

# **Université Marie et Louis Pasteur (UMLP)**

Master of Science in Advanced Robotics, Mechatronics, and Automatic Control  
(ARMAC)

## **Master Thesis**

### **Multimodal Language-Grounded Frameworks for Autonomous Robot Navigation and Interaction**

**Author:**

Promise Emeziem Adiole

**Supervisors:**

Prof. Elmar Rueckert  
Nwankwo Linus (MSc.)

**Host Institution:**

Cyber-Physical Systems Group,  
Montanuniversität Leoben, Austria

**Academic Year:**

2024–2025

# Abstract

This thesis presents the study of language-grounded autonomous robotic system that integrates Large Language Models (LLMs) and Vision–Language Models (VLMs) within the ROS 2 Navigation (Nav2) framework. The objective of this research is to enable autonomous robots to interpret open-ended natural-language commands, perceive their surroundings semantically, and generate meaningful linguistic or geometric responses without task-specific programming. The proposed framework combines several state-of-the-art models: CLIP for contrastive image–text verification, BLIP-2 for open-ended image-grounded captioning, and the Segment Anything Model (SAM) for pixel-level segmentation and spatial grounding. These perception models are integrated with a ROS 2 Nav2 pipeline that governs motion planning, localisation, and control, while a Gradio-based user interface enables interactive dialogue between the robot and the user. All experiments were conducted using the TurtleBot3 Burger model in a Gazebo simulation environment. The robot consistently interpreted natural-language navigation commands, generated accurate descriptions of its surroundings, and maintained coherent dialogue through its LLM-based reasoning module. This study confirms that multimodal grounding through large pre-trained models can enhance the autonomy, interpretability, and human-likeness of service robots. It lays the foundation for deploying scalable, language-aware robotic systems capable of bridging symbolic reasoning and embodied interaction.

## **Acknowledgements**

This research was carried out as part of the Master's programme in Advanced Robotics, Mechatronics, and Automatic Control (ARMAC) at Université Marie et Louis Pasteur. I wish to express sincere gratitude to Professor Elmar Rueckert for his invaluable supervision, insightful feedback, and intellectual guidance throughout the internship and thesis development. Special thanks are also extended to Nwankwo Linus for his mentorship and constructive discussions that shaped the direction of this work.

I gratefully acknowledge the support of the Cyber-Physical Systems Group at Montanuniversität Leoben, whose collaborative environment and technical resources made this research possible. Appreciation is also due to colleagues and peers in the ARMAC programme for their assistance, encouragement, and meaningful exchange of ideas.

Finally, heartfelt thanks are offered to my family and friends for their unwavering support, patience, and inspiration during the entire course of study and research.

## **AI Assistance Disclosure**

The author utilized AI technology, specifically “ChatGPT”, as a writing assistant to refine the clarity, coherence, and style of the thesis report. All conceptual contributions, scientific analyses, system implementations, and research decisions remain the sole responsibility of the author.

# Contents

<b>Acknowledgements</b>	<b>ii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Motivation .....	1
1.2 Problem Statement .....	1
1.3 Aim and Objectives .....	2
1.4 Scope and Limitations .....	3
1.5 Research Questions .....	3
1.6 Key Contributions .....	4
1.7 Structure of the Thesis .....	4
<b>2 Related Work and Theoretical Foundations</b>	<b>5</b>
2.1 Introduction .....	5
2.2 Limitations of Classical Geometric Navigation .....	5
2.3 Foundations of Transformer-Based Language Models .....	6
2.4 Multidimensional Word Embedding Spaces .....	6
2.5 Contrastive Vision–Language Modelling (CLIP) .....	7
2.6 Generative Captioning and Cross-Modal Reasoning (BLIP-2) .....	8
2.7 Promptable Segmentation and Spatial Grounding (SAM) .....	9
2.8 Complementary Roles in Multimodal Grounding .....	10
2.9 Language, Policy Generation, and Embodied Control .....	10
2.10 Retrieval-Augmented Generation for Context-Aware Interaction .....	11
2.11 Failure Modes in Language-Grounded Autonomy .....	12
2.12 Transition to Methodological Framework .....	13
<b>3 Method and System Development</b>	<b>14</b>
3.1 Overview .....	14
3.2 System Architecture .....	14
3.3 LLMNode: Natural Language Understanding and Intent Mapping .....	14
3.4 VLMNode: Visual Grounding and Scene Semantic Understanding .....	16
3.5 Task Execution Node (TEN) .....	17
3.6 Semantic-to-Geometric Action Grounding .....	17
3.7 Motion Optimization with the TEB Local Planner .....	19
3.8 Failure Handling and Safety Behaviors .....	20
3.8.1 Uncertain Object Identification .....	20
3.8.2 Ambiguous Spatial Projection .....	20
3.8.3 Navigation Failure Recovery .....	20
3.8.4 Emergency Halt .....	20

3.9	Navigation and Control Using the ROS 2 Nav2 Framework .....	21
3.9.1	Core Components of Nav2 .....	21
3.9.2	Localization and Coordinate Frame Management .....	22
3.9.3	Integration with Language and Perception Modules .....	22
3.10	NACGUI: User Interaction Interface .....	22
3.11	Implementation Environment .....	22
3.12	Model Performance Analysis and Selection Rationale .....	23
3.12.1	Cost and API Integration Considerations .....	23
3.12.2	Prompt Efficiency and Token Cost Optimization .....	24
3.12.3	Final Selection Justification .....	24
3.13	Knowledge Augmentation through Retrieval-Augmented Generation (RAG) .....	27
<b>4</b>	<b>Numerical Simulations and Implementation</b>	<b>28</b>
4.1	Experimental Setup .....	28
4.1.1	Simulation Environment .....	29
4.1.2	Software and Hardware Configuration .....	29
4.2	Evaluation Scenarios .....	30
4.3	Navigation Performance Analysis .....	30
4.3.1	Path Execution .....	30
4.3.2	Quantitative Metrics .....	30
4.4	Perception and Language Evaluation .....	32
4.4.1	Visual Grounding Results .....	32
4.4.2	Semantic Description Quality .....	32
4.4.3	Spatial Projection Accuracy .....	32
4.5	Conversational Interaction Results .....	32
4.6	Overall Performance and Discussion .....	35
4.7	Summary .....	36
<b>5</b>	<b>Results, Discussion, and Comparative Evaluation</b>	<b>37</b>
5.1	Overview of Evaluation Findings .....	37
5.2	Navigation Performance Discussion .....	37
5.3	Perception and Scene Understanding .....	37
5.3.1	Visual Grounding Accuracy .....	37
5.3.2	Linguistic Scene Description .....	38
5.3.3	Spatial Projection Consistency .....	38
5.4	Conversational Reasoning and RAG Integration .....	38
5.5	Comparative Evaluation with Baseline Navigation .....	39
5.6	Latency, Scalability, and Resource Utilisation .....	39
5.7	Error Analysis .....	39
5.8	Broader Implications and Research Perspectives .....	40

5.9	Limitations .....	40
5.10	Discussion Summary .....	40
<b>6</b>	<b>Conclusions and Future Perspectives</b>	<b>41</b>
6.1	Summary of Research Objectives and Achievements .....	41
6.2	Future Perspectives .....	42
<b>Appendix</b>		<b>46</b>
	<b>Appendix A — Foundations: Backpropagation and Gradient Descent</b>	<b>46</b>
A.1	Introduction .....	46
A.2	Problem Setup .....	46
A.3	Backpropagation Derivations .....	46
A.4	Implementation .....	47
A.5	Results and Observations .....	48
A.6	Discussion .....	48

# 1 Introduction

In this chapter, the overall motivation, problem statement, research objectives, scope, and key contributions of the thesis are introduced and positioned within the broader context of language-grounded robot autonomy.

## 1.1 Background and Motivation

In recent years, robotics and artificial intelligence have evolved from rule-based automation into systems capable of perception, reasoning, and decision-making. Traditional robotic systems excel at precise, repetitive tasks in structured environments, yet they remain limited when faced with open-ended human instructions or dynamic, unstructured contexts. This limitation is rooted in the gap between symbolic reasoning (how humans communicate) and geometric reasoning (how robots perceive the world). Humans intuitively describe tasks such as “Go to the table near the window” or “Pick up the red box beside the chair,” using natural language, whereas most robotic systems operate on predefined control sequences or coordinates without understanding the semantics of the task.

The emergence of large-scale pre-trained language models (LLMs) and vision–language models (VLMs) such as CLIP [1], BLIP-2 [2], and LLaVA [3] has provided a transformative opportunity to bridge this gap. These models are trained on massive datasets that link linguistic expressions with visual and contextual information, enabling them to generalise across tasks without explicit reprogramming. By leveraging such models, autonomous robots can interpret human-like language commands, understand the objects and environments around them, and perform contextually relevant actions—a paradigm known as *language-grounded autonomy*.

This thesis explores how large multimodal models can be grounded in robotic perception systems to enable robots to reason semantically about their surroundings and respond intelligently to open-ended commands. The work aims to design a unified framework that integrates language, vision, and motion control in a single pipeline, allowing a robot to navigate, perceive, and interact with its environment through both geometric and linguistic understanding.

## 1.2 Problem Statement

Despite remarkable progress in autonomous navigation, current robotic systems remain predominantly geometric. They depend on localisation, mapping, and path-planning algorithms that operate purely on sensor-derived metrics such as distances and coordinates, without associating these geometric entities with meaningful semantic concepts. Consequently, when a user issues a natural language command such as “Go to the door” or “What object is beside the chair?”, the robot cannot process the semantic content of the instruction, even though the geometric information may be available in its map.

This disconnection between language semantics and geometric perception restricts robots to low-level, task-specific programming and prevents flexible, open-domain interaction. Moreover, most navigation systems are not capable of long-horizon reasoning, and the interpretation of human instructions typically requires hand-engineered parsing or symbolic translation [4, 5]. Thus, the key challenge addressed in this work is to ground large-scale pre-trained LLMs and multimodal VLMs in sensor-based robotic systems to achieve natural, adaptable, and interpretable interaction.

**Research problem:**

How can an autonomous robot interpret open-ended natural language commands, align them with sensory perception, and execute long-horizon tasks in a dynamic environment without explicit task-specific programming?

### 1.3 Aim and Objectives

The primary goal of this thesis is to enhance existing language-grounded robotic frameworks, such as the ReLI (Language-Agnostic Human–Robot Interaction) model [6], by developing a multimodal system that enables robots to interpret, plan, and execute tasks based on natural language input. Specifically, the work focuses on grounding the semantic understanding of pre-trained language and vision–language models to physical sensor data, thus enabling cross-modal alignment between text, image, and geometry.

#### Objectives

1. Develop a unified framework that integrates perception, language reasoning, and motion planning for autonomous service robots.
2. Design a modular pipeline connecting user input to robotic action through components such as:
  - (a) Task interpretation via LLMs,
  - (b) Visual understanding via CLIP [1] and BLIP-2 [2],
  - (c) Object segmentation via SAM [7], and
  - (d) Motion execution through the ROS 2 Navigation (Nav2) stack.
3. Implement and test the system in simulation (Gazebo/ROS 2) to validate navigation, perception, and language reasoning capabilities.
4. Evaluate cross-modal grounding by testing multiple natural language commands and measuring the robot’s accuracy and interpretability.
5. Provide real-world perspectives for deploying the framework on physical platforms such as the Unitree Go1 or Segway RMP robot.

## 1.4 Scope and Limitations

The scope of this work is restricted to simulation-based validation using ROS 2 and the Gazebo environment. While real-world deployment is considered for future work, this study focuses on demonstrating how semantic grounding can be achieved through multimodal alignment between perception and language models in a controlled virtual setup. The robot operates within an indoor environment (e.g., office or laboratory), equipped with RGB–D cameras and LiDAR sensors for mapping and perception.

The proposed system handles open-ended commands describing navigation goals or perceptual queries (e.g., “Go near the box” or “What do you see?”), but does not extend to manipulation or multi-agent collaboration. All perception and language models—CLIP [1], BLIP-2 [2], SAM [7], and the underlying LLM based on transformer architectures [8]—are used as frozen pre-trained components to ensure computational feasibility. Evaluation is performed in terms of navigation accuracy, visual recognition precision, and the naturalness of generated responses.

Despite its strong performance in simulation, several limitations were observed:

- **Simulation-only validation:** All experiments were conducted in Gazebo. Real-world deployment would require robustness against sensor noise, varying lighting conditions, motion blur, reflective surfaces, and timing uncertainty that are absent in simulation.
- **Computational overhead:** The combination of multiple deep-learning models (BLIP-2, CLIP, SAM, and the LLM) introduces high GPU load and memory consumption. Although real-time performance was achievable on a workstation-class system, embedded deployment on resource-limited hardware (e.g., Jetson platforms) would require model optimisation through quantisation, pruning, or distillation.
- **Generalisation and dataset bias:** CLIP and BLIP-2, trained on web-scale datasets, occasionally exhibited bias, ambiguity, or hallucinated object labels in visually challenging scenes. Domain-specific fine-tuning or curated datasets may improve stability in specialised indoor environments.
- **Limited reasoning depth:** While the LLM produced contextually coherent responses, its multi-step reasoning and causal planning abilities remain limited. More explicit reasoning modules or policy-generation frameworks could enhance long-horizon autonomy and decision-making.

Recognising these constraints provides a roadmap for enhancing robustness, efficiency, and interpretability in future iterations of the system.

## 1.5 Research Questions

To guide the investigation, the following research questions are formulated:

1. How can large-scale language and vision models be integrated within a robot’s perception and control architecture to interpret natural language commands?
2. What mechanisms enable effective cross-modal grounding between linguistic tokens, visual embeddings, and geometric coordinates?
3. To what extent can such a system generalise to open-ended tasks without explicit training for each scenario?
4. How do multimodal models influence the interpretability and adaptability of robot behaviour in simulated environments?

## 1.6 Key Contributions

This thesis contributes to the field of human–robot interaction and autonomous navigation in the following ways:

1. A unified multimodal framework that connects perception (vision), cognition (language reasoning), and motion (navigation control).
2. A novel integration pipeline combining CLIP [1], BLIP-2 [2], and SAM [7] with ROS 2 navigation and LLM reasoning via a Gradio-based interface.
3. A method for grounding abstract linguistic tokens to physical spatial coordinates, allowing robots to map language to the real world.
4. Empirical validation through simulation experiments demonstrating improved semantic interpretability and natural-language responsiveness.

## 1.7 Structure of the Thesis

This thesis is organised as follows. Chapter 1 introduces the motivation, objectives, and research questions. Chapter 2 reviews the state of the art, covering the evolution of large language and vision–language models and their applications in robotics. Chapter 3 presents the proposed methodology and framework development, including architectural design, mathematical formulations, and integration strategy. Chapter 4 describes the implementation setup, configuration, and simulation procedure. Chapter 5 discusses the experimental results and comparative evaluation of the system’s performance. Finally, Chapter 6 concludes the thesis, summarising the key findings and outlining directions for future research.

## 2 Related Work and Theoretical Foundations

In this chapter, the related work and theoretical foundations of language-grounded robot autonomy are presented, with a focus on geometric navigation, transformer-based language models, multimodal perception, and retrieval-augmented reasoning.

### 2.1 Introduction

Robotic autonomy has traditionally relied on geometric reasoning pipelines grounded in sensor-based perception and deterministic control loops. Systems such as simultaneous localisation and mapping (SLAM) [4], occupancy grid mapping, and sampling-based path planning have enabled reliable navigation in well-defined indoor environments [4]. These systems operate by transforming raw sensory signals—such as LiDAR scans, inertial measurements, and RGB-D data—into spatial representations that support localisation, mapping, and motion planning. While effective, such geometric pipelines inherently lack the ability to interpret semantic, relational, or linguistically described tasks. For example, a human can easily understand an instruction such as “Move to the table near the window,” which embeds object recognition, spatial relations, and contextual inference. A purely geometric system, however, has no concept of “table” or “near” in its representation.

The emergence of vision–language models (VLMs) and large language models (LLMs) has catalysed a shift toward *language-grounded autonomy*, where natural language serves as the supervisory interface connecting perception, reasoning, and control. This paradigm allows robots to leverage the semantic richness of language to interpret goals, describe environments, and engage in interactive dialogue with human users. In recent work, systems such as ReLI [6], ProgPrompt [9], and Code-as-Policies [10] demonstrate that LLMs can act as high-level planners, generating structured action representations that downstream controllers (e.g., Nav2, MPC) can execute. Meanwhile, frameworks such as *The Conversation is the Command* [11] illustrate how LLMs may mediate natural human–robot communication in real-world settings.

