

Multimodal Ranking for Target Objects and Receptacles Based on Open-Vocabulary Instructions

Ryosuke Korekata, Kanta Kaneda, Shunya Nagashima, Yuto Imai, Komei Sugiura (Keio University)

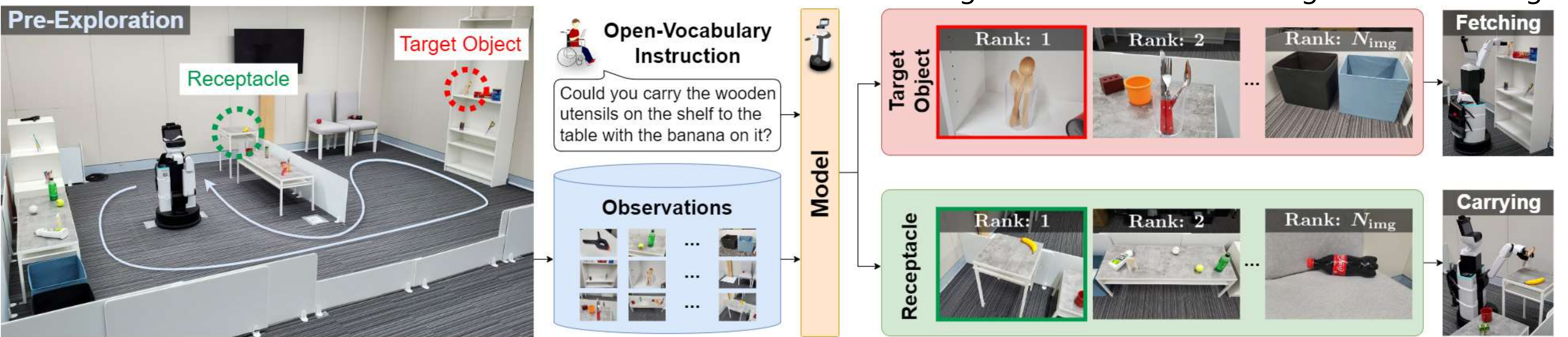
Abstract: Domestic service robots (DSRs)

- Task** Image retrieval-based **open-vocabulary** fetch-and-carry
- Novelty**
1. **LLM-based** Task Paraphraser
 2. Segment Anything Region Encoder
- Results**
1. Outperformed the baseline methods on the novel dataset
 2. Achieved a **success rate of 82%** in the physical experiments

Related work: Open-vocabulary manipulation

- RREx-BoT [Sigurdsson+, IROS23] Vision-and-language navigation based on the images collected through pre-exploration
- MultiRankIt [Kaneda+, RA-L24] Object fetching tasks based on the human-in-the-loop setting
- OVMM [Yenamandra+, CoRL23] Open-vocabulary mobile manipulation tasks SOTA method achieved SR of only 10%

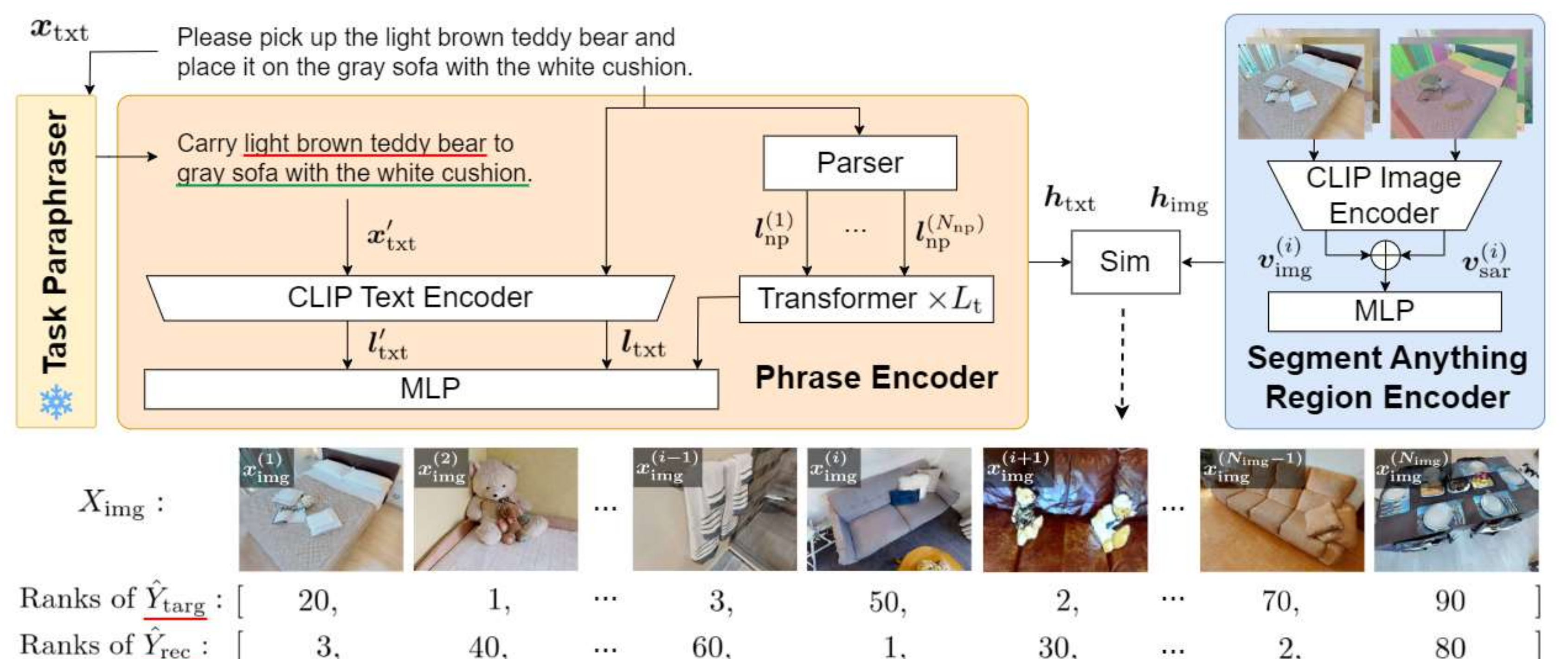
⊗ Few existing methods handle the image retrieval setting



