

Open Master Thesis Topic: Large Vision-Language Networks for Open-Vocabulary Robotic Manipulation



June 22, 2023

Supervisors: Snehal Jauhri (Email: snehal.jauhri@tu-darmstadt.de), Ali Younes, Georgia Chalvatzaki

1 Project Description

Robots are expected to soon leave their factory/laboratory enclosures and operate autonomously in everyday unstructured environments such as households. For this, robotic grasping and manipulation are important problems to solve. Significant progress has been made in recent years due to advancements in 3D Computer Vision and Deep Reinforcement Learning. Recent work at our iROSA lab [1] has successfully utilized learned neural geometric reconstruction of scenes for performing robotic grasping. However, most grasping methods only consider the problem of ‘de-cluttering’ a scene i.e. picking and removing objects placed in a scene. These methods do not consider the *semantics* of the objects. They also do not consider the problem of scene understanding and finding placing locations for the objects on tables/shelves etc. Semantic information is especially important when considering real-world robotic applications where the robot needs to re-arrange objects as per a set of language instructions or human inputs.

When it comes to semantic segmentation of scenes, many sophisticated neural networks exist [2]. However, a challenge when using such methods in the real world is that the semantic classes rarely align perfectly with the language input received by the robot. For instance, a human language instruction might request a ‘glass’ or ‘water’, but the semantic classes detected might be ‘cup’ or ‘drink’. Nevertheless, with the rise of large language and vision-language models, we now have capable segmentation models that do not directly predict semantic classes but use learned association between language queries and classes to give us ‘open-vocabulary’ segmentation [3]. Some models are especially powerful since they can be used with any arbitrary language query, as shown in Fig 1.

In this thesis, we aim to build on advances in 3D vision-based robot manipulation and large open-vocabulary vision models [3] to build a full pick-and-place pipeline for real-world manipulation. We also aim to find synergies between scene reconstruction and semantic segmentation to determine if knowing the object semantics can aid the reconstruction of the objects and, in turn, aid manipulation.

Very relevant to this thesis is the Open Vocabulary Mobile Manipulation (OVMM) Challenge [4], which was announced recently. The iROSA lab currently plans to participate in the challenge, and if the Master student is willing, this thesis could be a key contribution towards it. This could be a unique opportunity to directly apply the thesis work to a challenge at a top machine learning conference (NeurIPS) and compete against the best machine/robot-learning labs.



Figure 1: Recent work on open-vocabulary segmentation [3] can not only provide semantic masks for any classes but can also use any language query such as ‘I am hungry’ (left) or ‘Sports’ (right).