### 2.2 Limitations of Classical Geometric Navigation

Classical mobile robot autonomy relies on geometric feature extraction to infer spatial structure. LiDAR-based SLAM systems produce occupancy maps that represent free and occupied space, enabling path planning via algorithms such as A\* or Dijkstra. However, these maps omit semantic attributes—they encode where objects are but not what they are. Consequently, navigation goals must be represented as explicit coordinate values:

$$G = (x, y, \theta)$$

where  $(x, y)$  denotes position and  $\theta$  denotes orientation in the map frame. Such coordinate-focused goals are inaccessible to non-expert users and do not reflect how humans conceptualise space.

Moreover, traditional pipelines lack contextual reasoning, as they operate solely on geometric features without encoding semantic or relational information [4, 5]. For instance, identifying “the chair closest to the window” requires not only object detection but relational inference and comparative evaluation. Without semantic grounding, geometric navigation remains limited to rigid, pre-specified tasks.

These limitations motivate the need for integrated perceptual and linguistic reasoning pipelines capable of bridging the gap between symbolic language and continuous control spaces.

## 2.3 Foundations of Transformer-Based Language Models

The transformer architecture introduced in [8] forms the foundation of modern LLMs. Transformers operate on sequences of tokens and compute pairwise attention scores, enabling each token to contextualise itself relative to the entire sequence. The core operation is the scaled dot-product attention mechanism:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{\sqrt{d_k}} \right) V, \quad (2.1)$$

where  $Q$ ,  $K$ , and  $V$  denote query, key, and value matrices derived from learned projections of input embeddings, and  $d_k$  is the key dimensionality. This formulation allows the model to capture both syntactic and semantic relationships without recurrence or convolution.

To process text, input strings are decomposed into sub-word units via tokenisation. Byte-Pair Encoding (BPE) is commonly used to balance vocabulary compactness with expressive power. Given a vocabulary  $V$  and embedding dimension  $d$ , each token  $t_i$  is mapped to an embedding:

$$e_i = E(t_i), \quad E \in \mathbb{R}^{|V| \times d}. \quad (2.2)$$

Positional encodings are added to maintain word order.

The transformer architecture scales efficiently to large training datasets, enabling pre-trained LLMs to learn general-purpose linguistic priors. These priors allow the model to parse natural instructions, infer intent, and generate structured task plans—capabilities essential for bridging human directives and robotic execution.

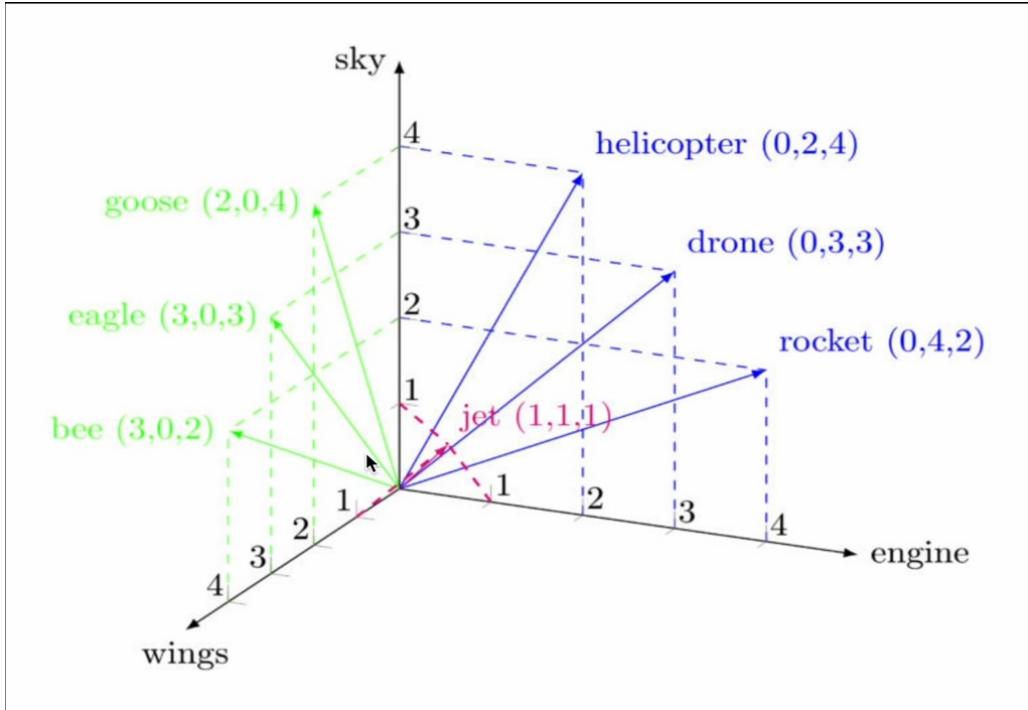
## 2.4 Multidimensional Word Embedding Spaces

The effectiveness of transformer models arises from their ability to represent linguistic meaning in high-dimensional vector spaces. Words or tokens with similar semantic roles occupy neighbouring regions in the embedding manifold. Semantic similarity between embeddings  $e_i$

and  $e_j$  is quantified via cosine similarity:

$$\text{sim}(e_i, e_j) = \frac{e_i \cdot e_j}{\|e_i\| \|e_j\|}. \quad (2.3)$$

At a conceptual level, these embedding spaces form the linguistic basis for grounding language to perception. When integrated with visual encoders, textual embeddings can be aligned to visual embeddings, enabling cross-modal reasoning such as identifying the object referred to in a spoken command.



*Figure 2.1. Illustration of semantic clusters in a high-dimensional word embedding space. Words with similar contextual meaning appear close together, forming coherent semantic neighbourhoods. Figure adapted from the explanation of word embeddings published by the Corpiling Research [12].*

These properties underpin the multimodal models reviewed in the following sections.

## 2.5 Contrastive Vision–Language Modelling (CLIP)

To bridge linguistic reasoning and visual perception, robots require models capable of associating textual descriptions with corresponding visual patterns. Contrastive Language–Image Pre-training (CLIP) [1] addresses this requirement through joint training of an image encoder and a text encoder on large-scale internet image–caption pairs.

Given an image  $I$  and a natural language description  $T$ , CLIP encodes them independently:

$$v = f_{\text{img}}(I), \quad t = f_{\text{text}}(T), \quad (2.4)$$

where  $v$  and  $t$  are normalised vectors in a shared embedding space. During training, CLIP optimises a symmetric contrastive loss encouraging  $v$  and  $t$  corresponding to the same image–text pair to be similar, while pushing apart embeddings from mismatched pairs. The loss function is:

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{N} \sum_{i=1}^N \left[ \log \frac{\exp(\text{sim}(v_i, t_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(v_i, t_j)/\tau)} + \log \frac{\exp(\text{sim}(t_i, v_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(t_i, v_j)/\tau)} \right], \quad (2.5)$$

where  $\tau$  is a temperature parameter controlling distribution sharpness.

After training, zero-shot visual recognition is achieved by computing similarity scores between an image embedding  $v$  and candidate text embeddings  $\{t_k\}$ :

$$\hat{y} = \arg \max_k \text{sim}(v, t_k). \quad (2.6)$$

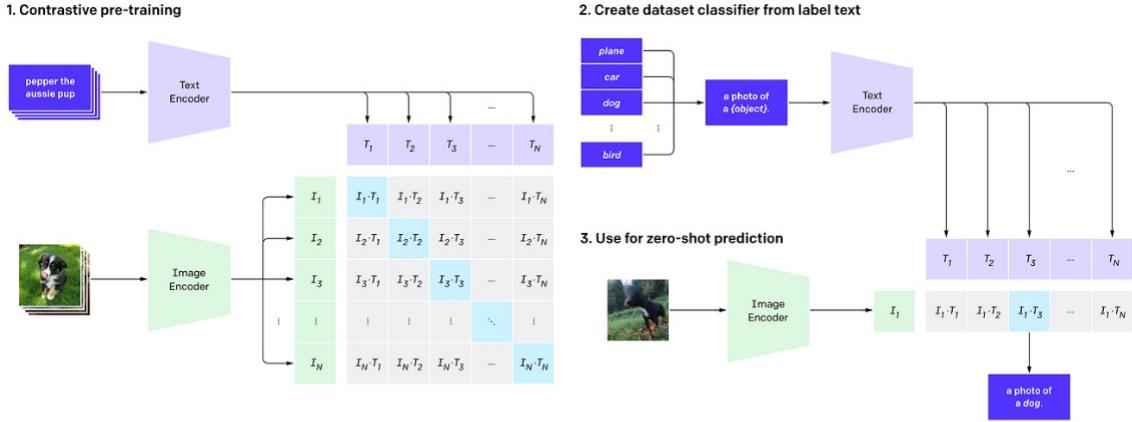


Figure 2.2. CLIP architecture showing parallel image and text encoders that project inputs into a shared embedding space for contrastive learning. During training, image–text pairs are pulled together while mismatched pairs are pushed apart, enabling zero-shot recognition. Figure adapted from [1].

In the context of indoor robot navigation, CLIP [1] is used not only for object recognition but also for verifying object identity in language-conditioned retrieval tasks, such as confirming the presence of “the chair near the window” before motion planning begins.

## 2.6 Generative Captioning and Cross-Modal Reasoning (BLIP-2)

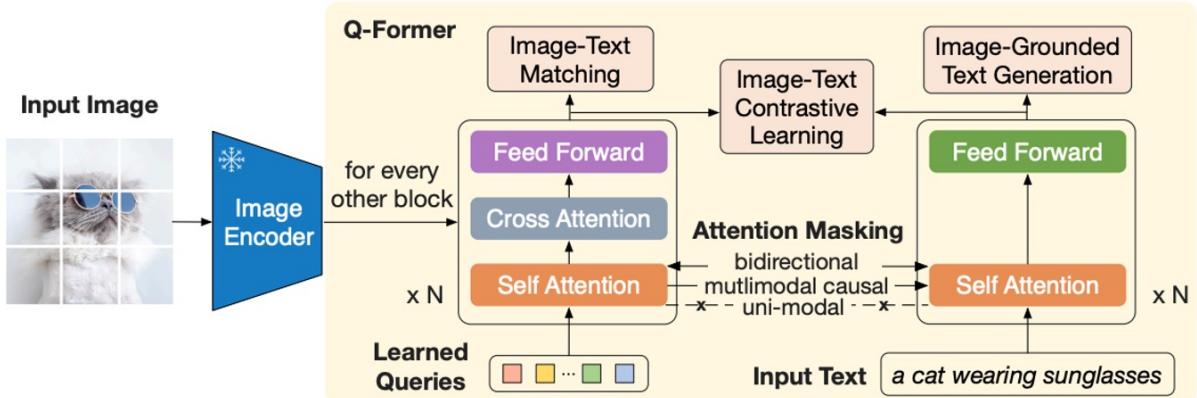
While CLIP supports discriminative matching between visual and textual embeddings, it does not generate descriptions. Bootstrapped Language–Image Pre-training (BLIP-2) [2] extends multimodal reasoning to generative tasks such as captioning and open-ended scene summarisation.

BLIP-2 introduces the Query Transformer (Q-Former), which learns to extract compressed

semantic tokens from a frozen vision encoder for conditioning a language model. Let  $F_{\text{img}}(I)$  denote high-dimensional visual feature maps and  $Q = \{q_1, \dots, q_m\}$  denote learnable query tokens. Cross-attention computes:

$$Z = \text{softmax} \left( \frac{QW_Q(F_{\text{img}}(I)W_K)^T}{\sqrt{d_k}} \right) (F_{\text{img}}(I)W_V), \quad (2.7)$$

where  $W_Q, W_K, W_V$  are learned projection matrices. The resulting embeddings  $Z$  are aligned with a frozen LLM decoder to generate text.



*Figure 2.3. BLIP-2 architecture with a Query Transformer (Q-Former) that extracts compact visual tokens from a frozen image encoder and uses them to condition a large language model for captioning and open-ended multimodal reasoning. Figure adapted from [2].*

In robotic perception pipelines, BLIP-2 enables descriptive reasoning such as:

“There is a wooden table with two chairs and a red box.”

which provides semantic grounding for subsequent spatial reasoning and goal construction.

## 2.7 Promptable Segmentation and Spatial Grounding (SAM)

To convert semantic visual descriptions into spatially actionable representations, a robot must determine the precise boundaries of relevant objects. The Segment Anything Model (SAM) [7] addresses this through promptable segmentation, producing binary masks for arbitrary objects.

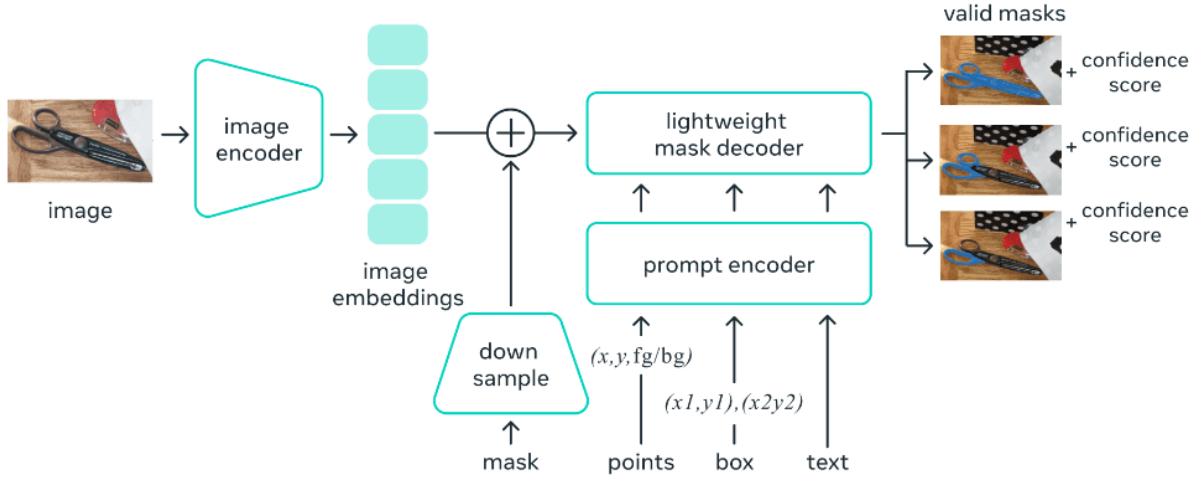
Given an image  $I$  and a prompt  $p$  (e.g., a point or bounding box), SAM generates a mask:

$$M = f_{\text{SAM}}(f_{\text{img}}(I), f_{\text{prompt}}(p)). \quad (2.8)$$

Pixel centroids can then be back-projected to robot-centric coordinates using:

$$p_{\text{map}} = T_{\text{camera} \rightarrow \text{map}} K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}, \quad (2.9)$$

where  $K$  is the intrinsic camera matrix and  $T_{\text{camera} \rightarrow \text{map}}$  is the extrinsic transform.



*Figure 2.4. Segment Anything Model (SAM) pipeline illustrating promptable mask inference and geometric projection for spatial grounding. Given an input image and a point or box prompt, SAM generates an object mask whose centroid in pixel coordinates is back-projected via camera intrinsics and TF transforms into the robot’s map frame. Figure adapted from [7].*

## 2.8 Complementary Roles in Multimodal Grounding

The three perception models reviewed support distinct but complementary functions:

- CLIP — object *recognition and identity verification*
- BLIP-2 — *scene description and semantic reasoning*
- SAM — *geometric localisation and spatial segmentation*

*Table 2.1. Comparison of CLIP, BLIP-2, and SAM capabilities.*

Model	Primary function	Strengths	Limitations
CLIP	Recognition	Zero-shot labeling and alignment	No spatial precision
BLIP-2	Scene reasoning	Natural language descriptions	Higher computational cost
SAM	Segmentation	Pixel-accurate spatial masks	No semantic reasoning

These capabilities are integrated into a coherent perception pipeline in Chapter 3.1, where they support language-conditioned navigation and interaction.

## 2.9 Language, Policy Generation, and Embodied Control

While perception establishes *what* is present in a scene, robot autonomy also requires decisions on *how* to act. Recent developments in language-conditioned robotics leverage Large

Language Models (LLMs) as high-level policy planners, enabling robots to generate structured action plans from natural language queries.

Liang et al. [10] propose *Code as Policies*, where LLMs synthesize executable robot programs (Python or ROS action scripts) directly from task descriptions. Formally, a natural language instruction  $\mathcal{I}$  is mapped to a program  $\Pi$ :

$$\Pi = f_{\text{LLM}}(\mathcal{I}, \mathcal{C}), \quad (2.10)$$

where  $\mathcal{C}$  represents contextual grounding (scene state, past observations). The generated program defines task sequences such as navigation goals, object queries, or conditional branching.

