# Enabling Novel Mission Operations and Interactions with ROSA: The Robot Operating System Agent

Rob Royce, Marcel Kaufmann, Jonathan Becktor, Sangwoo Moon, Kalind Carpenter,
Kai Pak, Amanda Towler, Rohan Thakker, Shehryar Khattak

NASA Jet Propulsion Laboratory
California Institute of Technology
4800 Oak Grove Dr
Pasadena, CA 91109

{rob.royce; marcel.kaufmann; jonathan.becktor; sangwoo.moon; kalind.carpenter;
kai.pak; amanda.towler; rohan.a.thakker; skhattak}@jpl.nasa.gov

*Abstract*—The advancement of robotic systems has revolutionized numerous industries, yet their operation often demands specialized technical knowledge, limiting accessibility for non-expert users. This paper introduces ROSA (Robot Operating System Agent), an AI-powered agent that bridges the gap between the Robot Operating System (ROS) and natural language interfaces. By leveraging state-of-the-art language models and integrating open-source frameworks, ROSA enables operators to interact with robots using natural language, translating commands into actions and interfacing with ROS through well-defined tools. ROSA's design is modular and extensible, offering seamless integration with both ROS1 and ROS2, along with safety mechanisms like parameter validation and constraint enforcement to ensure secure, reliable operations. While ROSA is originally designed for ROS, it can be extended to work with other robotics middle-wares to maximize compatibility across missions. ROSA enhances human-robot interaction by democratizing access to complex robotic systems, empowering users of all expertise levels with multi-modal capabilities such as speech integration and visual perception. Ethical considerations are thoroughly addressed, guided by foundational principles like Asimov's Three Laws of Robotics, ensuring that AI integration promotes safety, transparency, privacy, and accountability. By making robotic technology more user-friendly and accessible, ROSA not only improves operational efficiency but also sets a new standard for responsible AI use in robotics and potentially future mission operations. This paper introduces ROSA's architecture and showcases initial mock-up operations in JPL's Mars Yard, a laboratory, and a simulation using three different robots. The core ROSA library is available on GitHub [1].

## TABLE OF CONTENTS

[1] https://github.com/nasa-jpl/rosa

## 1. INTRODUCTION

In recent years, the field of robotics has witnessed significant advancements, with robots playing crucial roles in industries such as manufacturing, healthcare, and space exploration. As robotic systems become more complex and capable, the need for intuitive and efficient human-robot interaction has grown exponentially [1]. However, traditional methods of controlling and programming robots often require specialized knowledge and technical expertise, creating barriers for operators and limiting the accessibility of robotic technologies. It is not uncommon for operators to require advanced degrees and meticulous training to effectively operate even modestly capable robotic systems [2].

ROSA, the *Robot Operating System Agent* addresses these challenges by introducing a first of its kind AI Agent that enables human-robot interaction (HRI) in the form of natural language. ROSA leverages large language models (LLMs) and open-source software like the Robot Operating System (ROS) [3] to implement a Reasoning and Acting (ReAct) agent [4] capable of understanding and executing commands on its robotic host. By integrating with the ROS and ROS2 [5] ecosystems, ROSA provides easy access to a wide range of tools and functionalities that allow users to perform tasks such as system diagnostics, monitoring, and invoking existing navigation or manipulation tasks, without the need for extensive technical training. ROSA aims to make robotic systems more accessible and user-friendly, empowering a broader range of operators to interact with robots in a more intuitive manner.

The remainder of this paper is organized as follows: Section 2 situates ROSA within the context of related work, providing background on ROS, LLMs, and ReAct agents. Section 3 details ROSA's agent architecture, including its action space, memory modules, internal logic, and LLM-based tool-calling. Section 4 presents three real-world demonstrations of ROSA operating on different robotic systems across varied scenarios. Section 5 explores how ROSA enhances human-robot interaction by bridging knowledge gaps and enabling multi-modal and multi-lingual communication. Section 6 elaborates on the open-source implementation of ROSA, covering tool design, parameterization, ROS functionalities, and model selection. Ethical considerations are addressed in Section 7, with a discussion of foundational principles and community guidelines for responsible AI and robotics. Finally, Section 8 summarizes the contributions and outlines directions for future work.

(a) NeBula-Spot in the Mars Yard      (b) EELS in Laboratory Environment      (c) Nova Carter in Simulation

**Figure 1**: **ROSA has been deployed on various robots at NASA Jet Propulsion Laboratory, including NeBula-Spot in the Mars Yard (a), Exobiology Extant Life Surveyor (EELS) in a laboratory environment using synthetic terrain and obstacles (b), and Nova Carter in NVIDIA IsaacSim using a simulated Martian environment (c). Each of these integrations and a demonstration of their capabilities when equipped with ROSA are detailed in Section 4.**

## 2. BACKGROUND

The past decade has seen remarkable advancements in robotics, resulting in highly sophisticated systems with enhanced capabilities in perception, manipulation, and autonomy. Despite these advancements, interacting with robotic systems remains a significant challenge, particularly for users lacking extensive technical expertise. Frameworks such as ROS have streamlined development by encouraging modular and reusable robotic software components. However, the steep learning curve associated with controlling and interacting with robots persists, especially when relying on traditional command-line interfaces or programming languages. Recent breakthroughs in natural language processing (NLP), particularly through LLMs, present a compelling opportunity to bridge the gap between humans and AI-driven agents [6]. This section explores related work in robotics development and AI-powered agents while providing an overview of the foundational technologies that underpin ROSA.

**Related Work**

There have been several notable efforts aimed at using NLP for robotic control. Projects like Microsoft's LATTE [7] and research efforts like Google's "Language Models as Zero-Shot Planners" [8] explore how LLMs can be utilized for robotic task planning and execution. These initiatives showcase the potential of LLMs to interpret human instructions and translate them into robotic actions. However, direct integration of LLMs with robotic systems remains a challenge, primarily due to the inherent ambiguity of natural language and the difficulties in ensuring the consistent and reliable execution of tasks.

ROSA distinguishes itself from these efforts by adopting a more pragmatic and targeted approach. Instead of using LLMs for direct robot control, ROSA functions as an operator interface that leverages existing robot functionality within the ROS ecosystem. ROSA capitalizes on the strengths of LLMs to translate natural language into high-level commands while managing the intricacies of ROS behind the scenes. This method allows ROSA to provide a more reliable and user-friendly interface, avoiding some of the challenges inherent in direct LLM-based robotic control, while still benefiting from LLM's emergent capabilities.

Koubaa et al. investigate human-robot interaction (HRI) by utilizing the capability of LLMs to convert natural lan-guage queries into formal commands, enabling written or spoken language to be translated into structured commands for robotic systems [1]. While this method breaks down certain barriers in HRI, it also introduces new challenges. Specifically, it demands expertise in "prompt engineering" and "ontology-guided task elicitation", which require a deep technical understanding of LLM functionality, ontology creation, and formal specification. ROSA, by contrast, aims to simplify this process by utilizing pre-built frameworks and established methodologies to map natural language to robotic behavior. This allows developers to focus more on enhancing robotic capabilities rather than mastering the complexities of LLM operation and ontology design. Another important distinction between ROSA and [1] is that ROSA constitutes an **embodied agent**, as detailed in the following sections, which provides additional flexibility in the way queries are transformed into robotic actions.

**Preliminaries**

*Robot Operating System*—ROS is an open-source framework designed for developing robot software [3]. It provides a collection of tools, libraries, and conventions that streamline the creation of complex robotic behaviors across various platforms. By standardizing robotics programming, ROS encourages code reuse and collaboration, making it a cornerstone of modern robotics.

At its core, ROS enables modularity and scalability through a peer-to-peer network of processes called *nodes*, which can be distributed across multiple computers. These nodes communicate via a publisher/subscriber model using *topics*, or request/response interactions through *services*. Figure 2 shows a basic ROS setup, but real-world systems are often much more complex, involving dozens or even hundreds of nodes, topics, and services. Each ROS distribution is packaged with a set of tools that allow for monitoring, configuring, and interacting with the underlying system. As described in later sections, these tools are used directly by ROSA to carry out operator instructions, making it significantly easier to accomplish tasks that would otherwise require significant training and expertise.

