

# MASTER THESIS

Thesis submitted in fulfillment of the requirements for the degree of Master of Science in Engineering at the University of Applied Sciences Technikum Wien - Degree Program Mechatronics/Robotics

## **SAGE: Multi object semantic aware guided exploration with persistent memory**

By: Kevin Eppacher, BSc

Student Number: 2310331013

Supervisor: Simon Schwaiger, MSc

Vienna, July 17, 2025

This work was conducted in the context of the project “Stadt Wien Kompetenzteam für Drohnentechnik in der Fachhochschulausbildung” (project number MA23 35-02, financed by the Department MA23 for Economic Affairs, Labour and Statistics of the City of Vienna).

Funded by



Economic Affairs,  
Labour and Statistics



# Declaration

“As author and creator of this work to hand, I confirm with my signature knowledge of the relevant copyright regulations governed by higher education acts (see Urheberrechtsgesetz /Austrian copyright law as amended as well as the Statute on Studies Act Provisions / Examination Regulations of the UAS Technikum Wien as amended).

I hereby declare that I completed the present work independently and that any ideas, whether written by others or by myself, have been fully sourced and referenced. I am aware of any consequences I may face on the part of the degree program director if there should be evidence of missing autonomy and independence or evidence of any intent to fraudulently achieve a pass mark for this work (see Statute on Studies Act Provisions / Examination Regulations of the UAS Technikum Wien as amended).

I further declare that up to this date I have not published the work to hand nor have I presented it to another examination board in the same or similar form. I affirm that the version submitted matches the version in the upload tool.“

Vienna, July 17, 2025

Signature

# Kurzfassung

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

**Schlagworte:** Keyword1, Keyword2, Keyword3, Keyword4

# Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

**Keywords:** Keyword1, Keyword2, Keyword3, Keyword4

# Acknowledgements

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem Statement . . . . .	1
1.2	Thesis Structure . . . . .	2
<b>2</b>	<b>State of the Art</b>	<b>2</b>
2.1	Semantic Multi-Object Search Approaches . . . . .	3
2.2	Map Reconstruction and Persistent Semantic Mapping . . . . .	3
2.3	Object Detection and Promptable Models . . . . .	4
<b>3</b>	<b>Methods</b>	<b>4</b>
3.1	System Overview . . . . .	4
3.2	Semantic Frontier Mapping . . . . .	5
3.2.1	Frontier Detection and Calculation . . . . .	5
3.2.2	Value Map Generation using Vision-Language Models . . . . .	5
3.2.3	Navigation to High-Value Frontiers . . . . .	6
3.3	Persistent 3D Semantic Mapping . . . . .	6
3.3.1	Global Map Construction with Open-Fusion . . . . .	6
3.3.2	Semantic Clustering and Graph Node Generation . . . . .	6
3.4	Promptable Zero-Shot Detection . . . . .	6
3.5	Fusion Strategy . . . . .	6
<b>4</b>	<b>Implementation</b>	<b>7</b>
4.1	Simulation Environment . . . . .	7
4.2	Dataset . . . . .	7
4.3	Used Software . . . . .	7
4.4	Used Hardware . . . . .	8
4.5	Evaluation Metrics . . . . .	8
<b>5</b>	<b>Discussion and Results</b>	<b>9</b>
5.1	Results on Semantic Multi-Object Search Tasks . . . . .	10
5.1.1	Experiment 1: Single-Object Success Rate (SR) . . . . .	10
5.1.2	Experiment 2: Navigation Efficiency (SPL) . . . . .	10
5.1.3	Experiment 3: Multi-Object Success Rate (MSR) . . . . .	10
5.1.4	Experiment 4: Ablation of Exploitation (OpenFusion) . . . . .	10
5.2	Experiment 5: Improving Detection Robustness via Semantic Fusion . . . . .	10
5.3	Experiment 6: Real-World Deployment . . . . .	10
5.4	Experiment 7: Comparison of VL-Models for Frontier Scoring . . . . .	10