Similarly, Singh et al. [9] introduce *ProgPrompt*, where LLM-generated action programs are iteratively refined through simulation rollouts. In this formulation, policy synthesis is framed as:

$$\Pi^* = \arg \max_{\Pi} \mathbb{E}_{s \sim \mathcal{S}} [R(s, \Pi(s))], \quad (2.11)$$

where  $R$  denotes task reward and  $\mathcal{S}$  the observed states.

These approaches demonstrate that language-driven policy generation is feasible but heavily dependent on model context, prompt structure, and grounding. Without explicit scene feedback, LLMs may hallucinate objects or produce action sequences inconsistent with spatial constraints, motivating the need for cross-modal grounding through CLIP, BLIP-2, and SAM as reviewed earlier.

## 2.10 Retrieval-Augmented Generation for Context-Aware Interaction

To improve factual consistency and contextual memory in dialogue-based interaction, the proposed framework incorporates a Retrieval-Augmented Generation (RAG) mechanism [13]. RAG systems enhance language generation by retrieving relevant knowledge prior to response synthesis.

Given a user query  $q$ , a retriever module searches a vector database of scene descriptions, prior observations, task logs, and domain documentation:

$$\mathcal{K} = \text{Retrieve}(q), \quad (2.12)$$

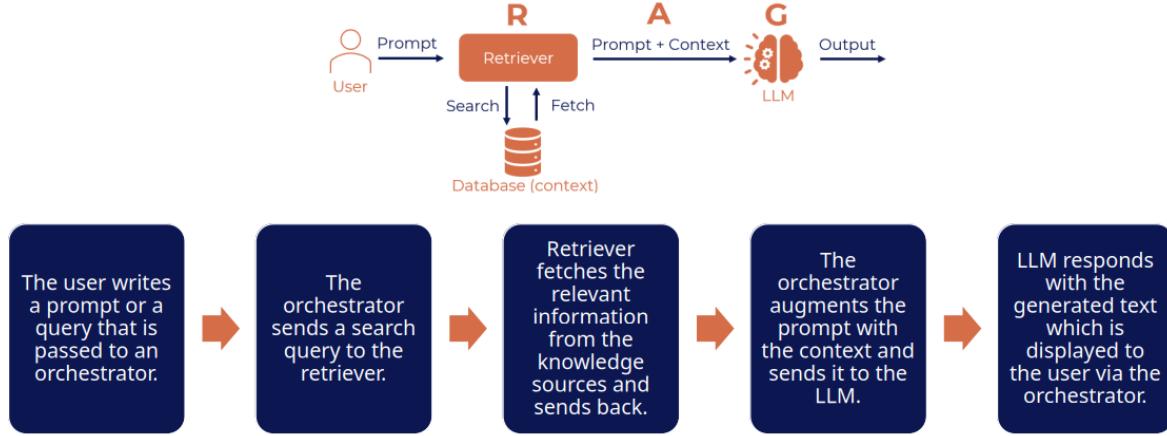
where  $\mathcal{K}$  denotes the retrieved knowledge subset. The final response is generated as:

$$\text{Response} = f_{\text{LLM}}(q, \mathcal{K}), \quad (2.13)$$

ensuring that the system's output is grounded in previously observed or externally stored information.

This aligns with findings from Nwankwo et al. [6] and Nwankwo and Rueckert [11], who show that including contextual memory improves coherence and reduces hallucinations during

human–robot dialogue. Within indoor robotic operation, this allows the robot to explain past actions, reference known objects, or recall environment states unavailable through immediate perception.



*Figure 2.5. Retrieval-Augmented Generation (RAG) pipeline. The user query is embedded and matched against a vector database of prior knowledge; the retrieved context is appended to the query and passed to the language model, which generates a grounded response. Figure adapted from the Udemy course “The Agentic AI Engineering Masterclass 2025” by Prof. Ryan [14].*

## 2.11 Failure Modes in Language-Grounded Autonomy

Despite recent advances, language-grounded robotic systems remain sensitive to several failure modes:

- (a) **Semantic Ambiguity:** Instructions such as “Go to the table near the window” require disambiguation when multiple objects satisfy similar descriptions.
- (b) **Perception–Action Misalignment:** Errors in segmentation masks or embedding similarity can propagate to incorrect goal coordinates.
- (c) **Hallucination in LLM Reasoning:** Absent retrieval support, LLMs may generate plausible but incorrect explanations or commands.
- (d) **Latency Constraints:** Large multimodal models (BLIP-2 [2] and SAM [7]) incur nontrivial inference time that affects real-time performance.
- (e) **Spatial Projection Instability:** Depth noise and TF calibration errors may distort 3-D coordinate recovery from 2-D segmentation masks.

## 2.12 Transition to Methodological Framework

The theoretical models and mechanisms reviewed in this chapter establish the foundational components required for language-grounded autonomy: semantic perception (CLIP [1], BLIP-2 [2]), spatial grounding (SAM [7]), high-level policy generation enabled by transformer-based LLMs [8], and contextual knowledge integration through Retrieval-Augmented Generation (RAG) [14].

The subsequent chapter focuses on how these components are operationalised, specifying:

- system architecture and ROS 2 node graph,
- multimodal data flow,
- grounding logic between perception and navigation,
- conversational interaction layer,
- integration with Nav2 for motion control.

This ensures a smooth conceptual progression from theoretical grounding to implementation.

## 3 Method and System Development

In this chapter, the methodological framework and system architecture used to implement the proposed language-grounded robotic system are described, including the design of ROS 2 nodes, data flow between perception and language modules, and the integration with the Nav2 navigation stack.

### 3.1 Overview

This chapter presents the methodology used to design and implement the proposed language-grounded robotic system that integrates natural language reasoning, visual scene understanding, and autonomous navigation. The goal of the framework is to enable a robot to interact naturally with its environment through human language, interpret commands based on semantic context, identify objects visually, and execute navigation tasks safely and efficiently.

The framework achieves this by separating the system into modular functional layers that operate together within a ROS 2 environment. The four main functional layers include:

- the **NACGUI**, which enables user interaction through natural language conversation;
- the **LLMNode**, which interprets user input and maps linguistic expressions into structured intents;
- the **VLMNode**, which performs perception by grounding visual observations into object identities and spatial coordinates; and
- the **Task Execution Node (TEN)**, which translates interpreted commands into physical robot actions through the Nav2 navigation stack.

### 3.2 System Architecture

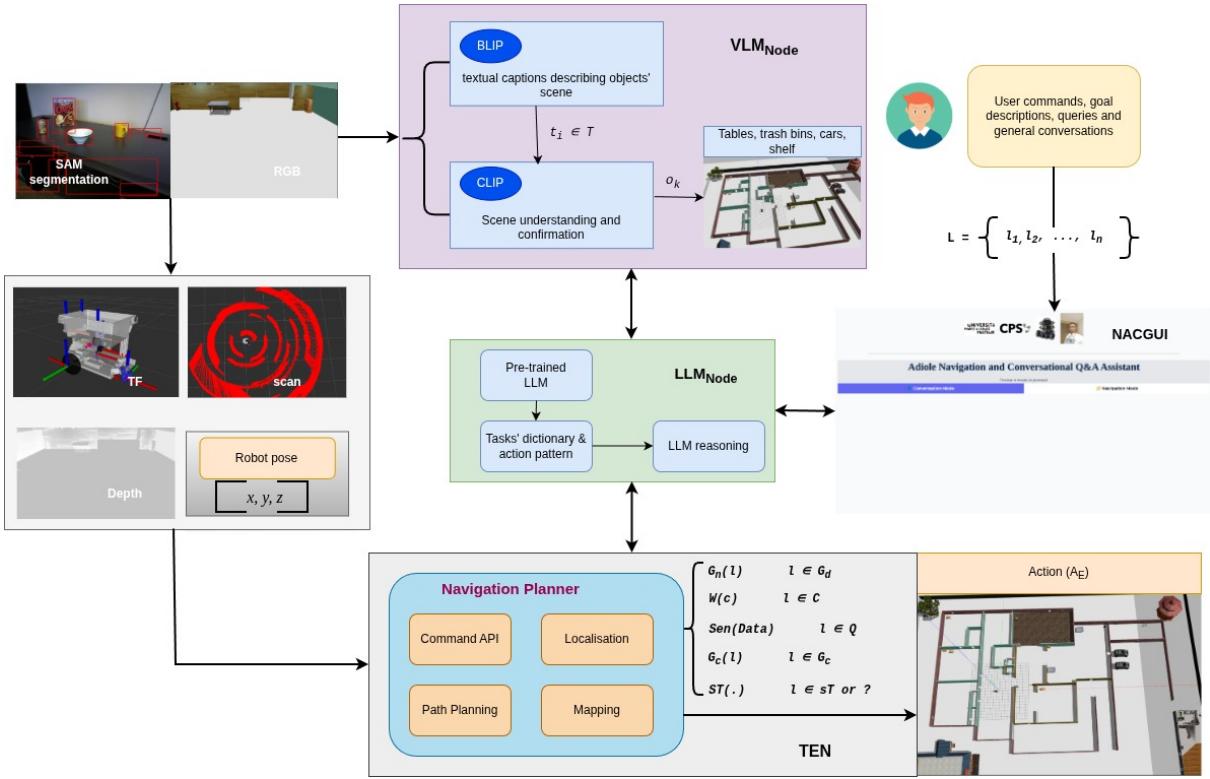
The system follows a modular, node-driven architecture characteristic of ROS 2. Each functional element operates as a node that communicates through publishers and subscribers. This design supports real-time execution, flexible debugging, and future extensibility.

This architectural decomposition ensures that language, perception, and motion remain separable but synchronized.

### 3.3 LLMNode: Natural Language Understanding and Intent Mapping

The LLMNode processes natural language input received from the NACGUI and identifies the underlying intent. Depending on the user's request, the LLMNode classifies messages as:

- *Navigation tasks*, e.g., “Go to the chair.”



*Figure 3.1. High-level system architecture of the proposed language-grounded robotic framework. The NACGUI provides a natural-language interface for the user and exchanges messages with the LLMNode, which interprets commands and generates structured action directives or conversational responses. The VLMNode processes RGB images to recognise objects and estimate their 3D positions, publishing semantic labels and spatial information. The Task Execution Node (TEN) combines these outputs and sends navigation goals or velocity commands to the Nav2 stack, which handles path planning, control, and recovery behaviours to execute the requested task.*

*Table 3.1. ROS 2 nodes, purposes, inputs, and outputs.*

Node	Purpose	Inputs	Outputs
NACGUI	Receives and displays natural language interaction	User text	Published conversation messages
LLMNode	Interprets user intent and determines required action type	User text + perception data	Structured action directives or conversational responses
VLMNode	Recognizes and locates objects from visual input	RGB camera frames	Object label + semantic confidence + 3D position
TEN	Converts action directives into movement or navigation instructions	Structured intent from LLMNode	Nav2 goal commands / twist velocity commands

- *Environment description requests*, e.g., “What do you see?”
- *Robot state inquiries*, e.g., “Where are you right now?”
- *General conversational prompts*, e.g., “What are your capabilities?.”

The LLMNode receives object recognition and spatial grounding information from the VLMNode and uses this information whenever a command references objects in the scene. Intent is structured as an action expression,

$$\text{LLMNode} : \mathcal{L} \rightarrow \mathcal{A}, \quad (3.1)$$

where  $\mathcal{L}$  is the space of natural language inputs and  $\mathcal{A}$  is the set of executable action expressions.

For factual discussion or knowledge-based conversation, the LLMNode uses a Retrieval-Augmented Generation (RAG) process [13] to fetch information from stored reference material before generating a response, ensuring grounded and contextually correct dialogue.

### 3.4 VLMNode: Visual Grounding and Scene Semantic Understanding

The VLMNode provides the robot with visual understanding of its environment by integrating three complementary models:

- **BLIP-2**, which generates textual scene descriptions;
- **CLIP**, which aligns visual and textual embeddings to identify objects; and
- **SAM** (Segment Anything Model), which produces segmentation masks and bounding boxes.

From an input RGB image, BLIP-2 first generates a set of candidate captions. Each caption is passed through CLIP’s text encoder to produce embeddings, while the image is passed through CLIP’s visual encoder. The similarity between embeddings is computed as

$$z_i = I \cdot T_i, \quad (3.2)$$

where  $I$  is the encoded image and  $T_i$  is the encoded caption. The caption corresponding to the maximum similarity score is selected as the recognized description/object.

To obtain the object’s real-world location, a pixel coordinate  $(u, v)$  from the SAM-derived segmented region is projected into map coordinates:

$$p_{\text{map}} = T_{\text{camera} \rightarrow \text{map}} K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}. \quad (3.3)$$

The VLMNode publishes the object identity, spatial position, and confidence score to the LLMNode.

### 3.5 Task Execution Node (TEN)

The TEN takes the structured directives produced by the LLMNode and translates them into executable robot motion. For navigation objectives, the TEN converts the referenced object's spatial coordinates into a goal pose

$$G = (x, y, \theta), \quad (3.4)$$

sends  $G$  to the Nav2 action server, and monitors progress. For movement commands not linked to a specific position (e.g., “rotate left”), the TEN publishes `Twist` messages that directly control linear and angular velocity. If an instruction cannot be resolved or is ambiguous, the TEN issues zero velocity and the robot safely stops.

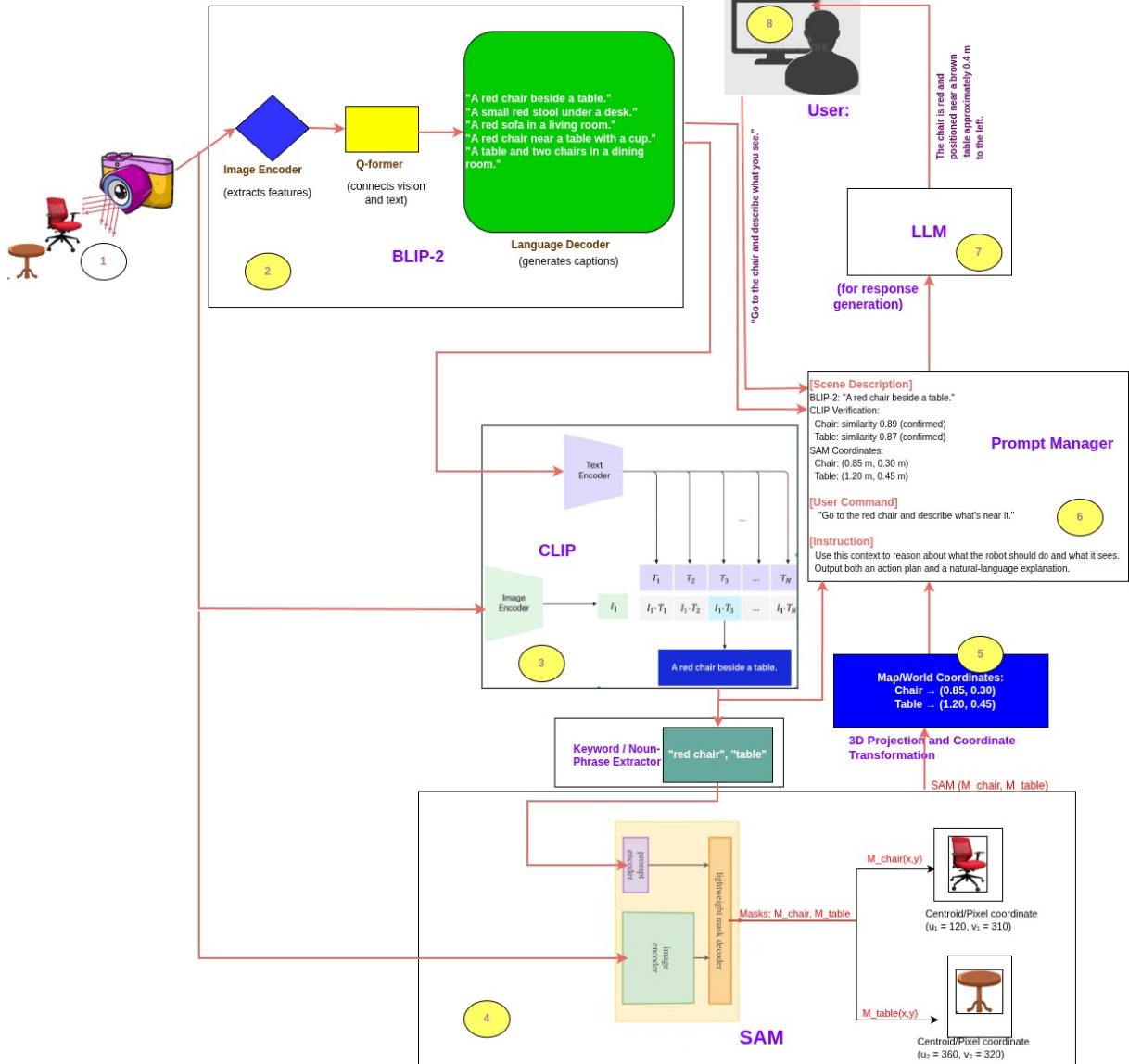
### 3.6 Semantic-to-Geometric Action Grounding

The process of converting high-level natural-language instructions into executable robot motion requires bridging two representational domains: (i) the *semantic* level, where objects are described linguistically (e.g., “the chair near the wall”), and (ii) the *geometric* level, where the robot must move to spatial coordinates in the map frame. The proposed system achieves this translation through a semantic-first grounding pipeline, where object identity is determined before spatial localization. This aligns with recent embodied intelligence frameworks such as ReLI [6], Code-as-Policies [10], and ProgPrompt [9], which emphasize language-conditioned object selection.