*LLMs and Natural Language Interfaces*—The emergence of LLMs like GPT-3 [9] has revolutionized NLP and enabled more intuitive human-computer interactions. In the context of robotics, LLMs have been proven capable of interpreting natural language commands and translating them into robotic
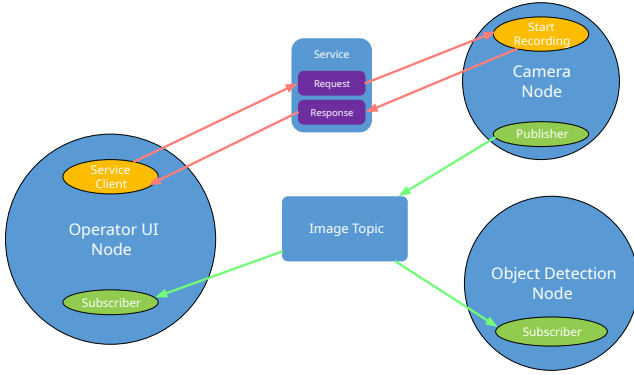
**Figure 2**: **A basic ROS system with three nodes (Operator UI, Camera, and Object Detection), one topic (Image), one service client, and one service provider (Start Recording). Image adapted from ROS Tutorials.**

actions [1], [7], [9]. For ROSA, the integration of LLMs is pivotal in bridging the gap between users and the robot's underlying capabilities. LLMs allow ROSA to process everyday language queries, interpret their intent, and generate appropriate actions within the ROS environment. This natural language interface eliminates the need for users to interact directly with command-line tools or custom scripts, and therefore significantly lowers the learning curve and empowers a wide array of people to interact with robotic systems. By handling linguistic nuances, context, and even ambiguous or incomplete instructions, LLMs make ROSA more accessible to non-expert users while maintaining robust control over the robot's functions.

*ReAct Agents*—To leverage the power of LLMs in robotics applications, ROSA utilizes the ReAct (Reasoning and Acting) paradigm [4]. ReAct agents operate through a reasoning-action-observation loop. The cycle begins with **Reasoning**, where the LLM interprets natural language input and formulates a plan. This plan may involve a single action, a sequence of actions, or follow-up questions to effectively address the query. Next is the **Action** phase, during which the agent executes the planned steps by invoking the necessary tools. In ROSA, these tools are represented as dynamically callable Python functions that the LLM can access as needed. Finally, in the **Observation** phase, the agent assesses the results of its actions to guide subsequent reasoning steps and contribute to constructing a final response. If the actions and observations do not resolve the query, the loop repeats until a resolution is reached or a predefined iteration limit is met.

Figure 3 presents an example reasoning trace in the context of ROSA. This trace illustrates the ReAct loop undertaken by an agent when responding to the query, "Provide a list of ROS nodes." During the reasoning phase, the agent interprets the user's request and determines that invoking a specific tool is necessary to address the query. The tool is then executed within the agent's operational environment. Following the action, the agent observes the tool's effects on the environment. In this instance, a single tool invocation is sufficient to fulfill the request, enabling the agent to provide a formatted list of ROS nodes as the response. In more complex scenarios, the reasoning trace may involve multiple iterations of the ReAct loop, with actions performed either sequentially or concurrently.

ROSA builds upon these foundational technologies by inte-



**Figure 3**: **An example reasoning trace that demonstrates the reasoning-action-observation loop characteristic of ReAct agents. This trace demonstrates a typical response from ROSA when given the query "Provide me with a list of ROS nodes."**

grating LLMs with ROS to facilitate a ReAct agent. The key innovation of ROSA lies in this integration, along with the suite of custom tools specifically designed for interacting with ROS and ROS2 systems. The tools developed for ROSA are used to invoke ROS utilities like *rosnode / ros2 node*, *rostopic / ros2 topic*, and *rviz*, but designed to be more interpretable by LLMs, allowing for more effective and reliable usage. By abstracting the complexities of ROS commands, these tools enable the LLM to generate appropriate actions in response to user inputs. The core tools provided by ROSA are outlined in section 6. ROSA also supports custom agents, which allow robotics developers to define their own tools in order to augment ROSA's action space. Details on how to create custom agents are omitted here, but are well documented in the open-source ROSA repository.

## 3. AGENT ARCHITECTURE

ROSA's overall architecture is composed of four primary components: the action space, a set of memory modules, internal logic for decision making, and the underlying LLM, as seen in Figure 4. Hence, ROSA follows what is known as the *Cognitive Architecture for Language Agents* [10]. Cognitive language agents come in a variety of forms, but are not typically considered "embodied agents" unless they operate within some physical, as opposed to digital, environment. While there are examples of embodied language agents in robotics [11], these agents typically focus on direct integration with planning, navigation, and locomotion capabilities. In contrast, ROSA focuses on human-robot interaction and utilizes existing robotic functionality rather than attempting
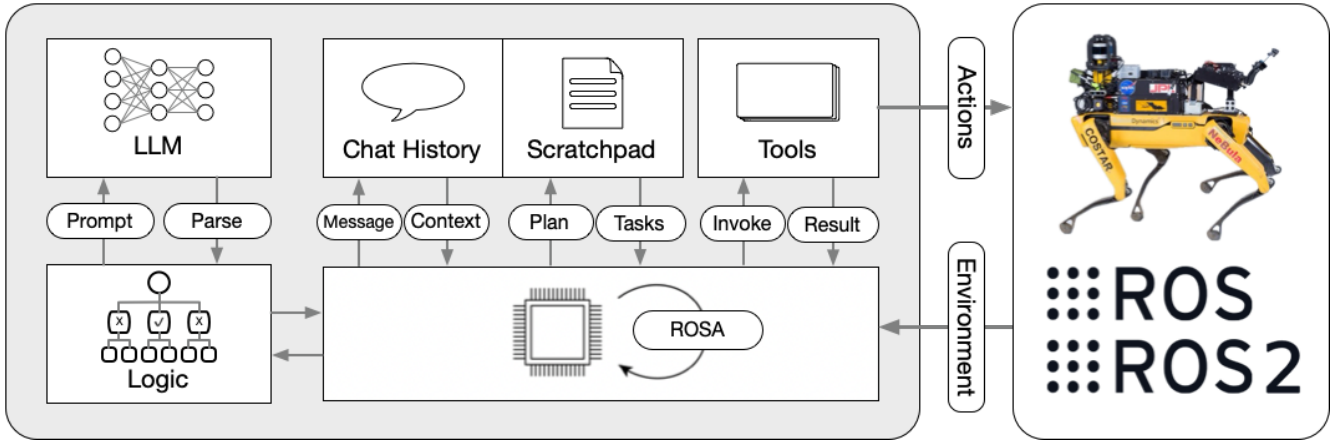
**Figure 4**: **The ROSA architecture includes memory (chat history and scratchpad), logic, tools (action space), and a tool-calling Large Language Model. Arrows in the diagram represent the general connectivity between components and the external robotics environment.**

to replace it with novel learning methods. The following subsections provide details on each of the major components in ROSA's architecture, including shortfalls and proposed solutions for future work.

**Action Space**

The action space in ROSA is defined as the set of tools[2] the agent can call within its working environment. ROSA's action space includes many of the standard ROS tools like `rosnode` and `ros2 topic`, as well as a set of utility tools for working with log files, capturing images and sensor readings, and so on. Developers can also extend ROSA's action space by adding custom functions tailored to specific robots or applications, thereby allowing ROSA to execute robot-specific commands and behaviors. ROSA's action space incorporates mechanisms to prevent the invocation of unsafe or unauthorized actions, ensuring that the robot operates within its intended parameters. All tools used by ROSA must follow a specific format, as defined in subsequent sections.

*Tool Invocation and Multi-Tool Usage*—Tool invocation is a core mechanism by which ROSA executes actions within the ROS environment. ROSA can invoke single tools at a time, or multiple tools either sequentially or in parallel, enabling it to handle complex tasks requiring coordinated or concurrent actions. Sequential execution ensures that dependent tasks—such as running diagnostics before calibrating sensors—occur in a logical order. Parallel execution, illustrated in the bottom half of Figure 5, enables ROSA to manage independent tasks simultaneously, such as retrieving telemetry data from multiple sensors at the same time. Moreover, ROSA can adapt its plan based on the outcomes of previous actions. If a tool fails or produces unexpected results, ROSA recalibrates its approach and proceeds accordingly. This multi-tool capability makes ROSA both efficient and flexible, ensuring it can respond dynamically to changing requirements.

Figure 5 provides an example sequence diagram illustrating both single- and multi-tool usage. (1) Suppose the user types the query: "Give me a status report". (2) In response, ROSA forwards the query, along with a list of all available

tools, prompts, and chat history, to the LLM. (3) In this case, the LLM decides to call the `get_robot_status` tool. (4) When ROSA receives the LLM's response, it invokes the associated `get_robot_status()` Python function and returns its output to the LLM. (5) At this point, the LLM calls two additional tools in parallel, specifying that ROSA should invoke `get_battery_status` and `get_cpu_status` concurrently to gather information about subsystem status. (6) Once again, ROSA returns the tools' outputs to the LLM, which can then formulate a final response for the user. This process typically completes in just a few seconds, allowing the user to quickly follow up with additional questions as needed.

