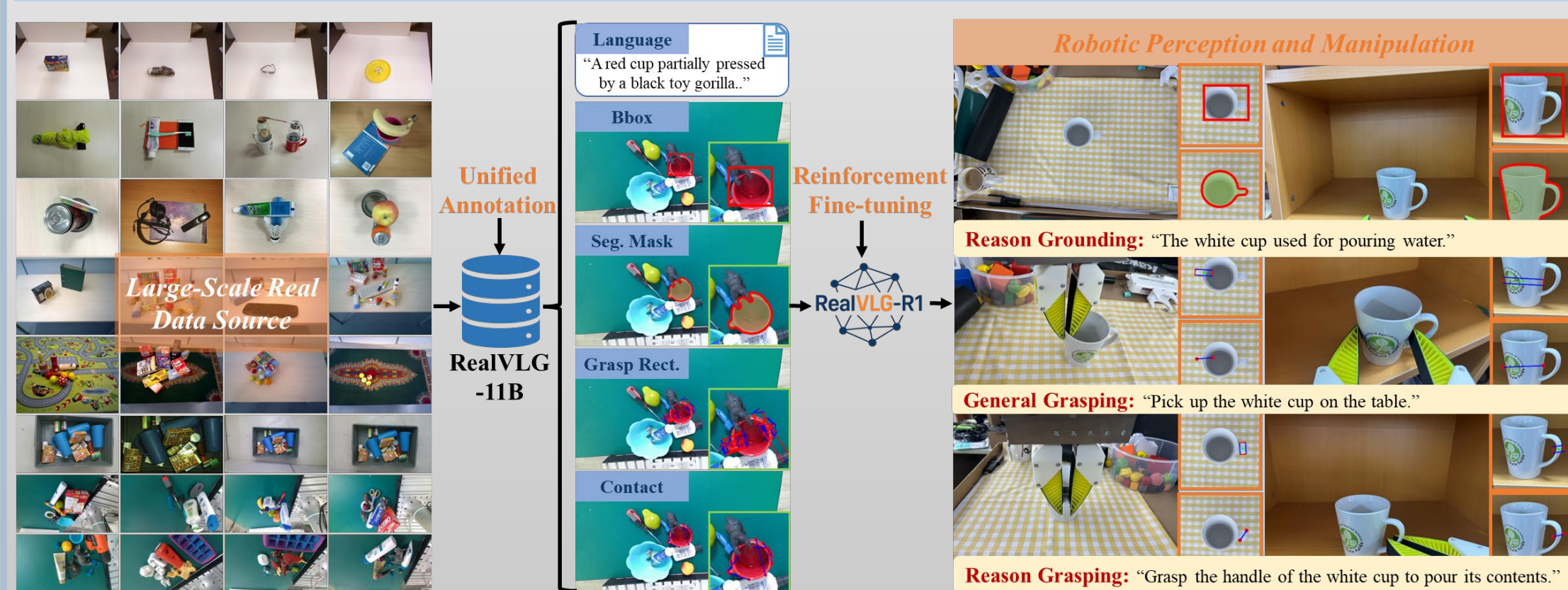


1. Motivation & Problem

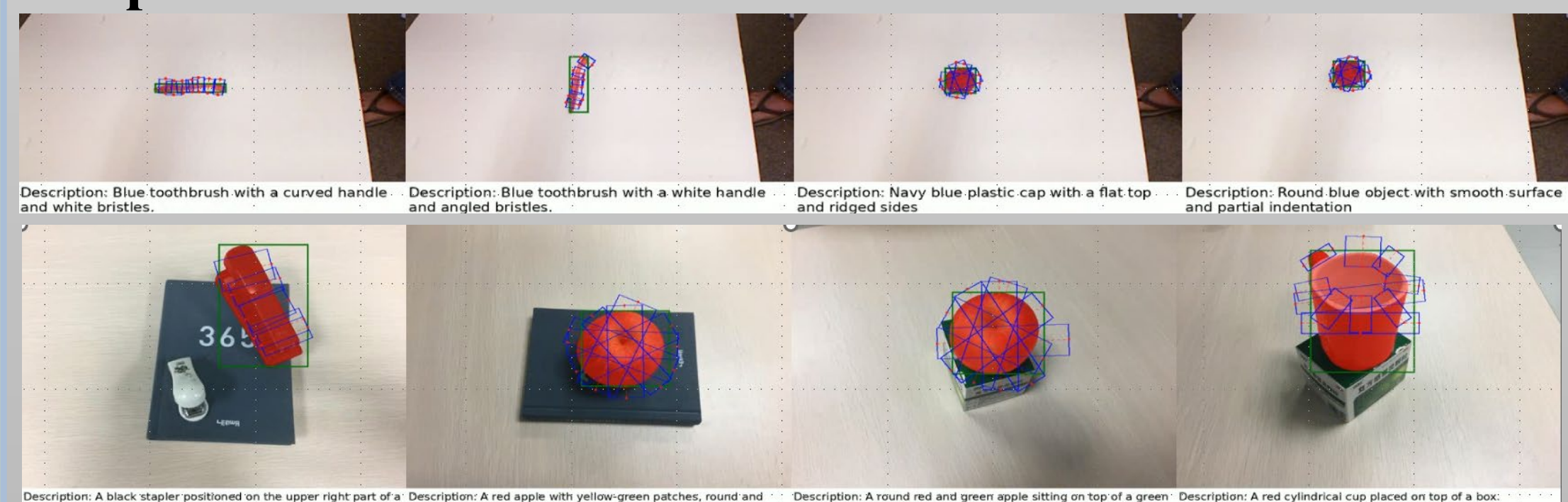


Gap: Coarse VLG vs. Geometric Grasping
Goal: Semantic-to-Manipulation Reasoning

Framework Overview



Samples:



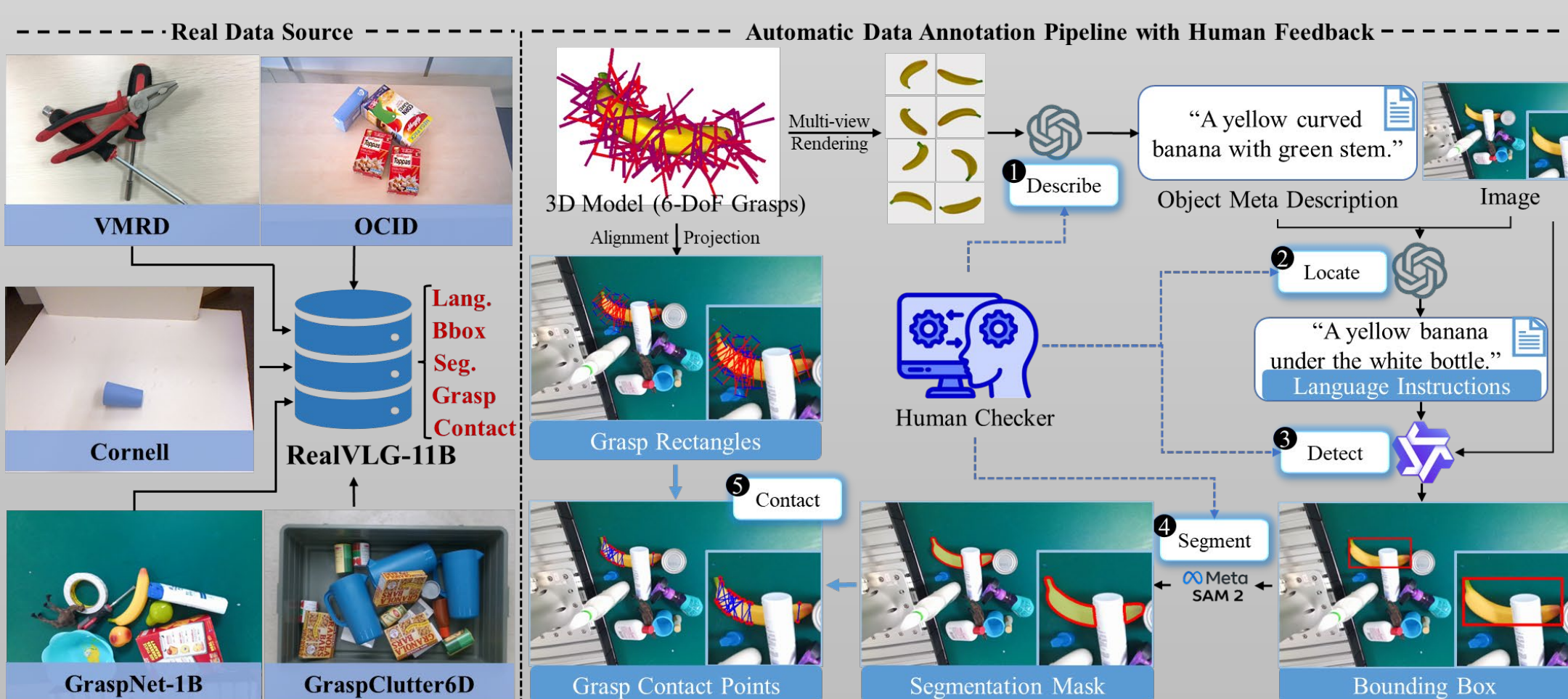
2. RealVLG-11B Dataset

DATASET STATS

165K Images	1~20 Obj/Image	15+ Grasp/Obj
1.3M Annotations	800+ Instances	11B Grasps