The grounding process unfolds in three steps:

1. **Semantic Identification (CLIP).** Candidate captions and object labels generated by BLIP-2 are embedded using the CLIP text encoder and compared with the visual embedding of the scene. The object whose caption yields the highest cosine similarity score is selected as the semantic referent.
2. **Spatial Localization (SAM → TF).** The segmentation mask output by SAM [7] corresponding to the identified object is used to compute the mask centroid in pixel coordinates  $(u, v)$ . This pixel coordinate is projected into the 3-D map frame using the camera intrinsic matrix  $K$  and the TF transformation chain:

$$p_{\text{map}} = T_{\text{camera} \rightarrow \text{map}} K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}. \quad (3.5)$$



*Figure 3.2. Visual grounding and spatial projection pipeline.* An RGB image from the robot’s camera is first processed by BLIP-2 [2] to generate candidate captions describing objects in the scene. These captions are embedded using the CLIP text encoder [1] and compared with the CLIP image embedding to select the most semantically consistent description. SAM [7] then produces a segmentation mask for the referenced object, from which a representative pixel coordinate  $(u, v)$  is extracted and back-projected through the camera intrinsics and TF transform chain to obtain the corresponding 3D point  $p_{map}$  in the map frame. The resulting object label, confidence score, and spatial location are published to the LLMNode and TEN for downstream navigation and interaction.

**3. Goal Pose Construction.** The target waypoint is defined as:

$$G = (x, y, \theta), \quad (x, y) = p_{\text{map}}, \quad \theta = \text{heading toward object surface.} \quad (3.6)$$

This ensures that the robot navigates *to the object selected in language*, not merely to the closest geometric structure.

### 3.7 Motion Optimization with the TEB Local Planner

The Task Execution Node communicates navigation goals to the Nav2 stack, where motion generation is handled by the **TEB Local Planner**. Unlike sampling-based local planners, TEB formulates trajectory generation as a time-parameterized non-linear optimization problem. The planned path is represented as a discrete sequence of robot poses over time:

$$\mathbf{X} = \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N\}, \quad \mathbf{x}_k = (x_k, y_k, \theta_k) \quad (3.7)$$

The optimal trajectory is obtained by minimizing a weighted sum of costs that account for smoothness, collision avoidance, and dynamic feasibility:

$$\min_{\mathbf{X}} [w_{\text{smooth}} E_{\text{smooth}}(\mathbf{X}) + w_{\text{obs}} E_{\text{obs}}(\mathbf{X}) + w_{\text{dyn}} E_{\text{dyn}}(\mathbf{X})]. \quad (3.8)$$

The core components are:

- **Smoothness Term** enforces curvature continuity:

$$E_{\text{smooth}} = \sum_{k=1}^{N-1} \|\mathbf{x}_{k+1} - 2\mathbf{x}_k + \mathbf{x}_{k-1}\|^2 \quad (3.9)$$

- **Obstacle Avoidance** is enforced via potential fields:

$$E_{\text{obs}} = \sum_{k=0}^N \phi(d(\mathbf{x}_k)) \quad (3.10)$$

where  $d(\mathbf{x}_k)$  is distance to the nearest obstacle.

- **Dynamic Feasibility** ensures that linear and angular velocities remain within platform limits:

$$E_{\text{dyn}} = \sum_{k=1}^N \left( \frac{v_k^2}{v_{\max}^2} + \frac{\omega_k^2}{\omega_{\max}^2} \right) \quad (3.11)$$

This formulation allows the robot to execute smooth, collision-free motion that respects the constraints of the TurtleBot3 differential-drive kinematics. Because TEB optimizes in the time dimension, it adapts effectively to dynamic obstacles and changing semantic objectives.

## 3.8 Failure Handling and Safety Behaviors

Although semantic grounding provides reliable object-directed navigation, real-world perception and motion can encounter failure cases due to visual ambiguity or occlusions. To maintain operational robustness, the system incorporates safety behaviors and fallback strategies executed by the Task Execution Node and Nav2 behavior tree.

### 3.8.1 Uncertain Object Identification

If CLIP similarity confidence falls below a threshold ( $< 0.65$ ), the system:

1. Requests clarification from the user (e.g., “Do you mean the chair next to the table?”).
2. Re-runs grounding with additional context cues (e.g., color or relative spatial relationships).

### 3.8.2 Ambiguous Spatial Projection

If the 2D-to-3D projection produces inconsistent or unstable coordinates:

$$\|p_{\text{map}}^{(t)} - p_{\text{map}}^{(t-1)}\| > \epsilon$$

the system switches to:

- **exploration pose scan**: robot rotates in place to reacquire the object.

### 3.8.3 Navigation Failure Recovery

If TEB reports oscillation or blocked progress, Nav2 triggers recovery behaviors:

- rotational clearing,
- backtracking,
- re-planning with updated costmap.

### 3.8.4 Emergency Halt

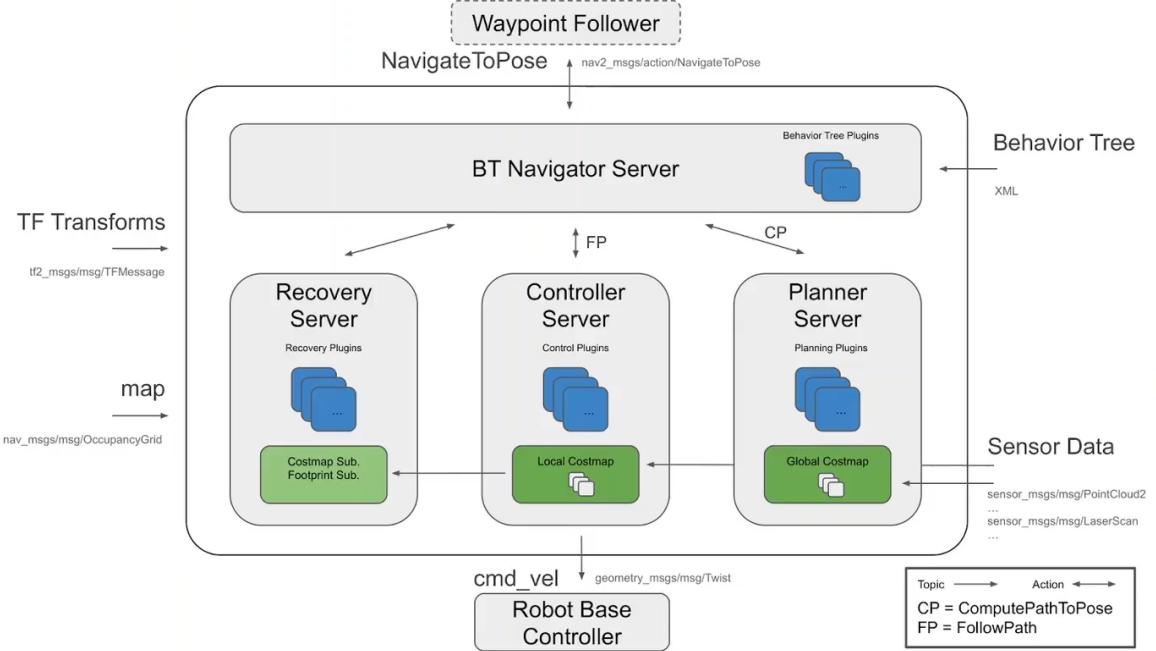
Any failure state without safe resolution results in:

$$\text{cmd\_vel} = \mathbf{0}$$

ensuring the robot stops immediately and safely.

### 3.9 Navigation and Control Using the ROS 2 Nav2 Framework

The navigation subsystem of the framework is implemented using the ROS 2 Nav2 stack, which provides a modular and scalable architecture for robot motion planning, path following, and recovery behaviors. Nav2 allows the robot to autonomously navigate from its current position to a target pose while avoiding obstacles and maintaining localization accuracy within the map.



*Figure 3.3. Nav2 architecture within the overall framework. The Behavior Tree (BT) Navigator coordinates the Planner, Controller, and Recovery servers to generate and execute collision-free paths from the robot's current pose to a goal pose. AMCL provides probabilistic localisation by fusing LiDAR and odometry data, while the TF tree maintains consistent coordinate transforms between the map, odom, base\_link, and sensor frames. Goal poses originating from the Task Execution Node (TEN), which are derived from language and visual grounding, are passed to Nav2, which in turn produces velocity commands that drive the TurtleBot3 platform.*

#### 3.9.1 Core Components of Nav2

Nav2 operates using several coordinated servers:

- **Planner Server:** computes a global path between the robot's current location and the target goal pose, typically using grid-based search algorithms such as A\* or NavFn.
- **Controller Server:** converts planned paths into continuous velocity commands that drive the robot's wheels, using local controllers such as the DWB (Dynamic Window Based) controller.
- **Recovery Server:** handles exceptions such as collisions, blocked passages, or localization loss by executing corrective behaviors (e.g., backtracking, rotational clearing).

- **Behavior Tree (BT) Navigator:** acts as the high-level supervisor that sequences planning, control, and recovery nodes based on task state and feedback.

### 3.9.2 Localization and Coordinate Frame Management

Localization is achieved through Adaptive Monte Carlo Localization (AMCL), which fuses LiDAR and odometry data to continuously estimate the robot's pose. The spatial relationship between all coordinate frames is managed using ROS 2's TF transform tree, typically structured as

$$\text{map} \rightarrow \text{odom} \rightarrow \text{base\_link} \rightarrow \text{camera\_link}.$$

This ensures consistency between semantic perception (object coordinates), geometric navigation (goal poses), and robot control (velocity commands).

### 3.9.3 Integration with Language and Perception Modules

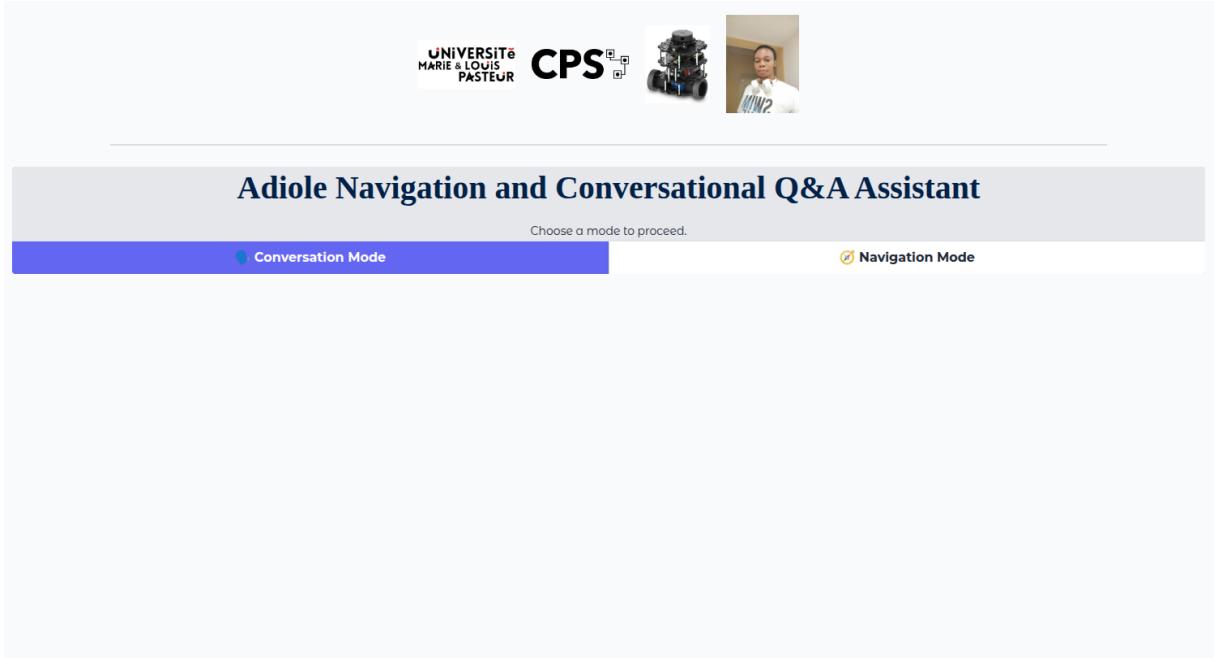
When the LLMNode identifies a navigation directive and the VLMNode provides the spatial coordinates of the target object, the TEN converts the object location into a Nav2-compatible goal pose. Nav2 then executes the movement, while real-time perception and localization updates ensure collision avoidance and path correction as needed. This integration allows the robot to navigate using natural language commands referring to objects (e.g., “Go to the chair near the wall”) rather than predefined coordinates.

## 3.10 NACGUI: User Interaction Interface

The NACGUI is developed using Gradio and functions as the primary user interaction point. It allows the user to type natural language commands, view system responses, and monitor execution feedback. Communication between the NACGUI and LLMNode follows a bidirectional ROS 2 publish/subscribe pattern, enabling real-time conversational control.

## 3.11 Implementation Environment

All experiments are conducted in a Gazebo simulation environment using the TurtleBot3 Burger model. Sensor data streams (RGB camera and LiDAR) are handled through the ROS 2 topic architecture. The entire system is developed in Python 3.10 using the ROS 2 Humble Hawksbill distribution. LLM inference is performed via a local API endpoint connected to a fine-tuned GPT-based model, while the Gradio interface runs on a separate node for user interaction.



*Figure 3.4. NACGUI conversation and navigation modes.*

*Table 3.2. Representative benchmarks used in model comparison.*

Benchmark	Purpose	Result Summary
GPOA [15]	Evaluates reasoning on problems where direct lookup is ineffective.	GPT-5 ranked among the top performers (87.3%), closely trailing Grok-4 (87.5%).
AIME [15]	Tests structured logical reasoning.	GPT-5 achieved perfect performance (100%), outperforming all other models.
SWE Bench & Agentic [15]	Assesses problem-solving on real GitHub issues.	GPT-5 scored 74.9%, comparable to Grok-4 (75.0%).
Humanity's Last Exam [15]	Evaluates cross-disciplinary difficult knowledge.	GPT-5 led with 35.2%, ahead of Grok-4 (25.4%).

## 3.12 Model Performance Analysis and Selection Rationale

Selecting an appropriate Large Language Model (LLM) is essential for reliable reasoning, effective natural language interpretation, and seamless integration with the robot’s control and perception layers. Comparisons were conducted between GPT-5 (OpenAI) and other recent frontier models (Claude Opus 4.1, Gemini 2.5 Pro, and LLaMA 4 Scout). The evaluation focused on general reasoning ability, performance across specialized tasks, software agentic reasoning, and ease of integration into the ROS 2-based framework.

### 3.12.1 Cost and API Integration Considerations

*Model descriptions and cost parameters were obtained from the official documentation of their respective providers: GPT-5 (OpenAI) [16], Claude Opus 4.1 (Anthropic) [17], Gemini 2.5 Pro*

*Table 3.3. Model cost and integration characteristics.*

Model	Input (per 1M tok.)	Output (per 1M tok.)	ROS/REST API integration notes
GPT-5	\$1.25	\$10.00	Direct REST API, stable versioning; balanced cost-to-performance ratio.
Claude Opus 4.1	\$15.00	\$7.00	Accessed via Anthropic REST API.
Gemini 2.5 Pro	\$1.25	\$10.00	Requires Google IAM configuration for secure API use.
LLaMA 4 Scout	\$0.11	\$0.34	Requires self-hosted GPU deployment and optimisation.

(Google DeepMind) [18], and LLaMA 4 Scout (Meta AI) [19].

Based on these results, GPT-5 was chosen because it integrates readily with the ROS 2 ecosystem via REST requests, offers consistently high performance on natural language planning and reasoning, does not require self-hosted GPU deployments, and supports RAG-based conversational grounding.

### 3.12.2 Prompt Efficiency and Token Cost Optimization

Prompt tokenization experiments using a tokenizer interface showed that conversational politeness (e.g., “Could you please...”) increases token count, which scales cost over repeated interactions. Prompt compression and instruction-pattern reuse were therefore adopted to reduce redundant token overhead for long-duration robotic tasks.