**Memory**

While there are many forms of memory for LLM-based agents, ROSA operates on two of the most simple forms: chat history and agent scratchpad. The *chat history* is simply a list of messages that follow the open message format adopted by OpenAI, Anthropic, and other model providers. Each message has a `role` and `content` field, where the role can be either `user`, `assistant`, or `system`, and the content is the message itself. These messages are managed by the logic module to ensure they fit within the context window[3], as described in the next sub-section. The *scratchpad*, on the other hand, is where the agent keeps its most recent plans, especially for multi-step tasks. In practice, it is simply a reserved and annotated section of the context window where the agent is allowed to store text. While these two memory stores are sufficient for shorter interactions, we have found that they are limited in value for more complex scenarios, such as keeping track of long-running plans, and interacting with specific users over time. Future versions of ROSA will make use of semantic and episodic memory defined in [10].

**Logic**

ROSA's internal logic mechanisms are used to (1) manage the LLM context window, and (2) parse LLM output into query responses and tool invocations. The task of managing the context window is relatively straightforward. We combine a

---

[2]Note that an action is defined as the invocation of a tool, not the tool itself.

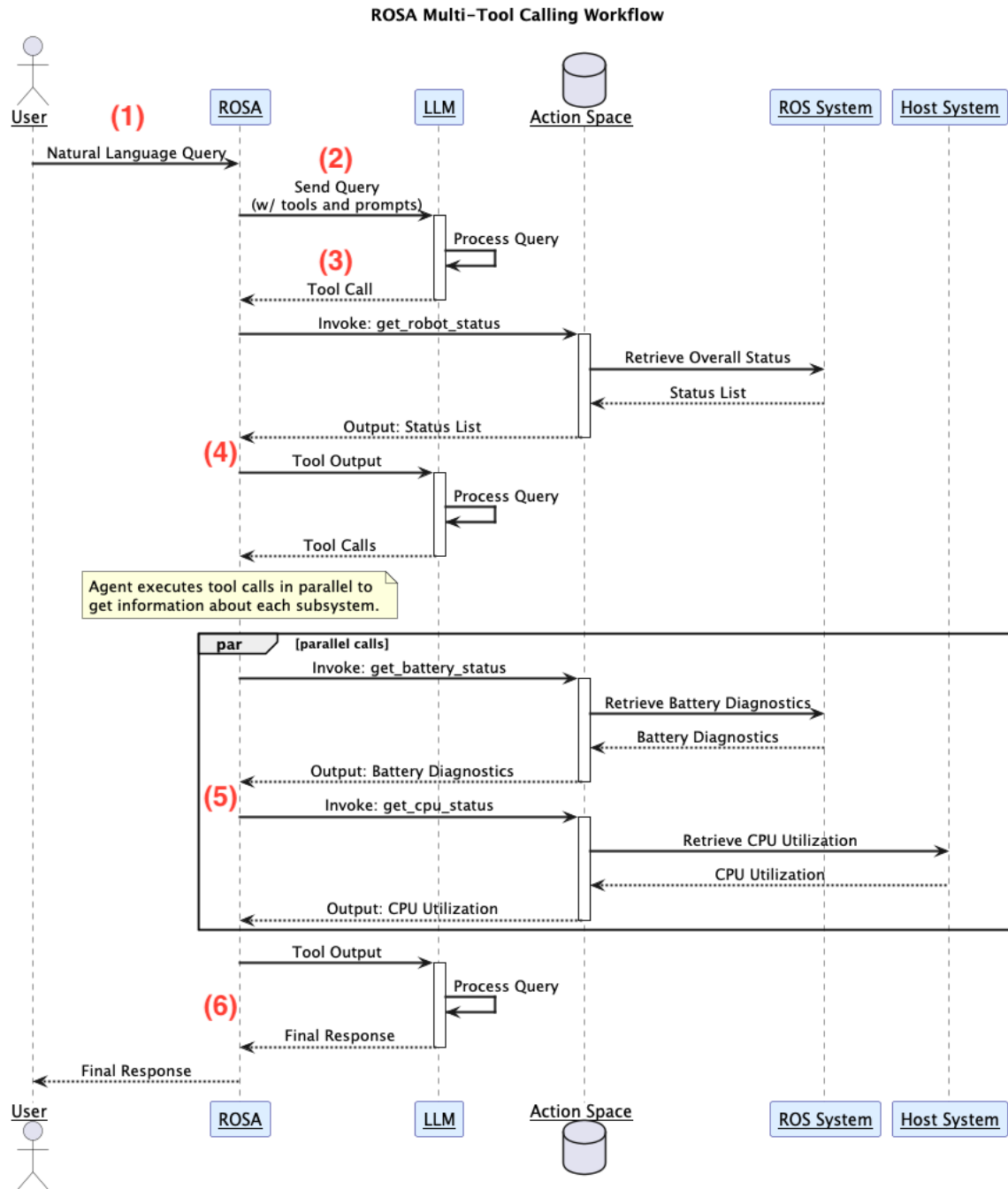[3]Text payload processed by an LLM during a single transaction.

**Figure 5**: **Sequence diagram demonstrating sequential, parallel, single-, and multi-tool calling. Initially, a single tool is called to retrieve a status overview. The model infers that certain subsystems need to be inspected further, so it retrieves detailed status information in parallel for each subsystem before returning a final result.**

set of "Robot System Prompts" (RSPs), tool specifications, agent scratchpad, and the on-going chat history, in that order, before sending a query to the LLM. As the chat history grows in length, we omit the least recently used user and assistant messages in order to stay within the maximum context length of the underlying LLM. We utilize the open-source LangChain library to parse LLM outputs, as it remains up-to-date with evolving specifications, ensuring adaptability and reliability.

*Robot System Prompts*—RSP's are a set of prompts that help guide ROSA's behavior. They are critical for ROSA's ability to maintain contextual awareness and provide coherent and relevant interactions. These prompts supply the language model with essential insights about the robot's identity, environment, and operating conditions, ensuring that ROSA's responses are aligned with the robot's specific capabilities and constraints. Furthermore, by providing detailed descriptions of the physical and digital environments the robot inhabits, ROSA can better anticipate and react to environmental changes in real-time. For instance, knowing that the physical robot is intended for use in a subterranean environment can better inform the invocation and inference of a tool used for scene understanding.

RSPs are instrumental in guiding ROSA's behavior, but they are inherently imperfect and function more as heuristics than definitive rules for the agent to follow. This imperfection is a natural consequence of the probabilistic nature of auto-regressive models, which makes it challenging to fully codify intention and establish consistent ground truth. During development, we observed discrepancies in ROSA's behavior that can largely be attributed to limitations in the design or interaction of system prompts. Contradictory or ambiguous prompts, in particular, can lead to unexpected or inconsistent responses, highlighting the sensitivity of LLMs to input framing. Future work should explore strategies to address these issues, such as uncertainty quantification and output validation techniques, to enhance the reliability and predictability of ROSA's behavior.

### Large Language Models

LLMs play a central role in enabling intelligent interaction and decision-making in ROSA's architecture. They are used for interpreting user queries, processing RSPs, and generating actionable plans. They also drive the reasoning and observation steps in the ReAct loop, allowing ROSA to infer user intentions, plan and execute appropriate actions, and evaluate outcomes based on observed results. Furthermore, LLMs enable tool calling by dynamically producing structured outputs, such as *JSON* or *XML*, that conform to predefined tool interfaces. A key design feature of ROSA is its flexibility to work with any LLM that supports tool calling, ensuring compatibility with a wide range of models, whether open-source or closed-source.

While closed-source LLMs typically deliver superior performance across diverse tasks, open-source models can perform effectively when operating within a well-defined and constrained action space. However, a major challenge in deploying open-source models for ROSA is the computational resource requirement for on-board inference. For instance, a 70-billion parameter model, even when quantized at 4-bits, demands 43GB of memory and must rely on specialized hardware, such as GPUs, to achieve processing speeds suitable for robotics. At present, the largest consumer-grade GPUs, with

a maximum of 24GB of memory, cannot meet these demands, making on-board deployment infeasible for all but the most advanced and costly systems. This limitation underscores the need for alternative solutions, such as the development of Small Language Models (SLMs) fine-tuned for ROSA-specific tasks [12], which could offer a more efficient and cost-effective approach for real-time robotics applications.