<b>6 Summary and Outlook</b>	<b>10</b>
<b>Bibliography</b>	<b>11</b>
<b>List of Figures</b>	<b>12</b>
<b>List of Tables</b>	<b>13</b>
<b>List of source codes</b>	<b>14</b>
<b>A Appendix A</b>	<b>15</b>
<b>B Appendix B</b>	<b>16</b>

# 1 Introduction

## 1.1 Problem Statement

1. Motivation & Relevance
2. Domain and Motivation: e.g. Search and Rescue (SAR), semantic Navigation, exploration in unknown environments, etc.
3. Which solutions exist? Comparison of general properties and performance metrics:
  - a) One Map to Find Them All (OneMap) [1]
  - b) Visual Language Maps for Robot Navigation (VLMaps) [2]
  - c) Vision-Language Frontier Maps for Zero-Shot Semantic Navigation (VLFM) [3]
  - d) Open-Vocabulary 3D Scene Graphs for Perception and Planning (ConceptGraphs) [4]
  - e) Object Goal Navigation using Goal-Oriented Semantic Exploration (SemExp) [5]
  - f) Learning Generalizable Feature Fields for Mobile Manipulation (GeFF) [6]

Methods	Persistent Memory	Real-Time <u>Cabability</u>	<u>Representation</u> (2D/3D)	<u>Open-Vocab Cabability</u>	Explicit Frontiers
VLFM	x	✓	2D	✓	✓
VLMaps	✓	X (offline)	2.5D	✓	x
OneMap	✓	2 Hz ( <u>drops heavily proportional to map size</u> )	2.5D	✓	✓
ConceptGraphs	✓	X (offline)	3D	✓	x
SemExp	✓	✓	2D	x	✓
GeFF	✓	✓	3D	✓	x

Figure 1: Comparison of state-of-the-art methods regarding common properties.

4. What is the technical problem? A combination of the following:
  - a) No persistent memory
  - b) Not Real-Time
  - c) 3D-Representation
  - d) Sensibility to false positives for zero shot object detection
5. Scientific Contribution: This work contributes to the state-of-the-art by introducing a hybrid semantic exploration framework that combines open-vocabulary frontier scoring [3]



Methode	Success Rate (SR)	Success Rate (Multi-Obj.)	SPL (Pfadlänge)	Semantic Precision	Robustheit kleine Objekte
VLFM	~81%	~60%	~60%	mittel	niedrig
VLMaps	~63%	~46%	~41%	mittel (~65%)	mittel
OneMap	89%	70%	70%	hoch (~85%)	mittel (~65-70%)
ConceptGraphs	hoch (geschätzt ~85-90%)	mittel (~65-70%, geschätzt)	mittel (~60%)	hoch	hoch
SemExp	64%	32%	35%	mittel	niedrig
GeFF	hoch (~85%, geschätzt)	mittel (~70%, geschätzt)	mittel/hoch (~65%)	hoch	mittel/hoch

Figure 2: Comparison of state-of-the-art methods regarding performance metrics.

with global semantic 3D scene representation [7], enabling multi-object search with open-vocabulary queries (text, image, audio) during autonomous exploration. The proposed method is evaluated for improvements in multi-object Success Rate (SR) as well as path efficiency measured by Success weighted by Path Length (SPL). The proposed method is evaluated with respect to multi-object SR and path efficiency, measured via SPL. This is achieved by streamlining existing frontier scoring systems that selectively choose expert models based on the current tasks into a singular vision-language model based on Segment Everything Everywhere All at Once (SEEM) [8], substantially reducing GPU memory requirements compared to traditional semantic exploration pipelines using multiple separate models. Exploration is complemented by a sensor-based fusion strategy that integrates semantic detections from SEEM [8] with clustered relevance fields from semantic 3D mapping [7], applying spatial confidence weighting to enhance robustness against false positives in zero-shot object detection. Experimental validation of the proposed system on a real mobile robot, focusing on practical aspects such as real-time performance and robustness to sensor noise including depth inaccuracies and changing lighting conditions.