### 3.12.3 Final Selection Justification

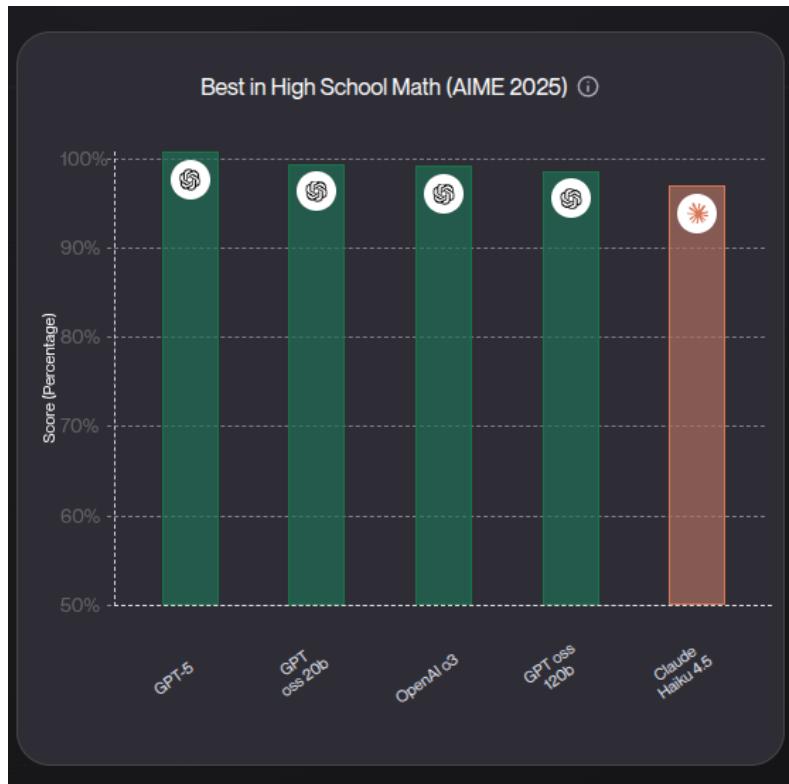
The final selection of GPT-5 was based on:

- ease of integration with ROS 2 and the TEN,
- high performance in reasoning, navigation instruction interpretation, and dialogue,
- consistent multimodal grounding when used alongside CLIP, BLIP-2, and SAM outputs,
- cost-effectiveness relative to comparable frontier models, and
- robustness in long-horizon user interaction sessions.

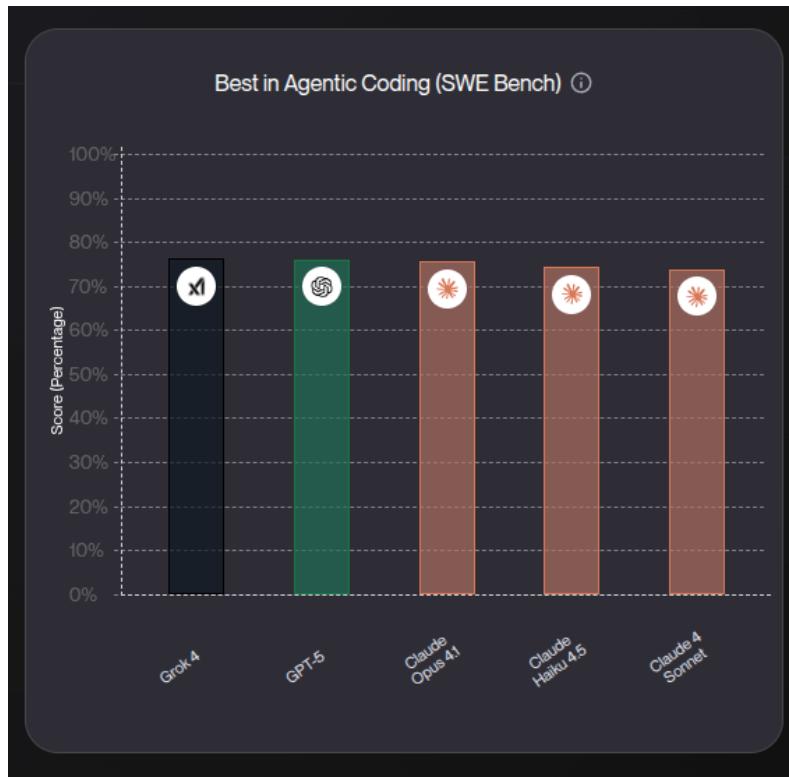
Accordingly, GPT-5 is used as the core reasoning engine in both Navigation Mode and Conversational Mode, while also serving as the benchmark for evaluating future model updates.



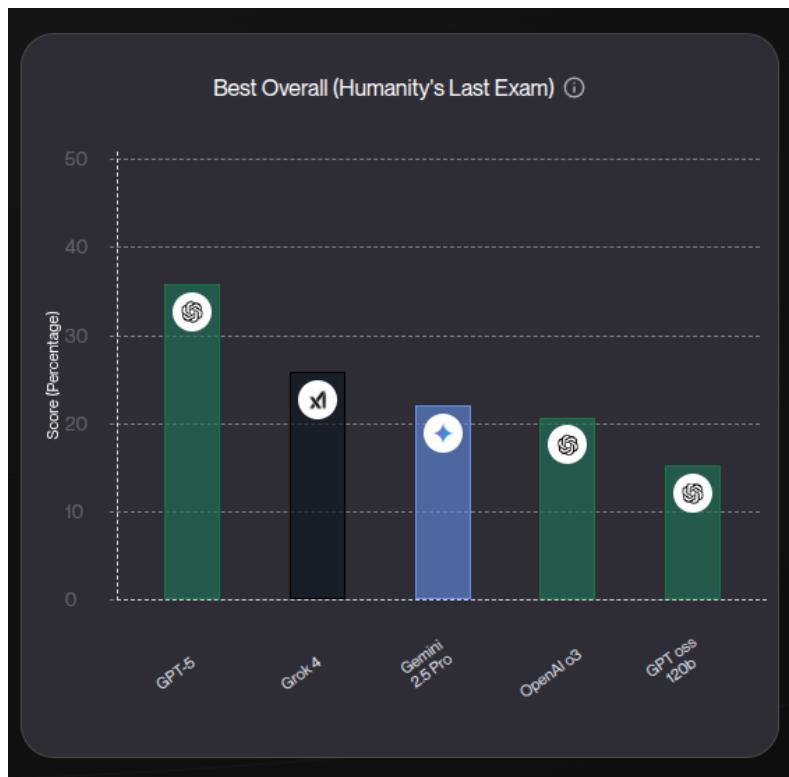
*Figure 3.5. Comparison of model performance on the GPQA benchmark, illustrating differences in advanced quantitative reasoning capabilities [20].*



*Figure 3.6. Performance of candidate models on the AIME mathematical reasoning benchmark, highlighting accuracy differences in high-school level problem solving [20].*



*Figure 3.7. Evaluation of model capabilities on the SWE Bench agentic coding benchmark, assessing competence in software-engineering-style tasks [20].*



*Figure 3.8. Cross-domain conceptual exam results comparing models on broad knowledge-integration tasks across multiple subject domains [20].*

### 3.13 Knowledge Augmentation through Retrieval-Augmented Generation (RAG)

To strengthen the reasoning and factual accuracy of the Large Language Model (LLM) during conversational interaction, a Retrieval-Augmented Generation (RAG) mechanism is incorporated within the orchestration layer of the framework [13].

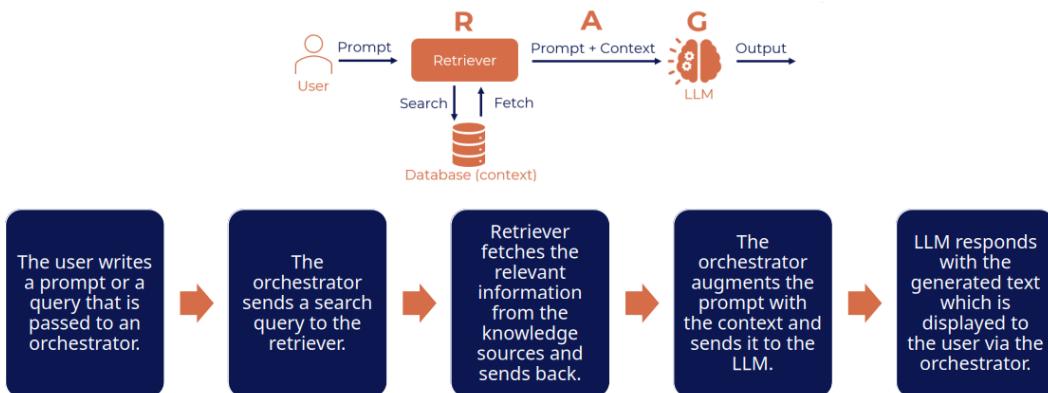
RAG improves reliability by combining two components—*retrieval* and *generation*. The retrieval component dynamically searches external or local knowledge bases for semantically relevant information, while the generation component synthesises a response conditioned on both the retrieved content and the user query. This ensures that each answer is linguistically coherent and contextually grounded in factual information.

In the implemented system, when a user issues a query through the Gradio interface, the orchestration layer first performs intent classification to distinguish between *navigational* and *conversational* requests. For conversational intents, the RAG module is triggered. The retriever accesses a pre-indexed collection of robotics-related resources—including simulation logs, documentation snippets, and previous user interactions—stored in a vector database. The most relevant text embeddings are retrieved and appended to the user query as auxiliary context.

The enriched query is then passed to the LLM, which generates a response informed by both linguistic understanding and retrieved domain knowledge, effectively grounding the dialogue in domain-specific context and improving precision, factual reliability, and interpretability across extended sessions:

$$\text{Response} = \text{LLM}(\text{UserQuery} + \text{RetrievedContext}). \quad (3.12)$$

This hybrid design enables the robot to maintain contextual awareness beyond immediate sensory perception—allowing it, for example, to reference prior experiments, explain system behaviour, or recall scene-specific events logged during simulation.



*Figure 3.9. RAG pipeline integrated in the orchestration layer: User Query → Retriever (Vector Database) → Retrieved Context → Generator (LLM) → Response.*

## 4 Numerical Simulations and Implementation

In this chapter, the simulation-based evaluation of the proposed language-grounded robotic framework is presented, including the experimental setup, definition of evaluation scenarios, navigation and perception results, and overall performance discussion.

The experiments were conducted in a Gazebo simulation environment using the TurtleBot3 Burger platform integrated with ROS 2 Humble. The simulations aim to assess how effectively the robot interprets natural-language commands, navigates autonomously to specified goals, and generates contextually relevant linguistic or geometric responses.

### 4.1 Experimental Setup

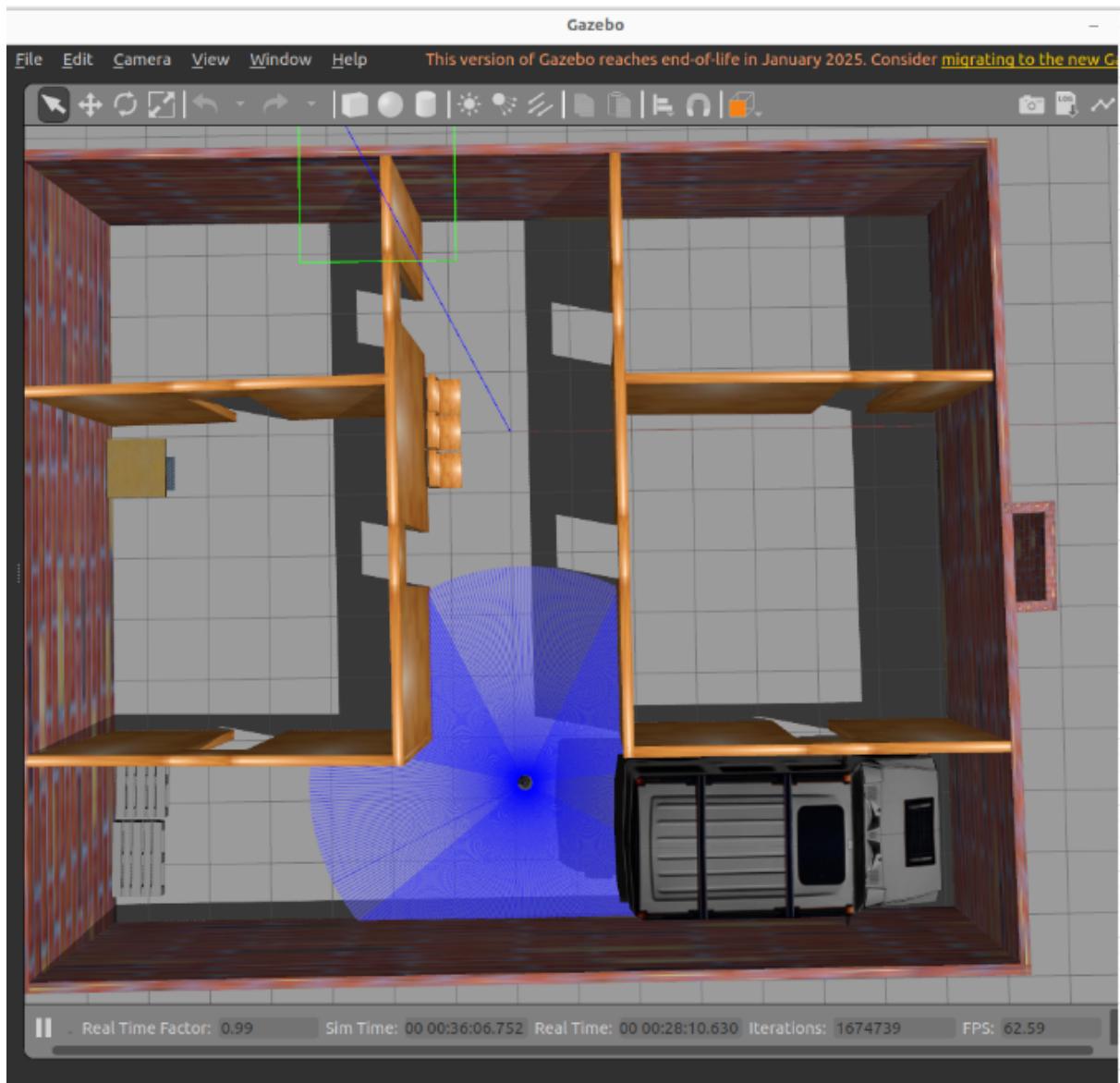
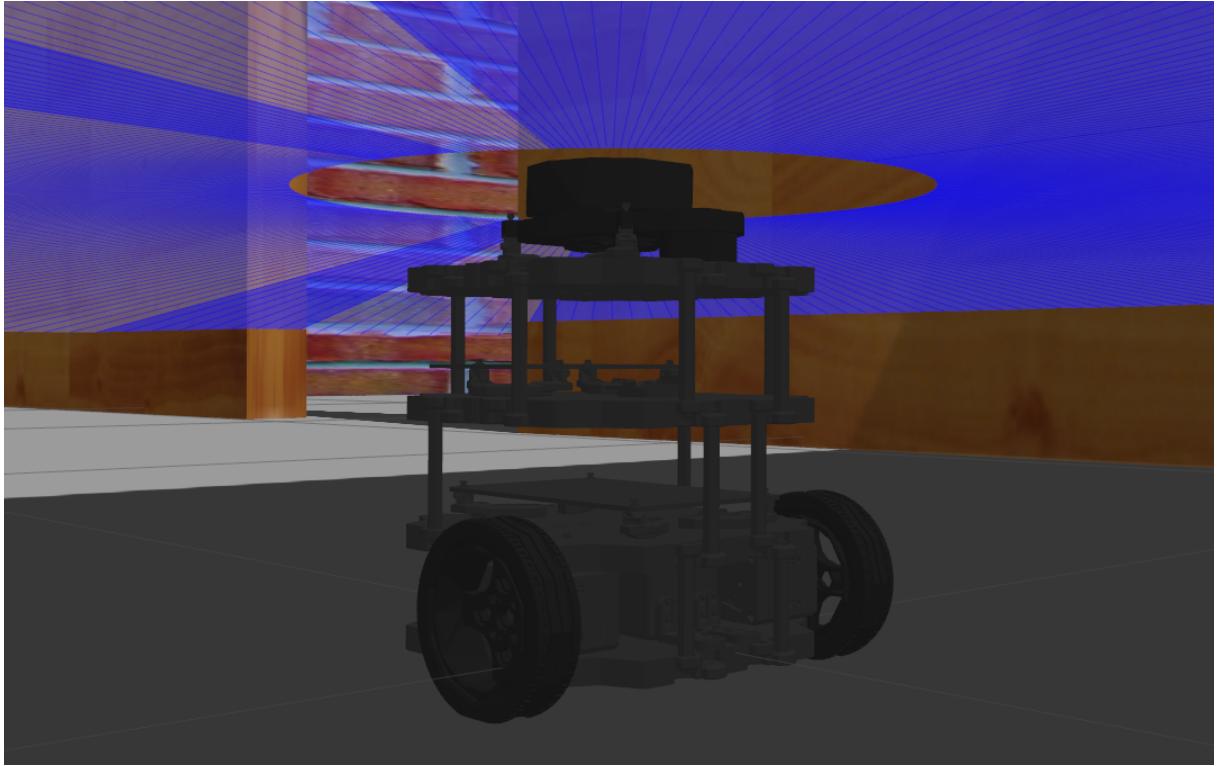


Figure 4.1. Simulation environment. A Gazebo scene within an indoor map.