## 4. ROBOT DEMONSTRATIONS

To illustrate ROSA's capabilities, we deployed it on three distinct robotic systems, each with unique operational characteristics. These demonstrations showcase how natural language commands, processed by LLMs and dynamically integrated with ROSA's internal toolset, can simplify operations for both novice and expert users. The user interface we used for these demonstrations is shown in Figure 6. All described interactions represent real-world engagements between humans and ROSA, with actions executed to completion without incident. Throughout these examples, we blend discussion of the agent's behavior with brief references to the underlying code, showing how specific tools interact with the different robotic platforms.

### Robot Demonstration #1: NeBula-Spot

NeBula-Spot is a quadruped robot developed at JPL in response to the DARPA Subterranean Challenge. It utilizes a Boston Dynamics Spot robot enhanced with JPL's NeBula autonomy stack [13]. For this demonstration, the system was equipped with a Visual Language Model (VLM) [14] for scene understanding in JPL's Mars Yard (see Figure 1a). The users could issue high-level commands such as instructing the robot to stand, walk, or interpret camera feeds. ROSA's uplink tools (for robot actuation) and downlink tools (for diagnostics and perceptual data) enabled these capabilities, but their details were abstracted behind the conversational interface shown in Figure 6.

The following transcript provides a sample interaction with ROSA. A video of the full interaction, along with ROSA's intermediate output and tool calls, is available to view on YouTube[4].

**User:** Hey Spot, please stand up.
**ROSA:** I am now standing, what else would you like me to do?

**User:** Go ahead and walk forward about a meter and turn $15°$ to the left.
**ROSA:** I'll move forward 1 meter and turn 15 degrees to the left. Please confirm.

**User:** Yeah, it's safe. Go ahead.
**ROSA:** I have moved forward 1 meter and turned 15 degrees to the left. What's next?

**User:** Alright, can you describe what you see in the camera?
**ROSA:** I see an open sandy area with large rocks scattered around. Trees are 20 meters to the left, and buildings are 25 meters ahead to the right.

**User:** Alright, that actually sounds pretty interesting, can you show me the camera feed?
**ROSA:** The camera feed is live. What's next?

---

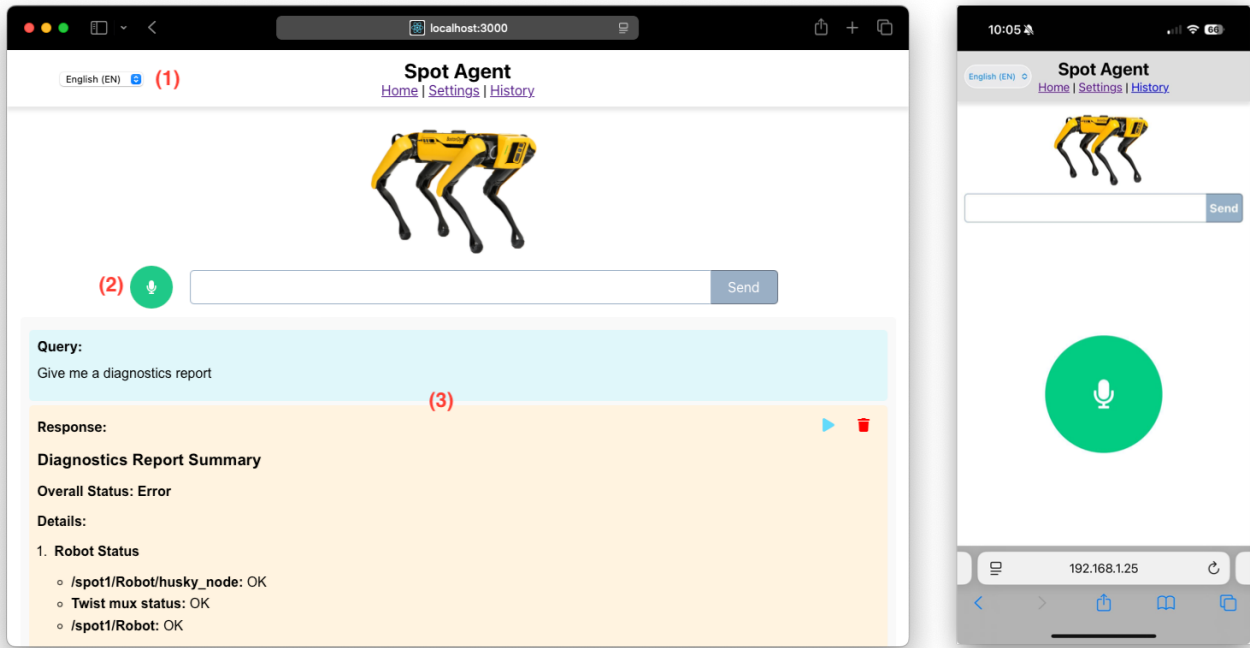[4]https://www.youtube.com/watch?v=mZTrSg7tEsA

**Figure 6**: **Locally-hosted web GUI used for interacting with ROSA during demonstrations. Users can select from English, Spanish, Polish, German, Korean, Dutch, and more (1), and query ROSA via text or speech (2). Full responses are displayed in the response field (3), or naturalized audio can be played back by pressing the play button.**

The code snippet below shows how we implemented the "stand up" functionality, effectively replicating the press of a button on a standard joystick teleop controller. From the user's perspective, it appears as though the robot is simply responding to natural language requests. This design hides the complexities of ROS topics and publishers, allowing operators to focus on mission objectives rather than lower-level implementation.

```
@tool
def stand_up() -> str:
    """Command the robot to stand."""
    args = get_default_args()
    args['buttons'][Button.B] = 1

    msg = Joy(**args)
    joy_publisher.publish(msg)

    # Validation omitted...

    return f"Spot is now standing up."
```

Additionally, at the end of the interaction, ROSA invokes the `display_camera_in_rqt` tool to bring up the ROS RQT interface, allowing the user to see a live view of the robots camera feed. Importantly, novice users and non-SME's need only specify their intent ("view the camera feed") and ROSA will take the appropriate actions to make it happen. This process essentially codifies subject-matter expertise into the agent, minimizing the need for costly and often incomplete knowledge transfer and onboarding procedures. This observation is also evident from our next demonstration on the EELS robot.

**Robot Demonstration #2: EELS**

EELS (Exobiology Extant Life Surveyor) is a serpentine robot designed for subterranean exploration, featuring adaptable gaits and advanced sensing for complex terrain [15]. For this paper, we deployed EELS in a laboratory setting (Figure 1b) and provided a more specialized toolset to handle waypoint navigation and head articulation. Although EELS can support more complex locomotion routines, we focused on simple tasks to illustrate ROSA's capacity for error handling and structured commands.

The following is a sample interaction that demonstrates ROSA's adaptability to commands that cannot be reasonably accommodated with its (purposely) constrained action space. In this case, we did not give ROSA the tools it would need to properly determine coordinates for "the rock in the corner." However, instead of refusing to carry out the task (or hallucinating the coordinates of some imaginary rock), the agent informed the user of an alternative way to specify the command:

**User:** Move toward the rock in the corner.
**ROSA:** I'm sorry, but I am not able to move toward a specific object. If you would like me to move, you can try giving me a waypoint in (x, y, $\theta$) format.

**User:** Move to the waypoint at (1, -0.2) with a $-90°$ heading.
**ROSA:** I have moved to the waypoint. However, it appears my heading is off by approximately $0.3°$. Would you like me to try again?

**User:** No, that's ok. Go ahead and raise your head.
**ROSA:** Ok, my head is raised. Would you like me to perform a scan?

**User:** Describe what you see in the camera feed.
**ROSA:** I see a laboratory environment with vaulted ceilings and flourescent lighting. To the left, I see metal shelves with electronic equipment. On the right, there is a desk, a chair, and a computer monitor. There is a person sitting at the desk, and another person nearby.

When the user asked the robot to perform an unsupported behavior, ROSA politely suggested a valid alternative (e.g., specifying a coordinate-based waypoint instead of a purely semantic location). This guided the user toward a supported command format, which reduced the need for expert intervention.

Compared to the NeBula-Spot scenario, EELS demonstrated waypoint-based navigation and fine-grained gait control. By interleaving motion commands with real-time feedback, such as self-reported heading errors, the agent showed its ability not only to execute tasks, but also to reflect the robot's internal state. This helps operators build trust in the system, as they can see it verifying the accuracy of each action and offering to correct minor deviations.

```
@tool
def raise_head():
    """Raise the head module."""
    if is_head_raised():
        return True
    head_raise_service = "/head_raise"

    rospy.wait_for_service(
        head_raise_service
    )

    service_proxy = rospy.ServiceProxy(
        head_raise_service, Empty
    )
    service_proxy()

    return True
```

When the user asked the robot to "raise [its] head," ROSA invoked a straightforward `@tool` function that calls an existing `rosservice`, as seen in the code snippet above. Invoking this functionality without ROSA would typically require the use of ROS command-line tools, an understanding of the service message format, and a verbose argument formatted in the *JSON* markup language. As such, using ROSA for these types of actions not only saves time, but also reduces the cognitive overhead required to recall and issue commands. Our next demonstration continues this trend by demonstrating how ROSA leverages the emergent reasoning capabilities of LLMs to interpret user intent and perform multi-step actions.