## 1.2 Thesis Structure

# 2 State of the Art

In this chapter, the current state of research in semantic multi-object search, map reconstruction, and object detection is reviewed. The goal is to identify strengths and limitations of existing methods and establish the technological context for the proposed hybrid approach. The chapter is divided into three key areas: approaches for searching multiple objects semanti-

cally, techniques for building and maintaining persistent semantic maps, and recent advances in object detection and promptable models for open-vocabulary tasks.

## 2.1 Semantic Multi-Object Search Approaches

- Review of methods targeting simultaneous or sequential search for multiple objects in unknown environments.
- Analysis of VLFM, VLMaps, OneMap, GeFF, . . . regarding:
  - Success Rate (SR) and Success weighted by Path Length (SPL) as key metrics.
  - Real-time capabilities and computational requirements.
  - Handling of open-vocabulary queries.
- Discussion of semantic exploration frameworks combining language models with spatial reasoning.
- Challenges of maintaining semantic context across multiple targets.

## 2.2 Map Reconstruction and Persistent Semantic Mapping

- Overview of approaches to build and update semantic maps during exploration.
- Techniques for fusing sensor data into persistent 2D/3D representations.
- Comparison of representations (Octomaps, point clouds, voxel grids) in terms of:
  - Memory efficiency.
  - Ability to store semantic labels persistently.
- Discussion of . . .
  - ConceptGraphs
  - ConceptGraph-Online
  - OpenFusion
  - Clio
  - OpenScene
  - GeFF
  - CLIP-Fields
  - ConceptFusion
  - VLMaps
  - LERF

as examples of global 3D semantic maps.

- Limitations in updating or correcting the map after wrong detections.

## 2.3 Object Detection and Promptable Models

- Review of traditional and open-vocabulary object detection methods.
- Analysis of grounding-capable detectors and segmentation models for zero-shot tasks.
- Specific evaluation of the following models for their suitability in semantic multi-object search:
  - YOLOv7
  - GroundingDINO
  - MobileSAM
  - GroundedSAM
  - SEEM
  - OWL-ViT
  - MaskDINO
- Discussion of promptable vision-language models supporting multi-modal queries (text, image, audio).
- Challenges with false positives in zero-shot settings and their implications for reliable multi-object detection.

## 3 Methods

This chapter details the methods developed for semantic exploration, persistent 3D mapping, promptable object detection, and robust fusion strategies for multi-object search.

### 3.1 System Overview

- Presentation of the overall architecture of the exploration, detection, mapping, and fusion pipeline.
- Description of data flow between exploration (frontier evaluation), detection (promptable models), and exploitation (persistent semantic mapping).
- Explanation of how exploration and mapping components interact to progressively build a semantic understanding of the environment.