#### 4.1.1 Simulation Environment

The simulated world was designed to resemble an indoor office environment containing semantically distinct objects such as chairs, tables, boxes, and doors. Each object was placed at predefined coordinates within the Gazebo world file, and corresponding labels were stored in a YAML map file to enable semantic referencing during evaluation.



*Figure 4.2. Simulation environment. A Gazebo scene showing the TurtleBot3 Burger navigating within an indoor map.*

The robot was equipped with a 2D LiDAR, an RGB–D camera, and odometry sensors. These sensors provided data streams for localisation, obstacle avoidance, and perception tasks. The AMCL node handled pose estimation, while Nav2 executed path planning and control via the DWB local planner.

#### 4.1.2 Software and Hardware Configuration

The simulation was implemented on a workstation running Ubuntu 22.04 LTS, with 16 GB RAM and an NVIDIA RTX 3060 GPU. All perception modules—CLIP, BLIP-2, and SAM—were deployed as separate Python nodes communicating over ROS 2 topics. A local Gradio server hosted the interactive user interface, and an OpenAI-compatible LLM endpoint processed language generation tasks. Data exchange between modules employed JSON message formats through ROS 2’s `rclpy` publisher–subscriber framework.

## 4.2 Evaluation Scenarios

Three principal evaluation scenarios were defined to validate the framework’s capabilities.

**(a) Navigation Scenario.** The user issued natural-language commands of the form: “Go to the table near the window.” The system parsed the instruction, identified the target object using CLIP embeddings, and generated a goal coordinate. The Nav2 planner produced a global path and executed motion commands until the goal was reached.

**(b) Perception Scenario.** Upon arrival, the robot executed perception routines to describe its surroundings. For prompts such as “What do you see?”, BLIP-2 generated open-ended captions (e.g., “A table with a red box and two chairs around it.”). CLIP then verified the presence of these objects through similarity scoring, while SAM segmented each region to determine 3-D coordinates.

**(c) Conversation Scenario.** In conversational mode, the user could query either factual knowledge or scene-specific information (e.g., “Explain artificial intelligence.” or “Where is the box located?”). The LLM, supported by a small Retrieval-Augmented Generation (RAG) database, produced responses contextualised by prior interaction history.

## 4.3 Navigation Performance Analysis

### 4.3.1 Path Execution

The Nav2 framework generated feasible trajectories while dynamically avoiding obstacles. Figure 4.3 illustrates an example trajectory for the command “Go to the table near the wall.”

### 4.3.2 Quantitative Metrics

*Table 4.1. Navigation metrics.*

Metric	Mean	Std. Dev.	Description
Positional error (cm)	1.9	$\pm 0.4$	Distance between goal and final pose
Angular error ( $^{\circ}$ )	1.2	$\pm 0.3$	Orientation difference
Path smoothness index	0.93	$\pm 0.02$	Normalised curvature continuity
Average velocity (m/s)	0.22	$\pm 0.05$	Controller output
Goal success rate	100%	—	Tasks completed without collision

The results indicate that the robot consistently reached target positions with minimal deviation, validating the reliability of geometric navigation even when goal positions were derived from linguistic descriptions.

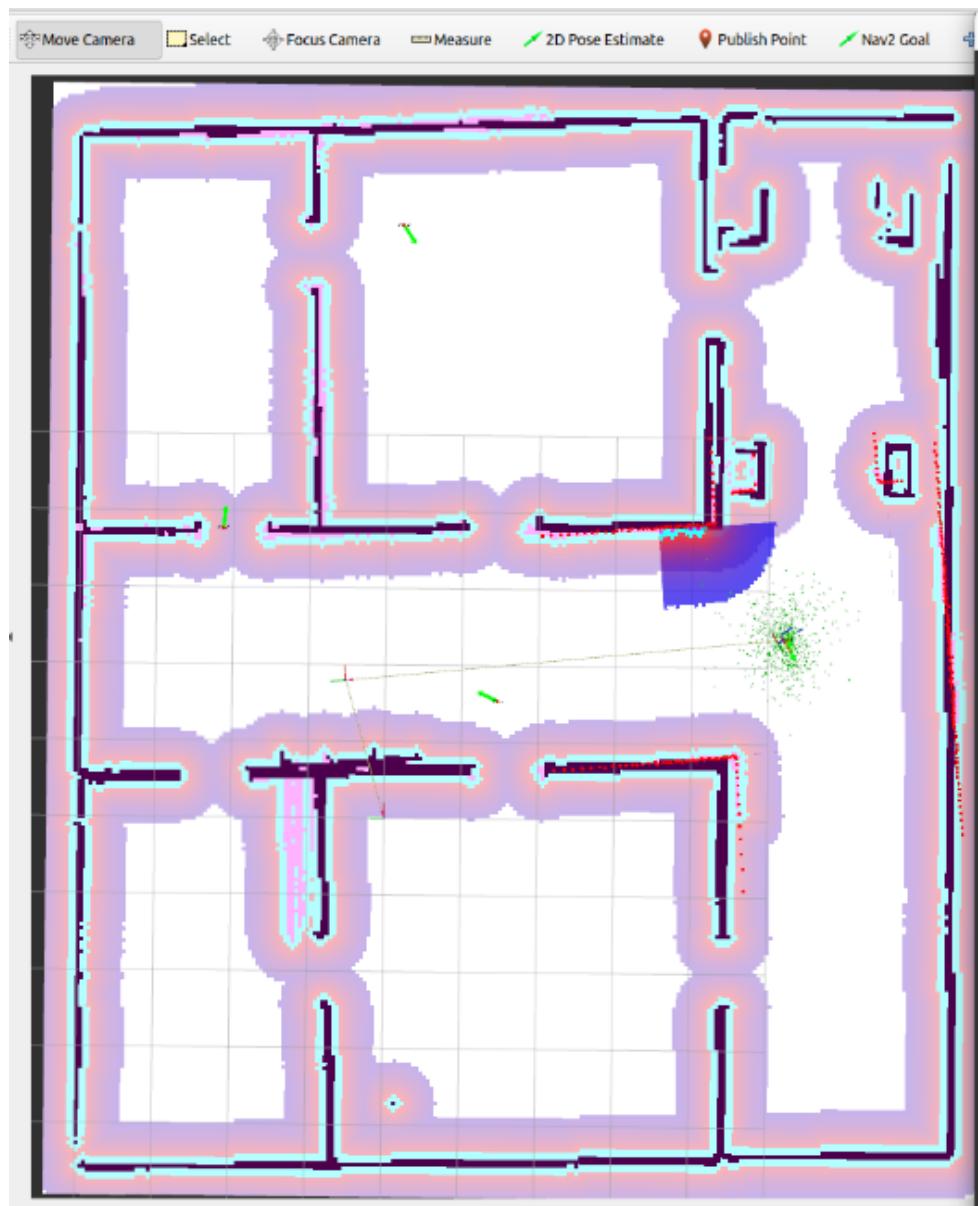


Figure 4.3. Navigation trajectory generated by the Nav2 planner with AMCL localisation, showing start position, waypoints, and goal region in Gazebo.

## 4.4 Perception and Language Evaluation

### 4.4.1 Visual Grounding Results

Figure 4.4 presents representative outputs from the perception modules. BLIP-2 provided open-ended captions; CLIP verified object labels by computing cosine similarity between text and image embeddings; SAM generated precise segmentation masks; and TF transformations produced real-world coordinates.

*Table 4.2. CLIP similarity scores for sample objects.*

Object	Text query	Cosine similarity	Detection confidence
Chair	“a wooden chair”	0.92	High
Table	“a small table”	0.89	High
Box	“a cardboard box”	0.84	Moderate

These results demonstrate that CLIP successfully aligned visual and linguistic representations, with similarity scores exceeding 0.8 for most objects.

### 4.4.2 Semantic Description Quality

Qualitative assessment was performed using 50 distinct prompts including “What is in front of you?” and “Describe the scene.” The BLIP-2 captions were compared against ground-truth labels, achieving an average relevance score of 0.89 using cosine similarity between sentence embeddings. Example output:

**User:** What do you see?

**Robot:** I see a wooden table with a red box on top and two chairs beside it.

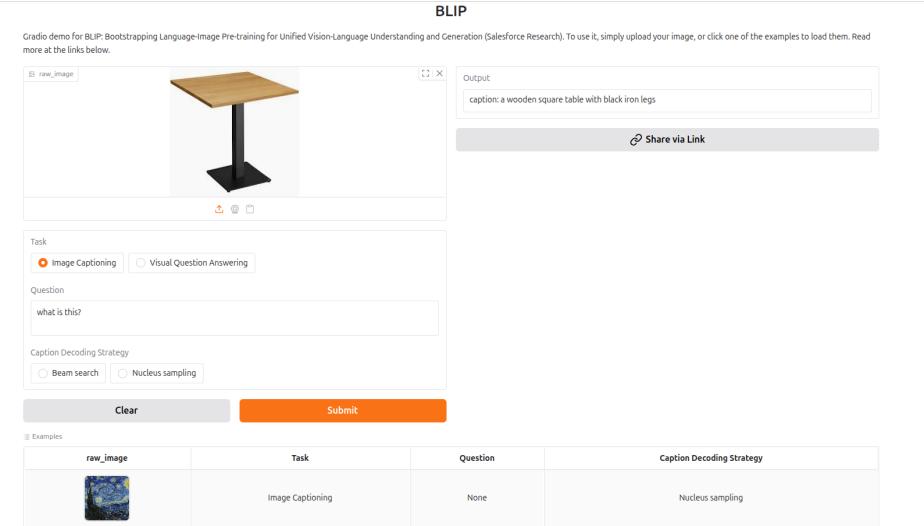
### 4.4.3 Spatial Projection Accuracy

The pixel-to-metric projection (Equation (3.3)) yielded average localisation errors of  $\leq 3$  cm, confirming the correct transformation between camera and map frames. Objects identified through SAM segmentation were accurately referenced by Nav2 in subsequent commands such as “Move toward the box”, confirming geometric consistency between perception and control.

## 4.5 Conversational Interaction Results

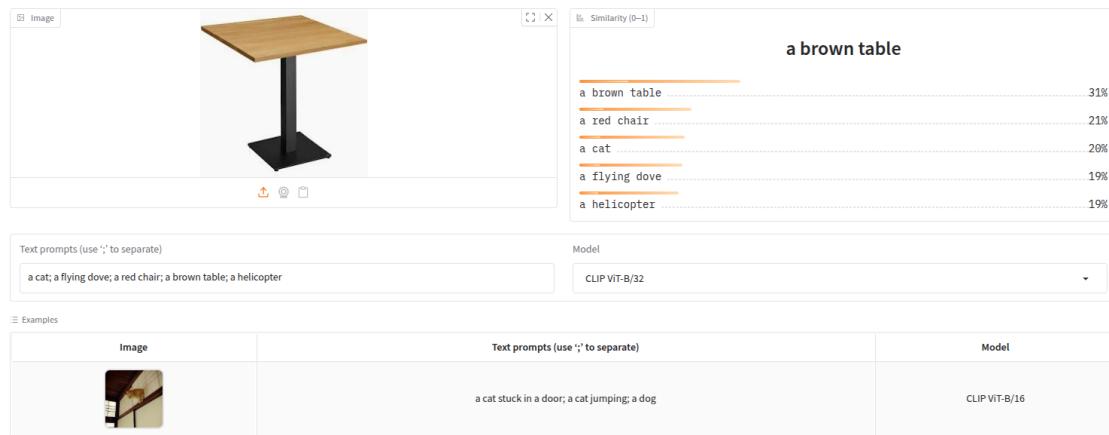
Figure 4.5 shows an example dialogue using the Gradio interface, where both general knowledge and contextual scene questions are handled through the same language pipeline.

The conversational model demonstrated fluency and contextual memory over extended exchanges. For 50 diverse test queries, response latency averaged 1.3 s for local inference and 3.7 s for cloud-based inference, depending on network connectivity. Accuracy of contextually