**Robot Demonstration #3: Carter in Simulation**

The Nvidia Carter is an autonomous mobile robot (AMR) platform that we operated entirely in simulation using NVIDIA IsaacSim (Figure 1c). We utilized a Martian terrain to explore extended sensor usage (e.g., LiDAR) and more creative camera manipulation workflows. In this environment, users issued composite commands, like scanning for obstacles before advancing, causing ROSA to chain multiple tools to achieve the desired result:

**User:** Move forward as far as you can.
**ROSA:** In order to move forward as far as I can, I need to check for obstacles and hazards. Would you like me to perform a LiDAR scan?

**User:** Yes.
**ROSA:** Ok, I have performed a LiDAR scan and see an obstacle approximately 4 meters ahead. Should I proceed?

**User:** Yes.
**ROSA:** Ok, I have moved forward by 4 meters.

**User:** Give me a 360 degree view of your surroundings.
**ROSA:** I can rotate my camera and take intermittent snapshots. How many snapshots should I take?

**User:** Assume the camera FoV is 90 degrees.
**ROSA:** Based on a 90 degree field of view, I have captured 4 snapshots that show a 360 degree view of my surroundings.

Below is a simplified snippet of the tool we created for rotating the on-board camera:

```
@tool
def rotate_camera(angle: float) -> str:
    """
    Rotate the on-board camera sensor by a
    ↪  specified angle.

    Args:
        angle: angle to rotate in degrees.
    """
    angle_rad = math.radians(angle)

    rotate_msg = CarterCameraRotate(
        angle_rad=angle_rad
    )

    carter_camera_pub.publish(rotate_msg)

    # Validation omitted...

    return f"Camera rotated by {angle}
    ↪  degrees."
```

This relatively trivial ROS wrapper led to a robust user experience: the agent reasoned about field-of-view (FoV) constraints, realized it needed to rotate four times by 90 degrees, and captured images after each rotation to create a mosaic of the environment. Even though the user's initial request was vague, ROSA adapted with a practical solution.

In simulation, Carter allowed us to put ROSA's reasoning capabilities to the test by combining LiDAR scans with collision checks and multi-step instructions. Because everything was virtual, novice operators felt comfortable experimenting with more advanced commands knowing they wouldn't damage physical hardware. The demonstration solidified how ROSA's abstraction layers can generalize to new operational constraints, whether real or simulated, and help users focus on high-level mission planning.

**Summary of Findings**

From terrestrial testyards (NeBula-Spot) and indoor lab spaces (EELS) to virtual Martian terrains (Carter), ROSA provided a consistent, natural-language interface for diverse platforms. Each demonstration showed how the agent can:

- Translate colloquial or imprecise language into valid ROS actions,

- Prompt for clarifications when commands exceed its known toolset,
- Offer context-dependent information to supplement the users operational understanding, and
- Reduce the cognitive load required to operate robotic systems of varying complexity.

Across these varied robotic systems and operational environments, we observed novices quickly grasp how to "talk" to ROSA, often requiring only minimal guidance on the possible commands. In almost all cases, ambiguity was resolved by simply conversing with the agent. Meanwhile, the agent seamlessly handled tool invocations in the background while maintaining a fluid conversation with operators. In the next section, we explore how ROSA enhances human-robot interaction, highlighting its current capabilities, future potential, and areas for further refinement in our methodology.

## 5. HUMAN-ROBOT INTERACTION

Effective human-robot interaction is essential for maximizing the utility and accessibility of robotic systems. Most robotic systems are inherently complex, often developed by interdisciplinary teams where each member specializes in specific domains [16]. This specialization can lead to fragmented knowledge, making it difficult for individual operators to fully comprehend or manage the entire system. Advanced robots that integrate multiple subsystems, sensors, and actuators further exacerbate this complexity, creating significant barriers for routine operation and troubleshooting.

ROSA addresses these challenges by serving as an intelligent intermediary, bridging the gap between operators and robotic systems. Using LLMs, ROSA consolidates diverse sources of documentation, requirements, specifications, and technical knowledge into a single, accessible interface. Through natural language interaction, ROSA empowers operators to access comprehensive system knowledge without requiring in-depth expertise. This capability provides step-by-step guidance for tasks, aids in troubleshooting, and enables a clear understanding of system functionalities, limitations, and operational parameters. By reducing reliance on subject matter experts (SMEs) for routine operations, ROSA streamlines workflows and enables SMEs to focus on high-value tasks like system development and optimization.

### Multimodal Interactions

ROSA's default mode of operation is through text-based interaction. This can take the form of a command-line interface (CLI) using Python's Read-Evaluate-Print-Loop (REPL) or a custom graphical user interface (GUI) as in Figure 6. Additionally, ROSA supports integration with dedicated supervisory interfaces for multi-robot systems [17]. While text-based interaction forms the foundation, ROSA's architecture is designed to accommodate multiple input and output modalities, enhancing HRI across diverse scenarios.

*Speech Integration*—Speech is a natural and intuitive medium for communication, making it a valuable extension of ROSA's interaction capabilities. By integrating speech-to-text (STT) and text-to-speech (TTS) libraries available in Python, ROSA allows operators to issue verbal commands and receive audible responses. This process involves capturing audio input, transcribing it to text with an STT module, and passing the resulting text to ROSA. The system's responses are then synthesized into speech using a TTS module.

This hands-free interaction is particularly useful in scenarios where operators' hands are occupied, such as when conducting field tests using wearable devices. Furthermore, it enhances accessibility for individuals with visual or motor impairments, making robotic systems more inclusive and practical in a broader range of contexts.

*Visual Perception and Image Recognition*—ROSA enhances situational awareness and decision-making by processing visual data from the robot's cameras. Using VLMs, ROSA can interpret and describe environments and identify objects, obstacles, and landmarks within the camera's field of view. Depth imaging allows ROSA to estimate distances and provide spatial insights, such as proximity to obstacles or the layout of surroundings, which is particularly useful for navigation in complex or unfamiliar environments.

For example, ROSA can analyze camera feeds to detect hazards in industrial settings or provide detailed descriptions of terrain for autonomous navigation. These capabilities are achieved by combining multimodal models, such as `gpt-4o`, with computer vision libraries like OpenCV. By transforming raw visual data into actionable insights, ROSA enables operators to make informed decisions with greater confidence and efficiency.

*Multi-Lingual Collaboration*—ROSA's design leverages the capabilities of modern language models, many of which support a wide range of languages, including English, Korean, Hindi, German, Spanish, and more. This multi-lingual functionality allows ROSA to operate in diverse linguistic contexts, making it adaptable for global applications. Operators can issue commands, ask questions, and receive responses in their preferred language, reducing barriers to adoption and broadening accessibility.

By tapping into these multi-lingual capabilities, ROSA can support international teams, cross-cultural collaboration, and deployments in regions where English is not the primary language. This feature ensures that advanced robotic systems are usable by a wide variety of operators, regardless of their linguistic background. Furthermore, ROSA's ability to operate across languages enhances inclusivity and expands its potential impact across industries and geographies.

### Limitations and Drawbacks

While ROSA introduces a streamlined and accessible interface for HRI, it also brings several important limitations and risks that must be acknowledged. One key limitation arises from the heavy computational resources required by LLMs. Even when leveraging local models optimized for edge devices, the memory footprint and real-time processing demands can be prohibitive, particularly in constrained spacecraft environments. These computational overheads may necessitate specialized hardware accelerators, adding weight, power consumption, and cost to already resource-limited missions. In addition, the inference time required by complex models may introduce latency, potentially impacting time-sensitive applications such as rapid maneuvers or collision avoidance.

Another drawback stems from the intrinsic nature of auto-regressive models. Despite impressive performance, LLMs can exhibit imprecise or "hallucinated" responses–

confidently providing incorrect or irrelevant information. This risk is amplified in high-stakes scenarios such as spacecraft operations, where errors can have serious or irreversible consequences. Mitigation strategies like human-in-the-loop verification, rule-based filters, or ensemble-based approaches can reduce but not entirely eliminate these risks. Furthermore, the difficulty in fully explaining or interpreting LLM decision processes can complicate system certification and operator trust, particularly in aerospace contexts with stringent reliability and safety requirements.