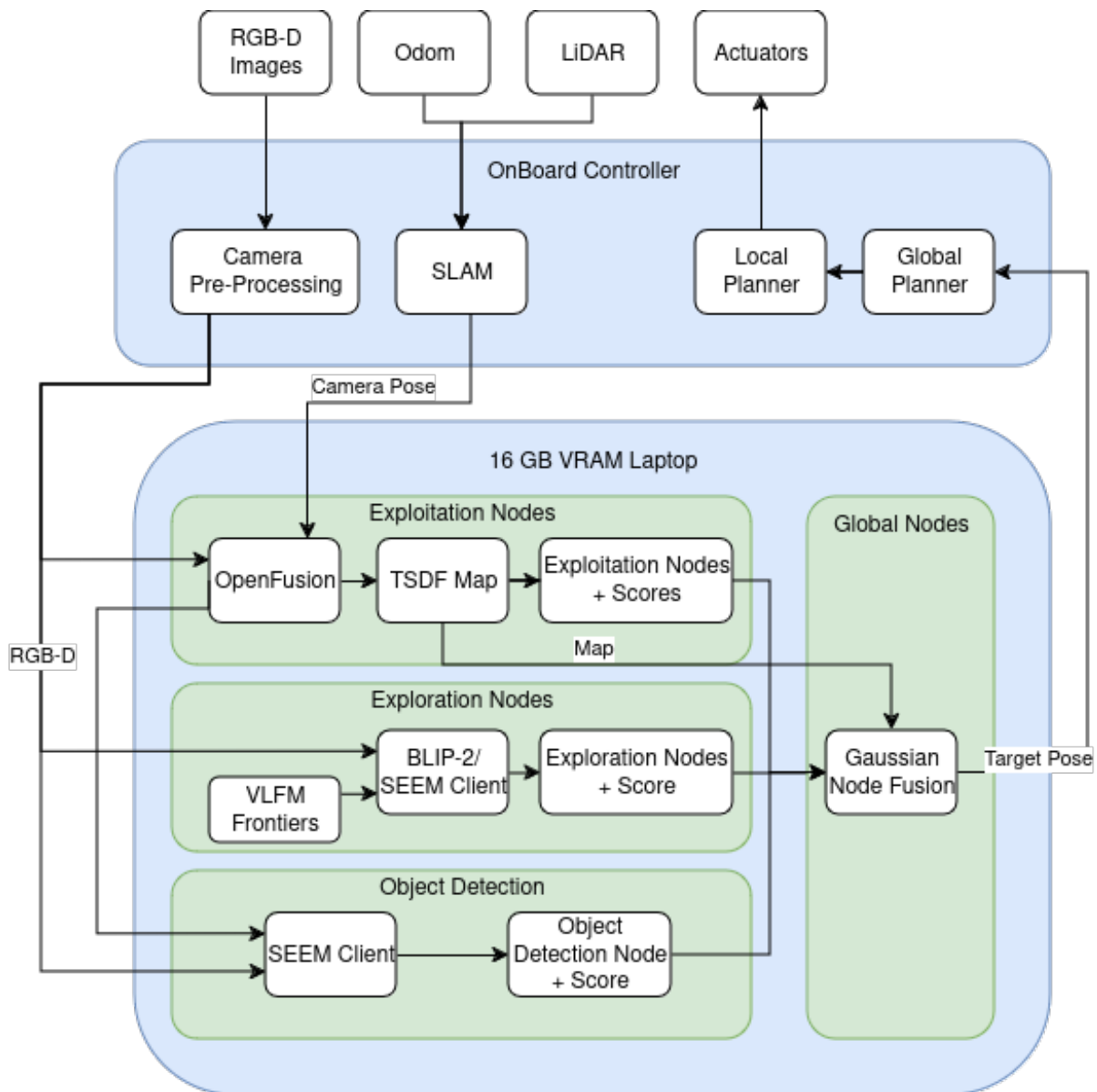


Figure 3: System architecture

## 3.2 Semantic Frontier Mapping

### 3.2.1 Frontier Detection and Calculation

- Detection of frontiers on a 2D occupancy grid to identify candidate regions for exploration.
- Application of classical frontier-based exploration algorithms extended with semantic information.

### 3.2.2 Value Map Generation using Vision-Language Models

- Computation of value maps by evaluating cosine similarity between multi-modal queries (text, image, audio) and scene observations.

- Use of promptable vision-language models (e.g., SEEM) to assign semantic relevance scores to each region.
- Dynamic update of value maps as new observations are integrated.

### 3.2.3 Navigation to High-Value Frontiers

- Selection of the frontier with the highest semantic relevance score.
- Continuous re-evaluation of frontiers during exploration to adapt to changing scene semantics.
- Strategy for balancing exploration efficiency and semantic search goals.

## 3.3 Persistent 3D Semantic Mapping

### 3.3.1 Global Map Construction with Open-Fusion

- Incremental creation of a global semantic point cloud map integrating RGB-D observations over time.
- Registration of observations using robot poses to maintain a consistent world representation.
- Association of semantic labels with 3D points based on query relevance scores.

### 3.3.2 Semantic Clustering and Graph Node Generation

- Clustering of points with similar semantic labels to form object-level hypotheses.
- Construction of semantic graph nodes representing detected object instances with aggregated confidence scores.
- Maintenance of the semantic graph as a persistent memory for multi-object search tasks.

## 3.4 Promptable Zero-Shot Detection

- tbd

## 3.5 Fusion Strategy

- Formulation of a fusion strategy combining frontier-based semantic relevance with persistent 3D semantic mapping.