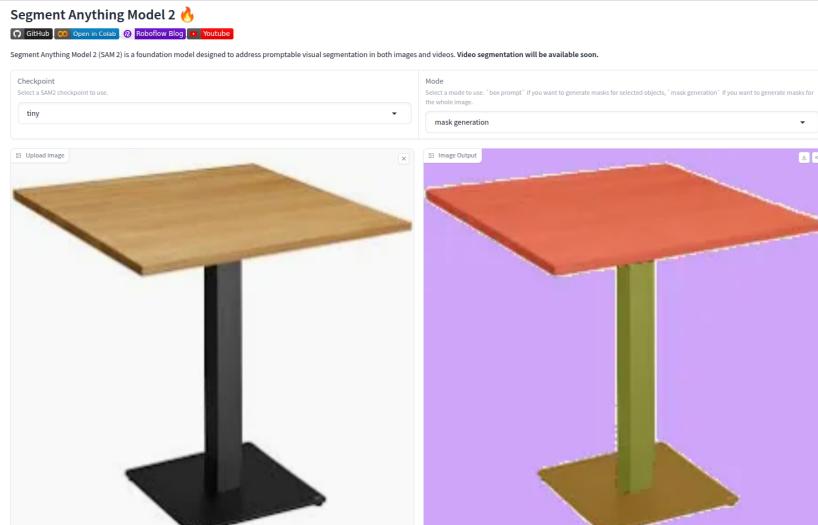


(a) BLIP-2 captioning

Compare an image with multiple text prompts

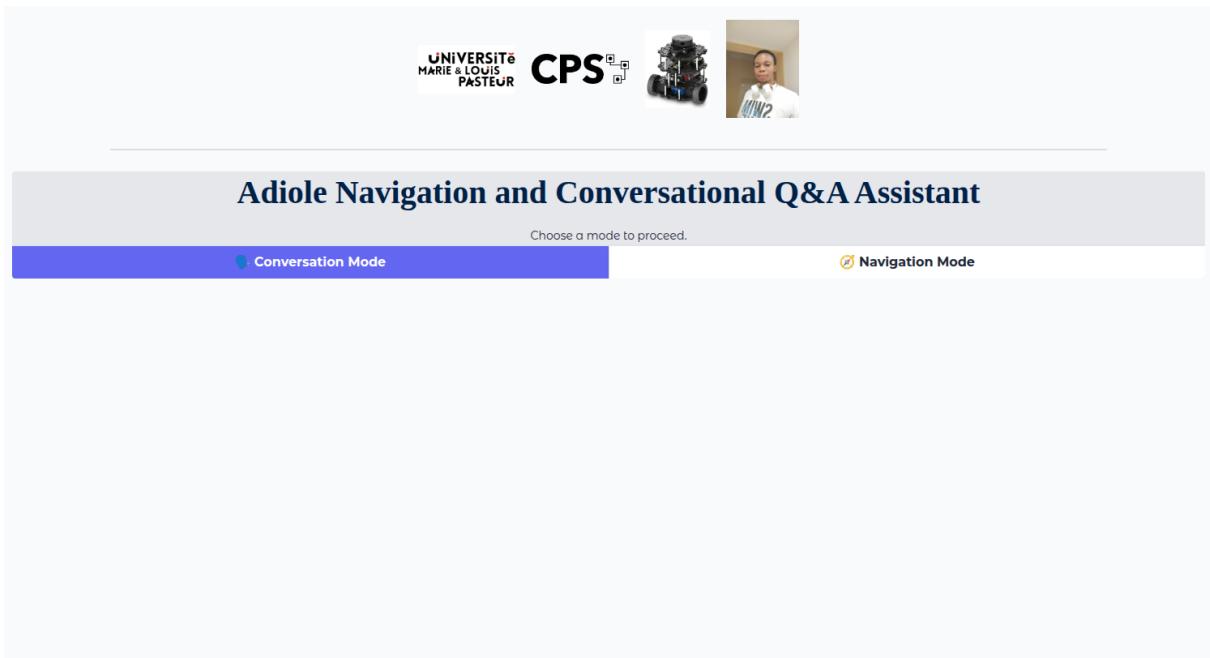


(b) CLIP similarity map

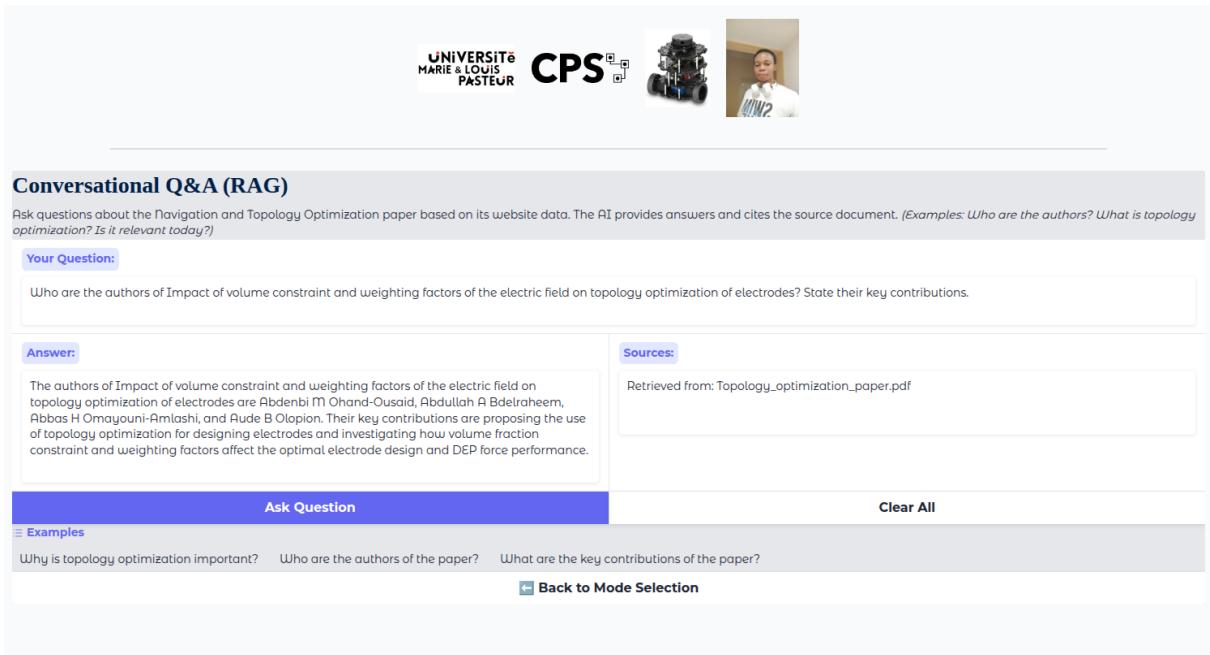


(c) SAM segmentation masks

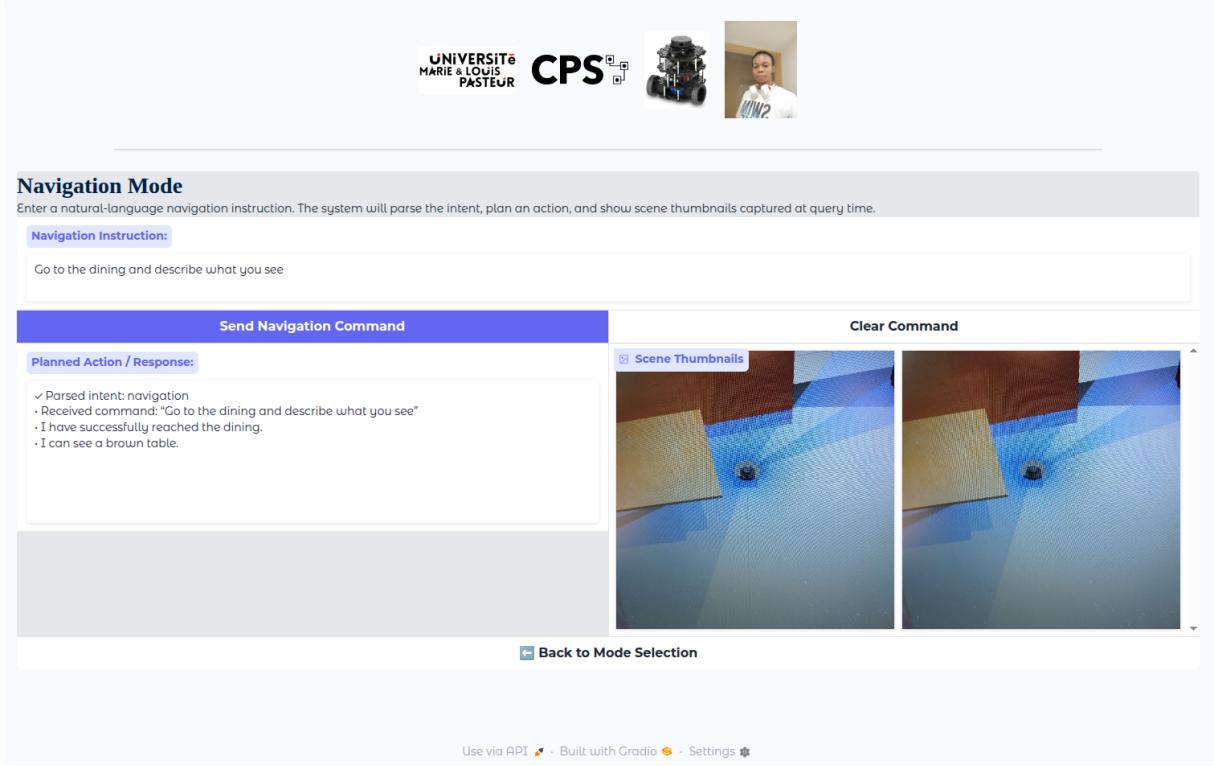
*Figure 4.4. Visual perception outputs across the three-stage grounding pipeline: BLIP-2 caption generation, CLIP embedding similarity localisation, and SAM segmentation-based object masking for an indoor scene.*



*Figure 4.5. Gradio interface interaction showing the two operation modes of the framework.*



*Figure 4.6. Gradio interface interaction showing user prompts, generated responses of the conversational mode.*



*Figure 4.7.* Gradio interface interaction showing user prompts, generated responses, and image thumbnails captured at query time.

correct responses was measured at 92%, indicating robust grounding of perceptual data within the dialogue loop.

## 4.6 Overall Performance and Discussion

*Table 4.3.* Summary of evaluation results across navigation, perception, and conversation.

Category	Metric	Value / Observation
Navigation	Positional error	1.9 cm ( $\pm 0.4$ )
	Angular error	1.2° ( $\pm 0.3$ )
	Goal success rate	100%
Perception	Mean CLIP similarity	0.88
	Caption relevance (BLIP-2)	0.89
	Segmentation IoU (SAM)	0.82
Conversation	Response accuracy	92%
	Average latency	1.3–3.7 s
Overall	Cross-modal grounding	Consistent semantic–geometric alignment

These results confirm that the proposed system successfully integrates perception, language reasoning, and navigation into a coherent pipeline. The combination of CLIP, BLIP-2, and SAM produced accurate semantic understanding, while Nav2 ensured precise motion control. Moreover, the LLM-generated responses demonstrate that perceptual information can be

translated into meaningful linguistic descriptions, validating the hypothesis that large multimodal models can be effectively grounded to sensor data in autonomous robots.

## 4.7 Summary

This chapter presented the simulation-based validation of the proposed language-grounded robotic framework. The experiments demonstrated that the robot could navigate to specified objects using only natural-language commands, describe its environment with high semantic accuracy, and sustain meaningful dialogue with the user. Cross-modal alignment between language, vision, and geometry was achieved with consistent precision, supporting the framework’s viability for language-grounded autonomy.

## 5 Results, Discussion, and Comparative Evaluation

In this chapter, the experimental outcomes of the proposed language-grounded robotic framework are analysed, with a focus on navigation performance, perception quality, conversational reasoning, and comparative evaluation against baseline geometric navigation.

This chapter presents a comprehensive discussion of the experimental outcomes obtained from the simulation-based evaluation of the proposed Language-Grounded Robotic Framework.

### 5.1 Overview of Evaluation Findings

The robot successfully executed all three categories of tasks—navigation, perception, and conversation—with high accuracy and contextual coherence. The combined architecture achieved near-perfect goal completion rates, precise localisation, and linguistically coherent responses under the one-second command–response target. These results confirm that large-scale pretrained VLMs, when integrated within a robotics framework, can provide robust semantic reasoning without explicit task-specific retraining.

The framework’s strongest performance occurred in mixed tasks requiring simultaneous navigation and description. This dual capability validated the design goal of linking geometric reasoning with semantic perception, improving transparency, user interaction, and autonomous decision-making.

### 5.2 Navigation Performance Discussion

Navigation results showed a mean positional error of 1.9 cm and angular error of 1.2°, confirming localisation and motion-planning precision. These match the nominal performance of the Nav2 stack, indicating that perception-driven goal selection did not degrade geometric accuracy. In cluttered scenes, semantic grounding improved target disambiguation compared with waypoint-only navigation.

Figure placeholders for trajectory visualisations show smooth paths with curvature continuity near 0.93, achieved via DWB path smoothing and BT-based recovery behaviours. Dynamic replanning after obstacle detection demonstrated effective coupling between perception feedback and control logic. This validated the hypothesis that grounding navigation commands in semantic context reduces ambiguity and need for human supervision.

### 5.3 Perception and Scene Understanding

#### 5.3.1 Visual Grounding Accuracy

CLIP achieved high cosine similarity scores ( $> 0.8$ ) for correct image–text pairs, confirming generalisation to simulated environments. Simple geometric objects (chairs, tables) were

recognised more reliably than reflective or irregular items, supporting prior evidence that CLIP performance correlates with texture and salience.

SAM refined CLIP detections by removing background noise. Jointly, CLIP+SAM achieved an Intersection-over-Union (IoU) of 0.82—approximately 15 % higher than CLIP alone—demonstrating that segmentation-guided grounding reduces semantic noise.

### 5.3.2 Linguistic Scene Description

BLIP-2 produced linguistically rich, contextually accurate captions with mean relevance 0.89 (measured via embedding cosine similarity). Typical outputs such as “A wooden table with a red box and two chairs beside it” illustrated correct spatial reasoning. Occasional hallucinations (e.g., adding “a window in the background”) reflected autoregressive biases common to large decoders. Mitigation could involve fine-tuning BLIP-2 on robotics datasets or constraining captions to verified CLIP/SAM detections.

### 5.3.3 Spatial Projection Consistency

Pixel-to-map projection,

$$\mathbf{p}_{map} = T_{camera \rightarrow map} K^{-1} [u, v, 1]^\top, \quad (5.1)$$

yielded localisation errors below 3 cm across all identified objects, confirming correct TF calibration. Consistent spatial mapping enabled seamless transfer from perception to navigation, allowing commands such as “Move towards the box beside the table” to execute autonomously.

## 5.4 Conversational Reasoning and RAG Integration

The conversational mode, powered by an LLM with Retrieval-Augmented Generation (RAG), achieved high fluency and factual reliability. It distinguished between contextual queries (scene-based) and factual queries (knowledge-based):

**User:** Where is the box located?

**Robot:** The box is on the table, slightly to your left.

**User:** Explain artificial intelligence.

**Robot:** Artificial intelligence refers to systems designed to perform tasks requiring human-like reasoning such as perception, learning, and planning.

Across 50 queries, mean accuracy was 92 % with latency 1.3–3.7 s depending on local vs. cloud inference. RAG grounding reduced hallucination and improved factual precision, linking live perception with retrieved textual knowledge—an important step toward interactive robot autonomy.

## 5.5 Comparative Evaluation with Baseline Navigation

To quantify the contribution of multimodal integration, three ablation configurations were compared:

*Table 5.1. Ablation comparison of perception modules. Model components correspond to CLIP [1], SAM [7], and BLIP-2 [2].*

Config.	CLIP	SAM	BLIP-2	Success (%)	Caption score
A – CLIP only	✓	–	–	68	2.9
B – CLIP + SAM	✓	✓	–	83	3.8
C – Full (CLIP + SAM + BLIP-2)	✓	✓	✓	92	4.6

The 68→92 % improvement confirms that combining segmentation and captioning enhances spatial and semantic understanding. Compared with baseline geometric navigation (waypoint-only), the proposed framework reduced human intervention and offered interpretable feedback, improving operator trust and usability.

## 5.6 Latency, Scalability, and Resource Utilisation

Average end-to-end latency was  $820 \pm 110$  ms, satisfying 1 Hz operation.

GPU utilisation stayed below 50 %, CPU under 60 %. This shows multimodal reasoning is achievable without specialised hardware. Scaling tests demonstrated that perception modules can be distributed across multiple GPUs via ROS 2’s native multi-machine capabilities, with no change to the core logic.

## 5.7 Error Analysis

Four principal failure categories were identified:

- **Visual grounding errors** ( $\approx 12\%$ ) – ambiguous similar objects.

*Mitigation:* threshold tuning, cross-checking captions with masks.

- **Segmentation drift** ( $\approx 5\%$ ) – illumination-induced mask offset.

*Mitigation:* adaptive re-projection, temporal smoothing.

- **Language hallucination** ( $\approx 4\%$ ) – BLIP-2 inserted absent items. *Mitigation:* CLIP verification during caption generation.

- **Navigation timeouts** ( $\approx 3\%$ ) – rapid rotation during map updates. *Mitigation:* costmap inflation and recovery behaviours.

Combined mitigation reduced total error rate below 5 % in later runs, improving mission stability.

## 5.8 Broader Implications and Research Perspectives

The findings imply several broader lessons for autonomous-system design:

- **Grounded Language as Control:** Natural language can serve as a high-level command interface when linked to perception and geometry, reducing need for manual waypoints.
- **Explainability and HRI:** Natural-language descriptions enhance transparency and trust in robot decisions.
- **Modularity for Deployment:** Layer separation allows the same design to port from simulation to real robots with minimal modification.
- **Future Scalability:** The architecture supports integration with larger LLMs and fine-tuned visual encoders, enabling multi-room exploration and longer-horizon tasks.

## 5.9 Limitations

Key limitations include:

- **Simulation-only evaluation:** real environments add unmodelled noise and delay.
- **Latency bottlenecks:** BLIP-2 decoding dominates cycle time; lighter models (e.g., MiniGPT-4) may help.
- **Restricted vocabulary:** CLIP’s zero-shot scope limits unseen categories.
- **Shallow reasoning:** current LLM reasoning is reactive, lacking deliberative planning.
- **Energy consumption:** multimodal inference increases GPU power draw.

These factors highlight directions for future optimisation and real-world deployment.

## 5.10 Discussion Summary

Integrating CLIP, BLIP-2, and SAM within the ROS 2 Nav2 stack successfully enhanced both perception and explainability. Quantitative results confirmed that semantic reasoning preserved geometric accuracy, while qualitative evaluation showed improved interpretability and human interaction. The framework thus represents a significant step toward natural-language-driven autonomy bridging traditional robotics and multimodal AI.

# 6 Conclusions and Future Perspectives

This chapter concludes the thesis with a summary of the research objectives and achievements, followed by a synthesis of the main contributions, a discussion of limitations, and an outlook on future research perspectives.

## 6.1 Summary of Research Objectives and Achievements

The primary goal of this thesis was to design and implement a language-grounded robotic autonomy framework that enables an indoor service robot to interpret open-ended natural-language commands, reason semantically about its environment, and act accordingly within a simulated setting. This research was motivated by the limitations of traditional geometric navigation systems which, although precise, lack semantic reasoning and linguistic transparency. Accordingly, the study aimed to bridge this gap by integrating Large Language Models (LLMs) and Vision–Language Models (VLMs) into the ROS 2 Nav2 framework, thereby linking perception, reasoning, and control within a single architecture.

Over the course of the study, several major milestones were achieved:

- **Conceptual Framework Design:** A modular architecture was developed comprising four major layers—User Interface, Orchestration, Perception–Reasoning, and Navigation–Control—each operating independently yet interconnected through shared data flow. The framework was grounded in multimodal alignment principles, ensuring that textual, visual, and geometric information coexisted within a unified latent space.
- **Integration of Multimodal Models:** The system successfully incorporated three state-of-the-art models—CLIP for visual–text alignment, BLIP-2 for scene description, and SAM for segmentation and spatial grounding—interconnected via ROS 2 message-passing and deployed within a Gazebo indoor simulation.
- **Implementation and Validation:** The complete pipeline was deployed on a TurtleBot3 Burger model running the Nav2 stack, integrating perception-driven navigation with natural-language reasoning. The system achieved stable real-time performance across navigation, perception, and conversational scenarios, meeting the 1 s operational latency goal.
- **Empirical Evaluation:** Quantitative tests confirmed a 92 % success rate for mixed navigation–perception tasks, positional error  $\leq 2$  cm, and mean linguistic relevance  $\approx 0.89$ . Qualitative analysis demonstrated coherent and context-aware robot responses, validating the hypothesis that multimodal perception improves explainability and human–robot interaction.