Method

Retrieves images of both target objects and receptacles using **multimodal foundation models**

- Task Paraphraser:**
Paraphrases instructions including redundancy into standardized format using LLM
- Phrase Encoder:**
Obtains fine-grained text features from instructions using CLIP [Radford+, ICML21]
- Segment Anything Region Encoder:**
Enhances visual features regarding shape and contour of objects by utilizing SAM [Kirillov+, ICCV23]



Settings: 1. Newly-built dataset, 2. Physical experiments in the standard environment [Okada+, AR19]

1. LTRRIE-FC dataset based on HM3D [Ramakrishnan+, NeurIPS21]

Instructions were collected by 226 annotators using a crowdsourcing service

#envs	#images	#instrs	Sentence length
774	7,148	6,581	15.69 (average)



2. Fetch-and-carry actions based on user instructions



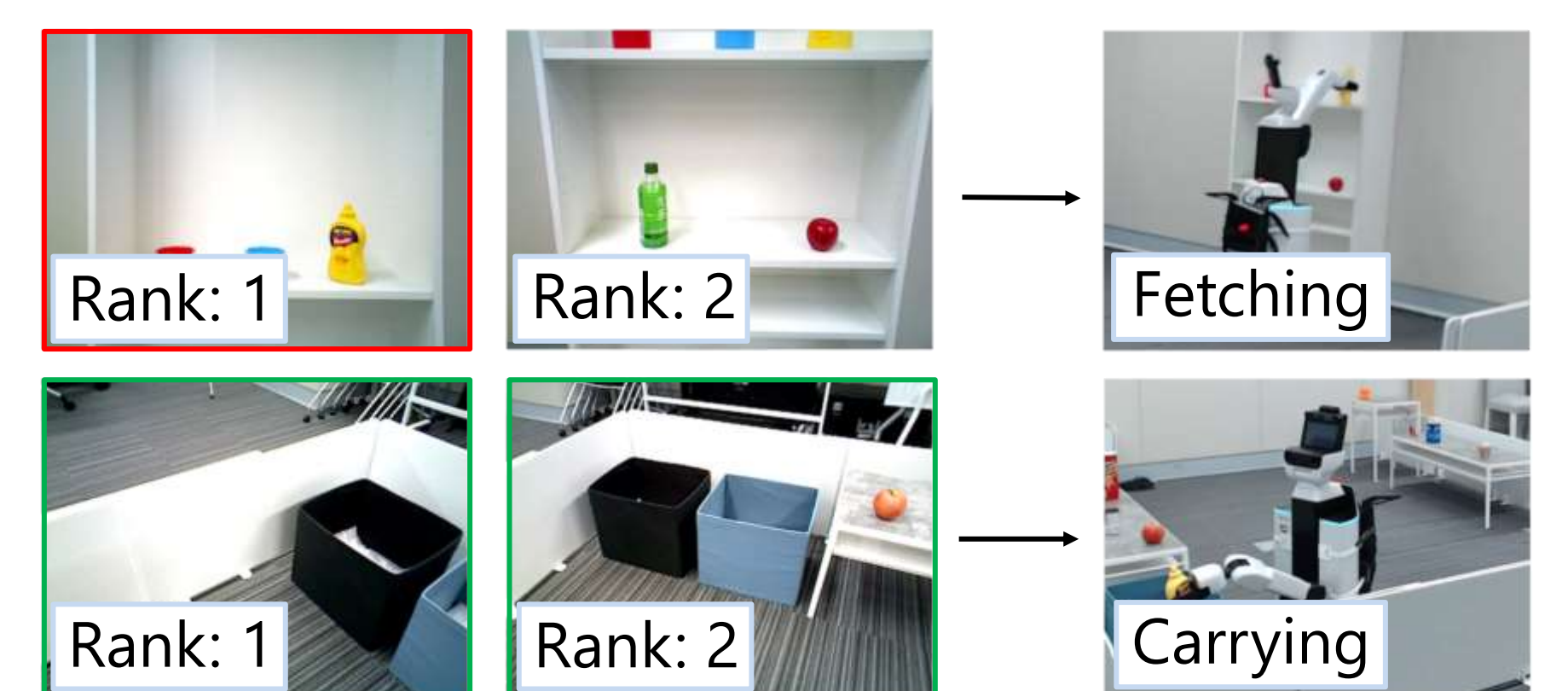
Results: 1. Outperformed the baseline methods, 2. Achieved a success rate of 82% in zero-shot setting

1. Quantitative results: Standard metrics for image retrieval
2. Qualitative result: Successful sample

	[%]	MRR↑	Recall@5↑	Recall@10↑
CLIP [Radford+, ICML21]		10.8	13.7	24.9
MultiRankIt [Kaneda+, RA-L24]		20.5 ± 2.3	30.1 ± 3.4	48.2 ± 1.4
Ours		22.5 ± 1.4	33.2 ± 1.8	53.0 ± 2.5

Improvements: +2.0 (MRR), +3.1 (Recall@5), +4.8 (Recall@10)

Target Object



Receptacle

"Can you take the mustard container on the shelf to the black box?"