Beyond these technical considerations, ROSA's reliance on robust training data raises challenges for knowledge maintenance and domain adaptation. Spacecraft systems evolve over time, and models must be regularly updated to reflect new design revisions, software patches, or mission parameters. Ensuring that the AI remains synchronized with the current state of the system demands continuous data curation and validation. These processes add complexity to mission planning and may require specialized personnel, effectively reintroducing some level of reliance on subject matter experts. While ROSA lowers operational barriers in theory, its real-world deployment involves careful planning and ongoing oversight to mitigate both computational and epistemic pitfalls. In the next section, we detail ROSA's implementation, including various techniques we use to mitigate some of these limitations.

## 6. IMPLEMENTATION DETAILS

ROSA is implemented in Python and utilizes either ROS1 or ROS2 for communication with robotic systems. This section outlines the core components of ROSA's implementation, including the design and structure of its tools, approaches to parameterization, integration with ROS functionalities, and considerations for selecting and configuring large language models. By focusing on these technical aspects, we aim to provide a comprehensive understanding of how ROSA operates in practice and how its modular design supports the diverse capabilities demonstrated in the previous section. A subset of the ROSA library is publicly available at `https://github.com/nasa-jpl/rosa`.

### Tool Structure and Organization

At the core of ROSA are the *tools*—Python functions that encapsulate various functionalities for robotics development and operation, including a large subset of standard ROS tools. These tools are designed to be modular, reusable, and easily integrated, enabling the agent to interact effectively with the ROS environment. To simplify the creation and management of tools, ROSA leverages the LangChain framework, which provides an intuitive abstraction for defining tools. Specifically, developers can use the `@tool` decorator provided by LangChain to create new tools with minimal effort, ensuring extensibility and consistency across the system.

Tools are organized into modules based on their functionality and ROS version compatibility. Separate modules, `ros1.py` and `ros2.py`, are used to support ROS1 and ROS2, respectively. This separation avoids conflicts between versions and facilitates maintenance and extension of the toolset. Additionally, common tools, such as utility functions, are provided in modules like `calculation.py` and `log.py`, ensuring applicability across both ROS versions.

```python
from langchain.agents import tool

@tool
def rosnode_list(
    pattern: Optional[str] = None,
    namespace: Optional[str] = None
) -> dict:
    """
    Returns a dictionary containing a list
    ↪   of running ROS nodes and other
    ↪   metadata.

    :param pattern: a regex pattern to...
    :param namespace: ROS namespace to...
    """
    # Implementation details...
```

**Figure 7**: **Example implementation of a basic ROSA tool, which must include the `@tool` decorator, a docstring, descriptive method and parameter names, and a return value that can resolve to a string.**

Each tool function is registered using the `@tool` decorator from LangChain (Figure 7). This decorator allows the function to be invoked by the language model agent. By standardizing tool definitions and registration, ROSA ensures consistent behavior and seamless integration with LangChain, and by extension, the underlying language model. This approach simplifies adding new tools and promotes a clear and maintainable code structure.

### Parameterization and Input Handling

To allow flexible and efficient interaction with the ROS environment, ROSA's tools are designed to accept well-defined parameters, often with default values and optional arguments. This approach enables the language model to perform complex queries and operations by selectively specifying the necessary parameters. Many tools accept a `pattern` parameter, allowing users to specify regular expressions to filter results. This is particularly useful when dealing with a large number of nodes or topics, as it enables targeted queries. The `namespace` parameter enables tools to operate within specific ROS namespaces, providing granular control over the scope of operations.

Additionally, a `blacklist` parameter can be used to exclude certain nodes, topics, or services from the output, enhancing safety and customization by preventing unintended interactions with critical system components. To streamline the inclusion of default blacklists and ensure consistent safety measures across tools, ROSA implements an `inject_blacklist` decorator (Figure 8). This decorator automatically injects any user-provided blacklist into tools that have a `blacklist` parameter, reducing the potential for human error and simplifying tool definitions. By centralizing safety mechanisms like a blacklist, ROSA maintains a high standard of operational safety without burdening individual tool implementations. This design decision enhances both usability and security.

### Coverage of ROS Functionalities

ROSA provides a wide range of ROS tools covering common tasks related to nodes, topics, services, and parameters. For nodes, it offers commands to list (Figure 9), inspect, and

```python
def inject_blacklist(global_blacklist):
    def decorator(func):
        @wraps(func)
        def wrapper(*args, **kwargs):
            # Implementation details here
            return func(*args, **kwargs)
        return wrapper
    return decorator
```

**Figure 8**: Implementation of the `inject_blacklist` decorator for ROSA tools. This decorator is automatically applied to any tools passed to ROSA during instantiation, provided they have a `blacklist` parameter. The agent may also choose to add additional values to the blacklist, in which case the union of the two lists (global and agent) is taken.

terminate active nodes. Topic tools support listing, echoing messages, and retrieving topic details. Services can be queried, inspected for types, and called with parameters. Parameter management allows listing, getting, and setting values for dynamic system configuration.

ROSA also includes commands for interacting with packages (e.g., listing launch files and starting processes) and for log handling (e.g., reading and filtering ROS logs). Additionally, utility tools enable data processing and numerical computations. These capabilities let users manage the ROS environment through concise natural-language commands, improving accessibility and efficiency.

### Structured Data Outputs

To ensure that the language model can generate accurate and contextually appropriate responses, ROSA's tools typically return structured data in the form of dictionaries or lists (Figure 10). This approach facilitates easy parsing and interpretation of the outputs. Structured tool responses augment the context window with important information that helps prevent confabulations by the language model.

For example, when a tool returns a list of nodes along with the total count, it provides explicit data that the language model can use to answer questions like "How many localization nodes are running?" without needing to infer or estimate the answer. This reduces the likelihood of the language model providing incorrect or fabricated responses. By including contextual information such as patterns used, namespaces, and counts, the tools provide comprehensive data that informs downstream tasks and interactions. This design choice enhances the reliability and usefulness of the language model's responses.

### Integration with LangChain and the ReAct Framework

ROSA integrates with the LangChain library to implement the ReAct (Reasoning and Acting) framework, allowing the language model agent to process natural language inputs, reason about appropriate actions, and invoke the corresponding tools. LangChain handles prompt management by facilitating prompt chaining, which is used to guide the language model's understanding of ROSA's capabilities, constraints, and the context of the interaction. It also manages conversational memory, enabling ROSA to maintain state across interactions and provide coherent, context-aware responses.

```python
@tool
def rosnode_list(
    pattern: Optional[str] = None,
    namespace: Optional[str] = None,
    blacklist: List[str] = None
) -> dict:
    """
    Returns a list of running ROS nodes
    ↪  with optional filtering.

    :param pattern: A regex pattern to
    ↪  filter nodes.
    :param namespace: ROS namespace to
    ↪  scope the search.
    :param blacklist: List of node names to
    ↪  exclude.
    """
    nodes = rosnode.get_node_names()
    # Apply namespace filtering
    if namespace:
        nodes = [
            n for n in nodes if
            ↪  n.startswith(namespace)
        ]
    # Apply pattern filtering
    if pattern:
        nodes = [
            n for n in nodes if
            ↪  re.match(pattern, n)
        ]
    # Apply blacklist
    if blacklist:
        nodes = [
            n for n in nodes if n not in
            ↪  blacklist
        ]
    return dict(nodes=nodes,
                total=len(nodes),
                namespace=namespace,
                pattern=pattern)
```

**Figure 9**: Implementation of the ROS1 tool that retrieves a list of ROS nodes, applies filters, and returns structured output. Note that the blacklist parameter is automatically populated when the ROSA class is instantiated.

The ReAct framework facilitates the logic of selecting and executing tools based on the language model's reasoning. When the model determines that an action is required to answer a query or fulfill a command, it can invoke the appropriate tool, pass the necessary parameters, and process the returned structured data. This integration streamlines the interaction between the language model and the ROS environment, allowing for seamless execution of commands and retrieval of information. It also abstracts the complexity of tool invocation, making it easier to extend ROSA's capabilities.

System prompts [18] (Figure 11) play a crucial role in this integration by providing the language model with instructions and context. These prompts define the agent's persona, capabilities, limitations, and critical instructions, ensuring that the language model operates within the intended parameters. An example of system prompts is shown in Figure-11.