- Design of decision-making algorithms leveraging combined semantic information to robustly answer multi-object search queries.

## 4 Implementation

This section details the practical implementation of the proposed approach, covering the simulation and real-world setup, datasets, software stack, and hardware configuration.

### 4.1 Simulation Environment

- Evaluation of simulation frameworks for indoor semantic navigation:
  - HabitatSim: Realistic Matterport3D-based environments with semantic annotations.
  - Isaac Sim / Isaac Lab: GPU-accelerated simulation, advanced physics, support for RTX ray tracing.
  - MuJoCo: High-speed physics engine, limited support for complex indoor scenes.
  - Ignition Gazebo: Modular simulator, ROS2 integration, good for real-robot transfer.
- ...

### 4.2 Dataset

- Use of **Matterport3D** scenes for realistic indoor environments with ground truth 3D reconstruction and semantic annotations.
- Incorporation of the **Habitat Navigation Challenge 2023** tasks to benchmark exploration and navigation performance (SR, SPL).

### 4.3 Used Software

- ROS2-based implementation (Humble Hawksbill) as middleware.
- Navigation stack: Navigation2 (Nav2) for frontier-based exploration and path planning.
- DDS communication layer for distributed communication between detection, mapping, and control nodes.
- Custom RobotDriver for interfacing with real robot hardware.

- Integration of promptable models (SEEM, GroundingDINO, etc.) for real-time zero-shot detection during exploration and exploitation.
- ...

## 4.4 Used Hardware

- **PC:**
  - CPU: AMD Ryzen 9 5950X 16-Core Processor
  - Motherboard: B550 Gaming X V2
  - GPU: ASUS TUF Gaming RTX 4090 24GB OC Edition
  - RAM: 64GB Corsair Vengeance LPX DDR4
- **Real Robot:** Configuration and components to be determined (tbd).

## 4.5 Evaluation Metrics

This section defines the evaluation metrics used throughout the experiments and assigns them to each corresponding experiment.

- **Experiment 1 – Success Rate (SR):** Measures the proportion of tasks in which the robot successfully reaches the queried single goal object. This metric reflects the system’s ability to semantically ground a user-specified object and to navigate toward it reliably. It serves as a fundamental indicator of task success and is essential for evaluating overall system effectiveness in basic search scenarios. *Evaluation against:* VLFM, VLMaps, OneMap, GeFF

$$SR = \frac{1}{N} \sum_{i=1}^N S_i$$

where  $S_i = 1$  if the goal was reached in episode  $i$ , and 0 otherwise;  $N$  is the total number of episodes.

- **Experiment 2 – Path Efficiency (SPL):** SPL measures the efficiency of successful navigation by comparing the shortest possible path to the actual path taken. It is defined only for successful runs and penalizes overly long trajectories. In the context of semantic exploration, SPL provides insight into how effectively the system prioritizes relevant regions and minimizes detours when searching for target objects.

$$SPL = \frac{1}{N} \sum_{i=1}^N S_i \cdot \frac{l_i}{\max(p_i, l_i)}$$

where  $S_i$  is the success indicator for episode  $i$ ,  $l_i$  is the shortest path length to the goal,  $p_i$  is the actual path length taken, and  $N$  is the total number of episodes.

- **Experiment 3 – Multi-Object Success Rate (MSR):** The average number of successfully found objects per episode (*Progress, PR*) captures partial success in multi-goal navigation. SPL is computed separately for each object in sequence, conditioned on the success of the previous one. This highlights the system’s ability to reuse semantic map information and improve efficiency across successive targets.

$$PR = \frac{1}{N} \sum_{i=1}^N C_i$$

where  $C_i$  is the number of successfully found objects in episode  $i$ , and  $N$  is the total number of episodes.

- **Experiment 4 – Ablation: Memory Component (OpenFusion):** Comparison of multi-object progress and SPL with and without the semantic 3D memory component. Highlights the contribution of global semantic mapping to task success and efficiency in the hybrid system.
- **Experiment 5 – Robustness to False Positives (Fusion Strategy): Metrics:** Evaluated using semantic precision and false positive rate. Assesses the effect of the fusion strategy on filtering erroneous detections to improve overall task success.