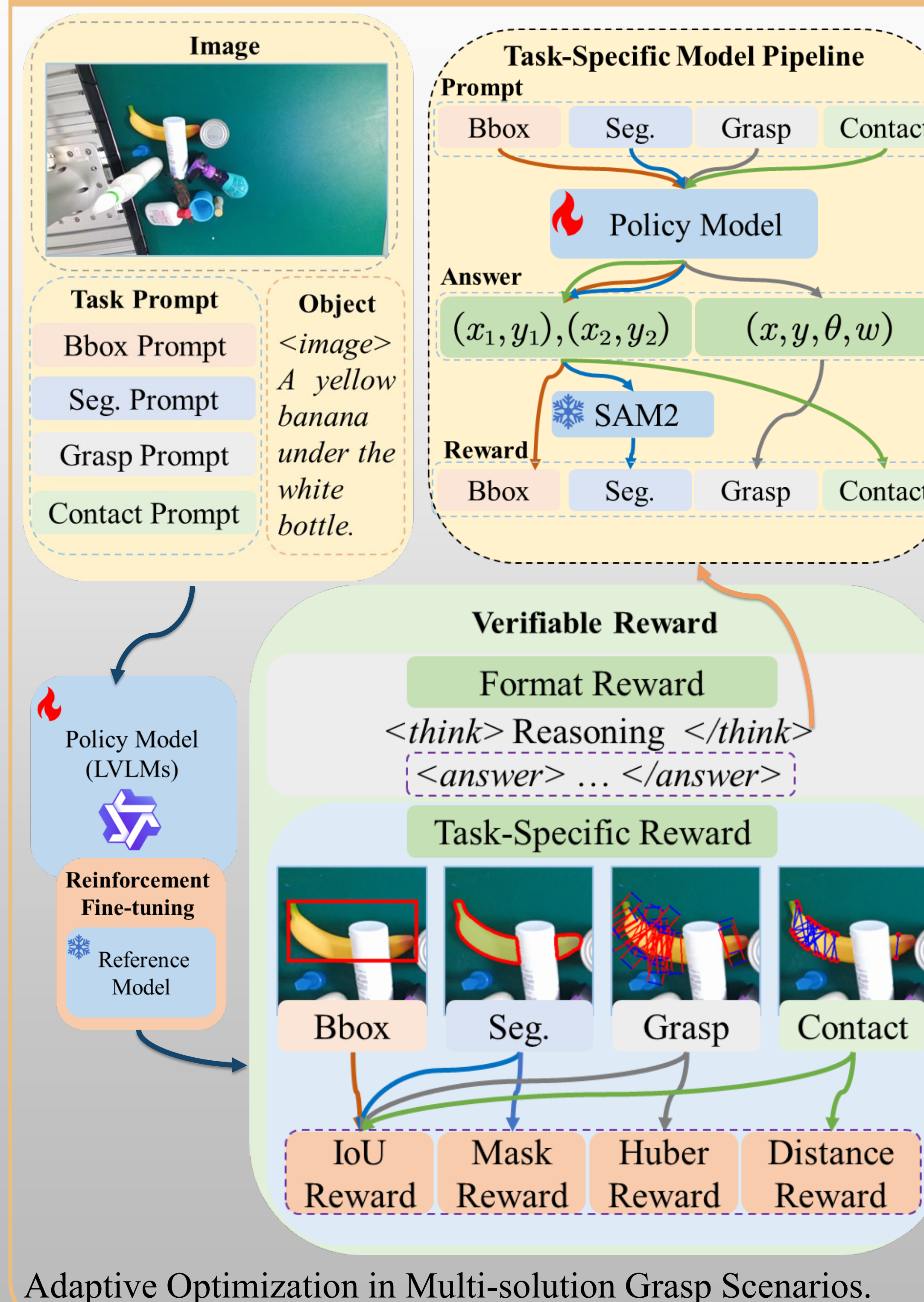
Type	Dataset	Grasp Label	Objects / Image	Scale			Grounding		Language Description		
				Images	Objects	Grasps	Seg.	Bbox	Object	Spatial	VLM-Check
Sim.	Dex-net [41]	Rect.	1	6.7M	1500	6.7M	✓	✓	✓	✓	✓
	Jacquard [14]	Rect.	1	54K	11K	1.1 M	✓	✓	✓	✓	✓
	VR-Grasping-101 [83]	6-DoF	1	10K	101	4.8M	✓	✓	✓	✓	✓
	Kappler et al. [28]	6-DoF	1	700	80	300K	✓	✓	✓	✓	✓
	6-DoF GraspNet [48]	6-DoF	1	206	206	7M	✓	✓	✓	✓	✓
	Epper et al. [115]	6-DoF	1	21	21	1B	✓	✓	✓	✓	✓
Synth.	EGAD! [47]	6-DoF	1	2231	2231	233K	✓	✓	✓	✓	✓
	ACRONYM [16]	6-DoF	-	-	8872	17.7M	✓	✓	✓	✓	✓
	MetaGraspNet [20]	6-DoF	~12	217K	82	-	✓	✓	✓	✓	✓
	Grasp-Anything [76]	Rect.	-	1M	~3M	~600M	✓	✓	✓	✓	✓
	Grasp-Anything++ [75]	Rect.	-	1M	~3M	~600M	✓	✓	✓	✓	✓
	Real	Cornell [26]	Rect.	1	1035	240	8019	✓	✓	✓	✓
Pinto et al. [57]		Rect.	-	50K	150	50K	✓	✓	✓	✓	✓
Levine et al. [32]		Rect.	-	800K	-	800K	✓	✓	✓	✓	✓
VMRD [88]		Rect.	~5	4683	~200	100K	✓	✓	✓	✓	✓
Multi-Object [11]		Rect.	~3	96	-	2904	✓	✓	✓	✓	✓
GraspClutter6D [3]		6-DoF	~15	52K	200	9.3B	✓	✓	✓	✓	✓
OCID-Grasp [11]		Rect.	~20	11K	~40	75K	✓	✓	✓	✓	✓
GraspNet [18]		6-DoF	~10	97K	88	~1.2B	✓	✓	✓	✓	✓
RealVLG-11B (Ours)	Rect.	1~20	165K	~800	11B	✓	✓	✓	✓	✓	

Annotation Pipeline with Human Verification



Closed-loop Verification Ensures Alignment of Language, Detection, Segmentation, and Grasping.