### Model Selection and Trade-offs

The choice of the language model significantly impacts ROSA's performance, capabilities, and resource require-

```
{
    "nodes": ["/talker", "/listener"],
    "total": 2,
    "namespace": "/",
    "pattern": ".*"
}
```

**Figure 10**: **Example tool response using structured output, which helps guide the model, reduce hallucinations, and provide traceability for agent actions and responses.**

```
system_prompts = [
    ("system", "You are ROSA, an AI agent
    ↪   that uses ROS tools to..."),
    # Additional prompts...
]
```

**Figure 11**: **Example system prompt following the OpenAI API standard. The list of system prompts can contain any number of prompts to help guide the agents behaviors, responses, etc.**

```
@tool
def read_log(
    log_file_directory: str,
    log_filename: str,
    level_filter: Optional[str] = None,
    num_lines: Optional[int] = None
) -> dict:
    """
    Reads a log file and returns entries
    ↪   matching the specified criteria.

    :param log_file_directory: directory...
    :param log_filename: name of the...
    :param level_filter: filter logs by...
    :param num_lines: number of lines...
    """
    # Implementation details...
```

**Figure 13**: **Example function signature for a tool that extends beyond ROS interactions. This particular tool can be used to read log files, allowing ROSA to find errors, warnings, and other commonly logged and informative information.**

ments. ROSA currently supports chat models like *GPT-4o*, *Claude 3.5 Sonnet*, and *Llama 3.2*. The selection between API-based and local inference models depends on the specific needs and constraints of the deployment environment. There are two requirements a model must satisfy in order to be used with ROSA:

1. It must support *tool calling* [19].
2. It must have a context length of no less than *8192* tokens

For API inference, models like *GPT-4o* are recommended due to their advanced reasoning capabilities and lower cost per token. These models offer superior performance in understanding complex queries and generating accurate responses. However, they require internet connectivity and may raise concerns regarding data privacy and dependency on external services. For local inference, models like *Llama 3.1 8B* provide a balance between performance and resource requirements. Running models locally enhances privacy and reduces reliance on external services, which is beneficial in secure or resource-constrained environments. However, local models may require optimization to achieve real-time performance and may have limitations in handling long context windows or complex reasoning tasks.

```
@tool
def add_all(numbers: List[float]) -> float:
    """Returns the sum of numbers."""
    return sum(numbers)

@tool
def mean(numbers: List[float]) -> dict:
    """Returns the mean and standard
    ↪   deviation of a list of numbers."""
    return {
        "mean": statistics.mean(numbers),
        "stdev": statistics.stdev(numbers),
    }
```

**Figure 12**: **Examples of custom tools that go beyond interacting with ROS. These tools help avoid hallucinations by giving the Agent a means to run calculations, rather than relying on fallible notions of calculation from the models training data.**

When selecting a model, trade-offs between context length, performance, resources, and deployment considerations must be evaluated. Larger models typically offer better performance but require more computational resources. Models with longer context windows can handle extended conversations but may consume more memory and processing power. Deployment considerations, such as the need for offline operation or data privacy requirements, may influence the choice of a local model over an API-based one.

**Design Choices**

Several key design decisions were made during the implementation of ROSA to enhance usability, safety, and effectiveness. The modular architecture–separating tool packages into ROS1, ROS2, calculation (Figure 12), logging (Figure 13), and other modules–simplifies maintenance and allows for clear separation of version-specific functionalities. This organization facilitates easier updates and scalability as ROS evolves.

By automating aspects like parameter injection and tool registration, the ROSA library reduces the potential for errors and streamlines the development process by design. Parameterization of tools with options like patterns, name spaces, and blacklists provides flexibility while maintaining control over operations. Allowing the language model to specify these parameters enables ROSA to perform precise actions as instructed by the user, enhancing the agent's usefulness in varied scenarios. Returning structured data enhances the language model's ability to generate accurate responses and reduces the risk of confabulation. This design choice is critical in ensuring that operators receive reliable information, thereby increasing trust in the system.

The integration with LangChain and the ReAct framework abstracts the complexity of tool invocation and reasoning processes, allowing developers to focus on extending ROSA's capabilities rather than managing low-level interactions. This abstraction layer simplifies the agent's architecture and promotes more rapid development and experimentation. Prioritizing safety through parameter validation, constraint enforcement, and centralized safety mechanisms addresses po-

tential risks associated with executing commands in the ROS environment.

# 7. ETHICS FOR EMBODIED AGENTS

The development and deployment of embodied agents, such as ROSA, raise significant ethical considerations. As robots gain autonomy and interact with complex environments, it becomes crucial to define frameworks that promote ethical behavior, protect human welfare, and maintain public trust. This section outlines principles, practices, and open questions that guide the safe and responsible adoption of such technologies, with particular emphasis on Asimov's Three Laws of Robotics.

**Foundational Ethical Principles and Asimov's Laws**

Isaac Asimov's *Three Laws of Robotics* [20] provide a widely recognized (albeit fictional) starting point for ethical guidelines in robotics:

**Law 1:** A robot may not injure a human being or, through inaction, allow a human being to come to harm.

**Law 2:** A robot must obey orders given by human beings except where such orders would conflict with the First Law.

**Law 3:** A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.

Practical implementations of these laws in real-world systems must contend with complexities such as interpreting user intent, balancing trade-offs in uncertain environments, and enforcing directives that may occasionally conflict. Below, we discuss how the ROSA framework addresses—or in some cases does not fully address—each of these laws.

*Addressing the First Law (Preventing Harm to Humans)*— ROSA is primarily an AI-agent framework that coordinates high-level decision-making. It does not, by default, include active "uplink" modules for physical control; only "downlink" modules for system observation are provided out of the box. Thus, the framework encourages developers to implement their own control pathways, including mechanisms to detect and avoid human injury. Some ways in which ROSA supports the First Law are:

- **Mandatory E-Stop and Overrides:** ROSA-based systems can (and should) be integrated with external e-stop devices or software overrides. These can immediately halt any action, mitigating the possibility of harm due to misinterpretation of data or erroneous commands.
- **Fail-Safe Defaults:** In ambiguous or hazardous conditions, ROSA can revert to a safe or idle state until it receives clarification from a human operator. This behavior aims to reduce the risk of unintended harmful actions.
- **Encouraged Integration with Safety Layers:** Although not enforced within ROSA itself, developers are strongly advised to employ lower-level ROS packages for human detection, collision avoidance, workspace monitoring, and other safety-critical functions.

However, ROSA alone cannot guarantee compliance with the First Law unless it is coupled with the appropriate sensing and hardware safety features. Responsibility ultimately lies with developers and operators to integrate these capabilities effectively.

*Addressing the Second Law (Obedience to Humans)*—ROSA supports a hierarchical command structure in which human operators retain ultimate control. By default, any directive from a human user supersedes the AI agent's plan. This can be realized through:

- **Command Multiplexers:** These modules listen to both ROSA-generated commands and human input (e.g., joystick or teleoperation). If a conflict arises, human commands take precedence and ROSA's commands are overridden.
- **Explainable Decision-Making:** ROSA logs and exposes its decision-making traces, enabling operators to review and revise or halt upcoming actions. This transparency ensures that human oversight remains effective and informed.

Despite these measures, it remains possible for humans to issue harmful or contradictory commands. ROSA cannot independently resolve deliberate misuse or ill-intentioned directives. Additional policy-level controls or organizational oversight are advised to address such scenarios.

*Addressing the Third Law (Preserving the System's Existence)*—ROSA includes a range of self-diagnostic and monitoring tools that enable an agent to detect and report anomalies in the environment or in the robot's hardware. Within an integrated robotic system, ROSA can:

- **Monitor Critical System Health:** The agent can be programmed to detect sensor failures, battery anomalies, or actuator malfunctions and then engage mitigation behaviors (e.g., moving to a safe state, shutting down power to certain subsystems).
- **Prioritize Self-Preservation Consistent with Human Safety:** In the event of certain high-risk conditions, ROSA can pause all non-critical tasks and focus on preserving the robot's functional integrity, as long as doing so does not conflict with higher-level safety concerns for humans.

This layered approach ensures that the agent can protect itself, but not at the expense of human safety or legitimate human instructions.

**Safety and Risk Mitigation**

*Human Override and E-Stop*—Human overrides and e-stop features are essential for mitigating accidental harm. ROSA allows for external e-stop signals that immediately suspend or cancel robotic motion. Command multiplexers or hardware interlocks can force all high-level commands to cease. These practices align with the First Law by preventing unwanted or potentially harmful actions.

*System-Specific Safety Considerations*— While ROSA can interface with standard ROS packages for collision avoidance and human-detection, it does not enforce these by default. Each deployment must integrate low-level safety features that match its specific application, such as sensors for proximity alerts, vision-based human tracking, and compliance-based motion controllers. ROSA's high-level autonomy thus complements rather than replaces existing safety modules.

*Gaps and Future Enhancements*—Although ROSA reduces risk by separating "downlink" (observation) from "uplink"

(actuation) capabilities, there is no absolute technical barrier to misuse if an operator creates modules that issue harmful commands. Future iterations of ROSA may introduce permission layers, code review pipelines, or curated tool repositories that clearly demarcate "trusted" from "untrusted" functionalities. More rigorous user authentication and code signing could also help mitigate malicious or accidental misuse.