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{FPR} = \frac{FP}{FP + TN}$$

where  $TP$  and  $FP$  are the numbers of true and false positive semantic detections, and  $TN$  is the number of true negatives (i.e., correctly rejected background regions or non-target classes).

- **Experiment 6 – Real-World System Performance: Metrics:**

- SR, MSR, SPL – for search performance under real-world conditions.
- System metrics – CPU/GPU usage, FPS, inference latency.

*Objective: Assess robustness, efficiency, and deployability in physical environments.*

- **Experiment 7 – Vision-Language Model Comparison (SEEM vs. BLIP2): Metrics:** Comparison of value maps generated using SEEM and BLIP2 for frontier scoring, focusing on their impact on exploration performance (SR and SPL) and GPU efficiency. Evaluates whether a lightweight unified model like SEEM can match or surpass BLIP2 in guiding semantic exploration through effective value map construction.

## 5 Discussion and Results

This chapter presents the experimental evaluation of the proposed hybrid semantic exploration system. Each experiment targets a specific research question and is evaluated using quanti-



tative performance metrics.

## 5.1 Results on Semantic Multi-Object Search Tasks

5.1.1 Experiment 1: Single-Object Success Rate (SR)

5.1.2 Experiment 2: Navigation Efficiency (SPL)

5.1.3 Experiment 3: Multi-Object Success Rate (MSR)

5.1.4 Experiment 4: Ablation of Exploitation (OpenFusion)

## 5.2 Experiment 5: Improving Detection Robustness via Semantic Fusion

## 5.3 Experiment 6: Real-World Deployment

## 5.4 Experiment 7: Comparison of VL-Models for Frontier Scoring

# 6 Summary and Outlook

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

# Bibliography

- [1] F. L. Busch, T. Homberger, J. Ortega-Peimbert, Q. Yang, and O. Andersson, One map to find them all: Real-time open-vocabulary mapping for zero-shot multi-object navigation, 2025. arXiv: [2409.11764](https://arxiv.org/abs/2409.11764) [cs.RO]. [Online]. Available: <https://arxiv.org/abs/2409.11764>.
- [2] C. Huang, O. Mees, A. Zeng, and W. Burgard, “Visual language maps for robot navigation,” in Proceedings of the IEEE International Conference on Robotics and Automation (ICRA), London, UK, 2023.
- [3] N. Yokoyama, S. Ha, D. Batra, J. Wang, and B. Bucher, “Vlfm: Vision-language frontier maps for zero-shot semantic navigation,” in International Conference on Robotics and Automation (ICRA), 2024.
- [4] Q. Gu et al., Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning, 2023. arXiv: [2309.16650](https://arxiv.org/abs/2309.16650) [cs.RO]. [Online]. Available: <https://arxiv.org/abs/2309.16650>.
- [5] D. S. Chaplot, D. Gandhi, A. Gupta, and R. Salakhutdinov, “Object goal navigation using goal-oriented semantic exploration,” in In Neural Information Processing Systems, 2020.
- [6] R.-Z. Qiu et al., “Learning generalizable feature fields for mobile manipulation,” arXiv preprint arXiv:2403.07563, 2024.
- [7] K. Yamazaki et al., “Open-fusion: Real-time open-vocabulary 3d mapping and queryable scene representation,” arXiv preprint arXiv:2310.03923, 2023.
- [8] X. Zou et al., Segment everything everywhere all at once, 2023. arXiv: [2304.06718](https://arxiv.org/abs/2304.06718) [cs.CV]. [Online]. Available: <https://arxiv.org/abs/2304.06718>.

# List of Figures

Figure 1 Comparison of state-of-the-art methods regarding common properties. . . . .	1
Figure 2 Comparison of state-of-the-art methods regarding performance metrics. . . . .	2
Figure 3 System architecture . . . . .	5

## List of Tables

## List of source codes

## A Appendix A

## B Appendix B