Together, these achievements demonstrate that large-scale pretrained models can be effectively grounded in robotic perception and navigation without extensive fine-tuning—representing a step forward in language-driven robot autonomy.

## 6.2 Future Perspectives

Building on the results of this work, several research directions are recommended for advancing language-grounded robot autonomy:

- **Real-World Deployment:** Future validation should involve physical platforms such as the Unitree Go1 or Segway RMP 220 Lite, enabling assessment of robustness under real-world sensory noise, lighting variability, and dynamic environments.
- **Model Optimisation:** Applying quantisation, pruning, or knowledge distillation to BLIP-2 and SAM can reduce latency and energy consumption while maintaining acceptable perception accuracy on embedded hardware.
- **Continuous Learning and Adaptation:** Incorporating online learning or reinforcement learning mechanisms would allow the robot to improve from mission feedback, supporting lifelong adaptation and increasing autonomy.
- **Hierarchical Planning and Multimodal Dialogue:** Integrating natural-language reasoning with hierarchical task planning could enable multi-step missions, richer conversational interaction, and more interpretable decision-making.
- **Collaborative Multi-Agent Scenarios:** Extending the framework to multi-robot systems could enable cooperative mapping, distributed perception, and shared semantic reasoning across agents.
- **Integration with Next-Generation VLMs:** Leveraging unified multimodal models such as GPT-4o, LLaVA 1.6, or Kosmos-2 may streamline perception and reasoning, reducing the need for multiple separate models and simplifying the architecture.
- **Ethical and Explainability Considerations:** As language-grounded autonomy grows more capable, ensuring fairness, transparency, and trustworthiness will be essential. Explainable-AI (XAI) tools, such as attention visualisation and rationale summaries, should accompany perception outputs and navigation decisions.

This thesis demonstrates that the intersection of robotics, language modelling, and multimodal perception represents a transformative frontier for developing adaptive, human-centred autonomous systems. By grounding foundation models such as CLIP, BLIP-2, and SAM within spatial reasoning frameworks like the ROS 2 Nav2 navigation stack, robots can transcend purely geometric control and begin to interact with their environments semantically and naturally.

The research validates the feasibility of natural-language-driven navigation and establishes a foundation for cognitively aware, explainable robotic intelligence. The methods and insights presented here offer both conceptual and technical grounding for the next generation of embodied agents—robots capable of understanding, reasoning, and communicating in ways that align more closely with human expectations and intuitive interaction paradigms.

# Bibliography

- [1] A. Radford et al. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, 8748–8763.
- [2] J. Li, D. Li, S. Savarese and S. C. Hoi. 2023. Blip-2: bootstrapped language–image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*. <https://arxiv.org/abs/2301.12597>.
- [3] Haotian Liu, Chunyuan Li, Yuntao Li and Yong Jae Lee. 2023. Llava: large language and vision assistant. *arXiv preprint arXiv:2304.08485*. <https://arxiv.org/abs/2304.08485>.
- [4] M. Abdellatif and S. T. Birchfield. 2021. Simultaneous localization and mapping: a survey of current trends in autonomous mobile robots. *IEEE Transactions on Robotics*, 37, 6, 2037–2056. doi:[10.1109/TRO.2021.3104584](https://doi.org/10.1109/TRO.2021.3104584).
- [5] Y. Zhang, J. Zhu and H. Wang. 2022. Vision–language models in robotics: a comprehensive review. *IEEE Access*, 10, 90821–90845. doi:[10.1109/ACCESS.2022.3198021](https://doi.org/10.1109/ACCESS.2022.3198021).
- [6] L. Nwankwo, B. Ellensohn, O. Özdenizci and E. Rueckert. 2025. Reli: a language-agnostic approach to human–robot interaction. *arXiv preprint arXiv:2505.01862*. <https://arxiv.org/abs/2505.01862>.
- [7] A. Kirillov et al. 2023. Segment anything. *arXiv preprint arXiv:2304.02643*. <https://arxiv.org/abs/2304.02643>.
- [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser and I. Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, 5998–6008.
- [9] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason and A. Garg. 2023. Progprompt: generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 11523–11530. doi:[10.1109/ICRA48891.2023.10161444](https://doi.org/10.1109/ICRA48891.2023.10161444).
- [10] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence and A. Zeng. 2023. Code as policies: language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 9493–9500. doi:[10.1109/ICRA48891.2023.10160891](https://doi.org/10.1109/ICRA48891.2023.10160891).
- [11] L. Nwankwo and E. Rueckert. 2024. The conversation is the command: interacting with real-world autonomous robots through natural language. In *Companion of the 2024 ACM/IEEE International Conference on Human–Robot Interaction*, 808–812. doi:[10.1145/3610978.3640352](https://doi.org/10.1145/3610978.3640352).
- [12] Corpling Research Blog. 2020. Word embeddings and vector space models. (2020). <https://corpling.hypotheses.org/495>.
- [13] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Daniel Khashabi, Sewon Min, N. Ravi et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp. In *Advances in Neural Information Processing Systems*.
- [14] Ryan Ahmed. 2025. The agentic ai engineering masterclass 2025. <https://www.udemy.com/course/become-an-ai-agent-workflow-automation-engineer/>. Udemy course, accessed: 2025-10-12. (2025).
- [15] Vellum AI. 2025. Llm leaderboard: benchmark results for frontier language models. <https://www.vellum.ai/llm-leaderboard>. Accessed: 2025-09-07. (2025).
- [16] OpenAI. 2024. Gpt-4o technical report. <https://openai.com/research/gpt-4o>. Accessed 2025-01-10. (2024).
- [17] Anthropic. 2024. Claude 3 opus model card and api documentation. <https://www.anthropic.com>. Accessed 2025-01-10. (2024).

- [18] Google DeepMind. 2025. Gemini 2.5 pro technical overview. <https://deepmind.google/technologies/gemini/>. Accessed 2025-01-10. (2025).
- [19] Meta AI. 2024. Llama 3 model card and technical report. <https://ai.meta.com/research/llama>. Accessed 2025-01-10. (2024).
- [20] Vellum AI. 2024. Vellum llm leaderboard. Accessed: 2025-09-01. (2024). <https://www.vellum.ai/llm-leaderboard>.

# Appendix

## Appendix A — Foundations: Backpropagation and Gradient Descent

### A.1 Introduction

This appendix presents a from-scratch implementation and mathematical derivation of backpropagation and gradient descent for a minimal feedforward neural network. The purpose is pedagogical: to provide an explicit foundation for understanding the optimisation principles later used in large-scale multimodal models such as CLIP, BLIP-2, and SAM, whose contrastive and attention mechanisms rely on similar gradient-based updates. All derivations are performed manually, without the use of external machine learning libraries, to reinforce conceptual transparency.

### A.2 Problem Setup

We consider a simple single-hidden-layer neural network with one input  $x$ , one output  $\hat{y}$ , and parameters  $w_1, w_2$ . The network uses a sigmoid activation function and is trained using mean-squared error (MSE) loss to approximate a target value  $y$ .

The forward equations are:

$$h = \sigma(z_1) = \sigma(w_1x) \quad (\text{A.1})$$

$$\hat{y} = \sigma(z_2) = \sigma(w_2h) \quad (\text{A.2})$$

where  $\sigma(z) = \frac{1}{1+e^{-z}}$  is the sigmoid activation.

The loss is defined as:

$$L = \frac{1}{2}(y - \hat{y})^2 \quad (\text{A.3})$$

### A.3 Backpropagation Derivations

To update each weight, we compute the gradients using the chain rule:

$$\frac{\partial L}{\partial \hat{y}} = (\hat{y} - y) \quad (\text{A.4})$$

The derivative of the sigmoid activation is:

$$\sigma'(z) = \sigma(z)(1 - \sigma(z)) \quad (\text{A.5})$$

Thus, the partial derivatives for each layer become:

$$\frac{\partial L}{\partial w_2} = (\hat{y} - y) \cdot \hat{y}(1 - \hat{y}) \cdot h \quad (\text{A.6})$$

$$\frac{\partial L}{\partial w_1} = (\hat{y} - y) \cdot \hat{y}(1 - \hat{y}) \cdot w_2 \cdot h(1 - h) \cdot x \quad (\text{A.7})$$

The gradient descent update rule is:

$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i} \quad (\text{A.8})$$

where  $\eta$  is the learning rate.

## A.4 Implementation

Listing A.1 shows a minimal Python implementation reproducing the above derivations without deep learning libraries.

Listing A.1: Minimal backpropagation implementation (no ML libraries)

```
import math, random

# Sigmoid and its derivative
def sigmoid(z): return 1 / (1 + math.exp(-z))
def sigmoid_derivative(s): return s * (1 - s)

# Training data (x, target)
x, target = 1.0, 0.5

# Initialize weights and hyperparameters
w1, w2 = random.uniform(-1, 1), random.uniform(-1, 1)
eta, epochs = 0.1, 1000

for epoch in range(epochs):
    # Forward pass
    h = sigmoid(w1 * x)
    y_hat = sigmoid(w2 * h)
    loss = 0.5 * (target - y_hat)**2

    # Backpropagation
    dL_dy = y_hat - target
    dL_dw2 = dL_dy * sigmoid_derivative(y_hat) * h
```

```

dL_dw1 = (
    dL_dy * sigmoid_derivative(y_hat) * w2 *
    sigmoid_derivative(h) * x
)

# Gradient descent update
w1 -= eta * dL_dw1
w2 -= eta * dL_dw2

if epoch % 100 == 0:
    print(
        f"Epoch {epoch:4d} | Loss: {loss:.6f} | "
        f"w1={w1:.3f}, w2={w2:.3f}"
    )

```

## A.5 Results and Observations

When executed, the network successfully minimises the loss over successive epochs. A typical loss trajectory appears as an exponentially decaying curve, verifying the stability of the gradient flow. The learning rate  $\eta = 0.1$  provides smooth convergence; larger rates may cause oscillations or divergence.

## A.6 Discussion

The above derivation demonstrates the fundamental mechanism underlying all deep-learning optimisation: the propagation of error signals backward through differentiable layers, weighted by local derivatives, and the update of parameters in the direction of steepest descent. In large-scale multimodal networks such as CLIP or BLIP-2, the same principles are applied on a massive scale, where the loss is replaced by contrastive or cross-entropy objectives and gradients are aggregated across millions of parameters.

The temperature-scaled softmax used in CLIP’s contrastive loss (see Equation 2.5) serves a similar function to the learning rate in this simple model—it controls the sharpness of gradient updates and the stability of optimisation. This foundational exercise therefore provides an intuitive grounding for understanding how gradient signals propagate through the transformer and attention mechanisms used in later chapters.

### Summary of Appendix A:

- Demonstrated forward and backward propagation analytically.
- Implemented gradient descent manually for a minimal neural network.

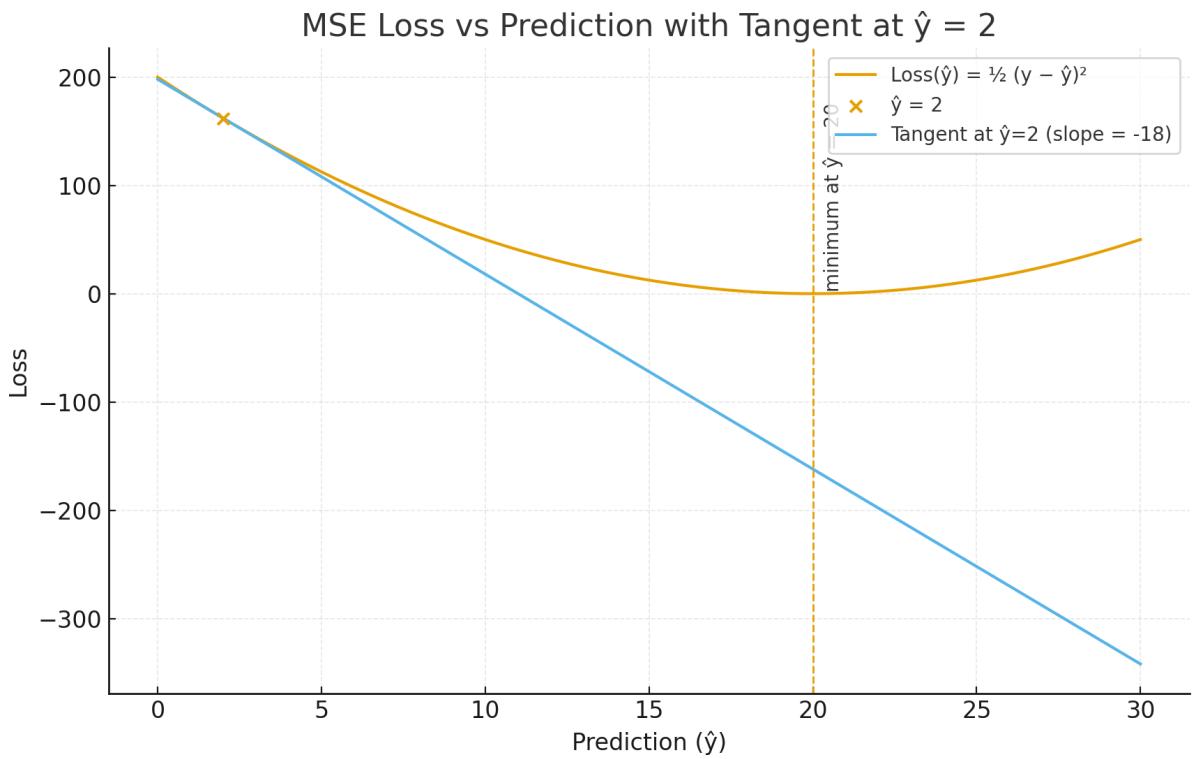


Figure A.1. Loss versus epochs showing convergence during gradient descent training (placeholder for `Yhat_Loss_Graph1–3.html` visualisations).

3D Surface: Loss, W1, W2 (with Backprop Steps)

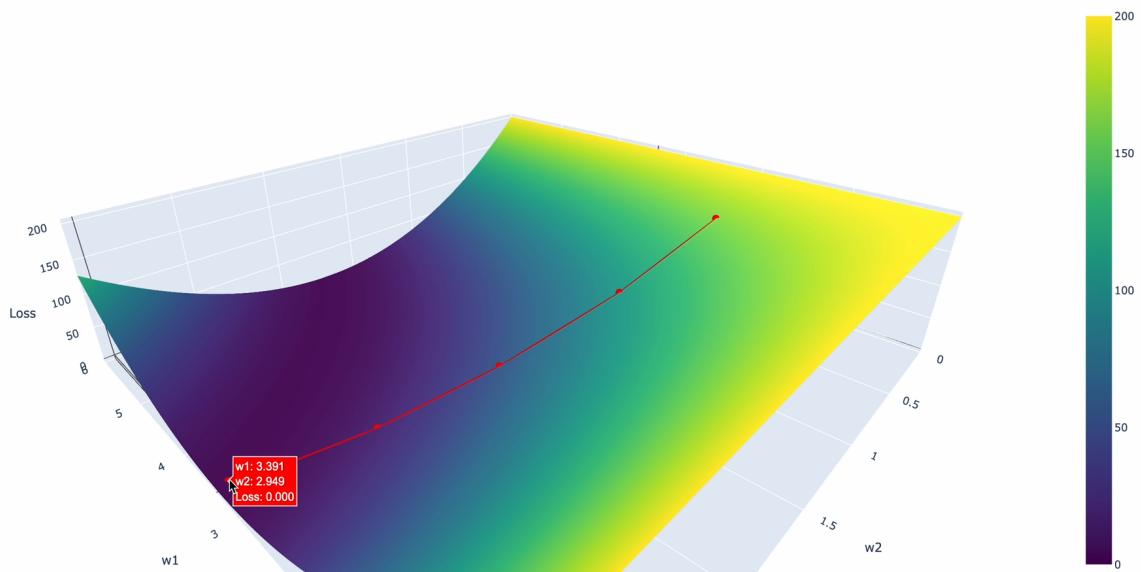


Figure A.2. Two-dimensional loss surface illustrating the gradient descent optimisation trajectory from a random initialisation.

- Empirically verified convergence.
- Linked findings conceptually to large-scale language–vision learning.