3. RealVLG-R1 Model



4. RealVLG-Benchmark

Benchmark Evaluation Summary

Benchmark

Method	Seen				Similar				Novel			
	Bbox	Seg.	Grasp	Contact	Bbox	Seg.	Grasp	Contact	Bbox	Seg.	Grasp	Contact
Qwen-VL-Max [13]	92.393.251.3	87.644.877.5	16.016.771.2	347.25.615.4	89.490.064.3	86.842.566.4	14.971.117.5	32.67.018.3	88.488.867.7	47.388.664.0	8.15.422.6	18.820.018.3
Gemin2.5-Flash [21]	60.161.508.4	3.314.183.4	3.212.210.0	3.472.610.0	59.260.580.9	3.813.180.9	2.52.210.0	6.347.110.0	62.663.795.1	3.914.095.1	3.401.910.0	4.04.310.0
Qwen2.5VL-72B [4]	50.051.909.2	75.143.509.2	4.32.994.1	25.217.787.0	51.252.494.5	75.142.494.5	4.01.492.3	27.822.985.5	52.953.796.9	45.938.797.0	3.11.395.7	14.58.698.8
Qwen2.5VL-72B + SFT	56.457.110.0	76.243.410.0	3.471.710.0	28.621.410.0	54.755.910.0	75.463.210.0	3.62.110.0	29.328.210.0	57.258.610.0	46.238.910.0	4.41.510.0	19.715.810.0
RealVLG-R1-7B (GRPO)	87.287.696.1	87.445.296.0	34.740.310.0	51.833.210.0	86.587.197.9	85.643.897.9	28.231.910.0	47.831.910.0	78.579.696.2	60.339.496.2	16.317.110.0	26.926.210.0
RealVLG-R1-3B (GSPO)	87.788.410.0	88.745.310.0	29.233.610.0	48.629.610.0	84.695.310.0	85.043.710.0	27.630.610.0	46.528.110.0	78.079.310.0	49.1396.10.0	15.39.110.0	26.917.710.0
Qwen2.5VL-72B [4]	50.051.909.2	75.143.509.2	4.32.994.1	25.217.787.0	51.252.494.5	75.142.494.5	4.01.492.3	27.822.985.5	52.953.796.9	45.938.797.0	3.11.395.7	14.58.698.8
Qwen2.5VL-72B + SFT	56.457.110.0	76.243.410.0	3.471.710.0	28.621.410.0	54.755.910.0	75.463.210.0	3.62.110.0	29.328.210.0	57.258.610.0	46.238.910.0	4.41.510.0	19.715.810.0
RealVLG-R1-7B (GRPO)	87.287.696.1	87.445.296.0	34.740.310.0	51.833.210.0	86.587.197.9	85.643.897.9	28.231.910.0	47.831.910.0	78.579.696.2	60.339.496.2	16.317.110.0	26.926.210.0
RealVLG-R1-3B (GSPO)	87.788.410.0	88.745.310.0	29.233.610.0	48.629.610.0	84.695.310.0	85.043.710.0	27.630.610.0	46.528.110.0	78.079.310.0	49.1396.10.0	15.39.110.0	26.917.710.0
Qwen2.5VL-72B [4]	56.858.310.0	80.643.810.0	4.11.843.5	33.726.399.2	57.357.999.6	78.042.499.6	3.20.039.4	32.823.599.6	55.656.810.0	48.339.210.0	1.60.064.6	20.012.810.0
Qwen2.5VL-72B + SFT	64.863.510.0	82.244.110.0	5.1/3.180.6	38.731.210.0	66.782.710.0	79.642.610.0	4.62.795.4	37.429.610.0	64.965.210.0	48.539.410.0	3.51.478.1	26.416.410.0
RealVLG-R1-7B (GRPO)	88.088.510.0	89.245.310.0	32.939.399.6	53.933.199.2	83.454.210.0	85.043.810.0	28.836.399.6	50.514.997.4	84.605.410.0	52.239.710.0	17.019.198.8	28.120.710.0
RealVLG-R1-3B (GSPO)	89.089.610.0	88.945.410.0	33.632.810.0	55.337.210.0	86.486.910.0	86.643.810.0	27.930.210.0	50.636.210.0	88.588.910.0	52.039.710.0	16.518.310.0	28.925.010.0

Data Quality

Real-World Vision-Language Grasping

Conclusion & Future Work

- Unified Real-World VLG/Grasping framework
- Zero-shot deployment in clutter
- **Future:** 3D VLG-to-Action, SmoVLM for realtime