*Open-Source Governance and Ethical Oversight*—ROSA is distributed as open-source software, which fosters collaboration and rapid innovation but also necessitates diligent governance. Currently, repository administrators review pull requests, and submissions that enable harmful behaviors are disallowed. Future plans include an ethics-focused discussion forum, more explicit contributor guidelines, and automated checks for dangerous functionalities. These measures aim to prevent the addition of modules that could violate Asimov's Laws or general ethical norms.

**Operability and Need for Extended Case Studies**

ROSA's present design has been validated in small scale proofs of concept, but more comprehensive multi-robot or long-duration experiments will be necessary to quantify improvements in efficiency, safety, and reliability. Proposed metrics include:

- **Mean Time Between Human Interventions (MTBHI):** A measure of how often operators must manually override the system.
- **Task Completion Rate:** The fraction of tasks successfully completed under normal autonomous conditions.
- **Incident Reporting Rate:** The frequency of safety-related events or near-misses.

Such metrics could clarify the impact of the ROSA framework on real-world applications and demonstrate how effectively it upholds Asimov's principles in practice. Longer-term or field-based case studies will also provide opportunities to refine system components, expand the safety feature set, and address evolving ethical challenges.

In summary, while ROSA embraces the spirit of Asimov's Laws by prioritizing human safety, respecting user directives, and enabling self-protective features, complete adherence to these laws depends on hardware integration, robust sensing, and diligent governance practices. Ongoing enhancements to framework features, coupled with rigorous field testing, will improve both operability and compliance with these foundational principles.

## 8. CONCLUSION

In this paper, we introduced ROSA, a novel AI agent that bridges the gap between robotic systems and human operators through natural language interfaces. By leveraging advanced language models and open-source frameworks, ROSA enables users to interact with complex robotic platforms using everyday language, thereby enhancing both accessibility and operational efficiency. We described ROSA's agent architecture, detailing its action space, system prompts, and tool invocation mechanisms. The modular and cross-compatible design supports integration with both ROS1 and ROS2, underscoring its adaptability. Furthermore, structured tools and safety mechanisms ensure that ROSA operates reliably and securely within diverse ROS environments.

ROSA represents a major step forward in human-robot interaction, offering operators immediate access to extensive system knowledge and multimodal communication capabilities, including speech, visual perception, and multilingual usage. By addressing the challenges arising from the complexity and specialization of robotic systems, ROSA empowers users with varying levels of expertise to interact more effectively with robots. We also explored the ethical considerations pertinent to AI agents in robotics, emphasizing the need for safety, transparency, privacy, and responsible usage. By integrating foundational ethical principles and proposing community guidelines, we aim to support the creation of robust ethical frameworks that can guide future innovations in AI and robotics.

Looking ahead, ROSA's development can be advanced in several ways. Future efforts include improving contextual understanding, refining disambiguation strategies, and strengthening the agent's ethical decision-making processes. Investigations into fine-tuned small language models may reduce reliance on hosted models, thereby enhancing privacy and expanding offline capabilities, although challenges related to latency, model performance, and the increased hardware costs of local deployment must be addressed. Furthermore, incorporating semantic and episodic memory, as inspired by the Cognitive Language Agent methodologies, could enable ROSA to maintain more comprehensive context, make better-informed decisions, and learn more effectively from its mistakes. By fostering cross-disciplinary collaboration and upholding ethical standards, we envision a future in which embodied agents like ROSA play an integral role in advancing robotics, benefiting society at large, and promoting responsible technological progress.

## REFERENCES

[1] A. Koubaa, A. Ammar, and W. Boulila, "Next-generation human-robot interaction with ChatGPT and robot operating system," *Software: Practice and Experience*, vol. n/a, no. n/a, 2024. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1002/spe.3377

[2] S. Hoque, F. F. Riya, and J. Sun, "HRI challenges influencing low usage of robotic systems in disaster response and rescue operations," 2024. [Online]. Available: https://arxiv.org/abs/2401.15760

[3] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, A. Y. Ng *et al.*, "ROS: an open-source robot operating system," in *ICRA workshop on open source software*, vol. 3, no. 3.2. Kobe, Japan, 2009, p. 5.

[4] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao, "React: Synergizing reasoning and acting in language models," 2023. [Online]. Available: https://arxiv.org/abs/2210.03629

[5] S. Macenski, T. Foote, B. Gerkey, C. Lalancette, and W. Woodall, "Robot Operating System 2: Design, architecture, and uses in the wild," *Science Robotics*, vol. 7, no. 66, 2022. [Online]. Available: https://www.science.org/doi/abs/10.1126/scirobotics.abm6074

[6] D. H. Hagos, R. Battle, and D. B. Rawat, "Recent advances in generative ai and large language models: Current status, challenges, and perspectives," *IEEE Transactions on Artificial Intelligence*, pp. 1–21, 2024.

[7] A. Bucker, L. Figueredo, S. Haddadin, A. Kapoor, S. Ma, S. Vemprala, and R. Bonatti, "Latte: Language trajectory transformer," 2022. [Online]. Available: https://arxiv.org/abs/2208.02918

[8] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch, "Language models as zero-shot planners: Extracting actionable knowledge for embodied agents," 2022. [Online]. Available: https://arxiv.org/abs/2201.07207

[9] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," 2020. [Online]. Available: https://arxiv.org/abs/2005.14165

[10] T. R. Sumers, S. Yao, K. Narasimhan, and T. L. Griffiths, "Cognitive architectures for language agents," 2024. [Online]. Available: https://arxiv.org/abs/2309.02427

[11] H. Zhang, W. Du, J. Shan, Q. Zhou, Y. Du, J. B. Tenenbaum, T. Shu, and C. Gan, "Building cooperative embodied agents modularly with large language models," 2024. [Online]. Available: https://arxiv.org/abs/2307.02485

[12] M. J. J. Bucher and M. Martini, "Fine-tuned 'small' LLMs (still) significantly outperform zero-shot generative ai models in text classification," 2024. [Online]. Available: https://arxiv.org/abs/2406.08660

[13] B. Morrell, K. Otsu, A. Agha *et al.*, "An Addendum to NeBula: Towards Extending TEAM CoSTAR's Solution to Larger Scale Environments," *IEEE Transactions on Field Robotics*, 2024. [Online]. Available: https://doi.org/10.1109/TFR.2024.3430891

[14] F. Bordes, R. Y. Pang, A. Ajay, A. C. Li, A. Bardes, S. Petryk, O. Mañas, Z. Lin, A. Mahmoud, B. Jayaraman *et al.*, "An introduction to vision-language modeling," 2024. [Online]. Available: https://arxiv.org/abs/2405.17247

[15] T. S. Vaquero, G. Daddi, R. Thakker, M. Paton, A. Jasour, M. P. Strub, R. M. Swan, R. Royce, M. Gildner, P. Tosi, M. Veismann, P. Gavrilov, E. Marteau, J. Bowkett, D. L. de Mola Lemus, Y. Nakka, B. Hockman, A. Orekhov, T. D. Hasseler, C. Leake, B. Nuernberger, P. Proença, W. Reid, W. Talbot, N. Georgiev, T. Pailevanian, A. Archanian, E. Ambrose, J. Jasper, R. Etheredge, C. Roman, D. Levine, K. Otsu, S. Yearicks, H. Melikyan, R. R. Rieber, K. Carpenter, J. Nash, A. Jain, L. Shiraishi, M. Robinson, M. Travers, H. Choset, J. Burdick, A. Gardner, M. Cable, M. Ingham, and M. Ono, "EELS: Autonomous snake-like robot with task and motion planning capabilities for ice world exploration," *Science Robotics*, vol. 9, no. 88, 2024. [Online]. Available: https://www.science.org/doi/abs/10.1126/scirobotics.adh8332

[16] T. H. Chung, V. Orekhov, and A. Maio, "Into the robotic depths: Analysis and insights from the DARPA subterranean challenge," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 6, no. Volume 6, 2023, pp. 477–502, 2023. [Online]. Available: https://www.annualreviews.org/content/journals/10.1146/annurev-control-062722-100728

[17] M. Kaufmann, R. Trybula, R. Stonebraker, M. Milano, G. J. Correa, T. S. Vaquero, K. Otsu, A.-A. Agha-Mohammadi, and G. Beltrame, "Copiloting autonomous multi-robot missions: A game-inspired supervisory control interface," *arXiv*, 2022. [Online]. Available: https://doi.org/10.48550/arXiv.2204.06647

[18] OpenAI, "Prompt engineering," https://platform.openai.com/docs/guides/prompt-engineering/strategy-write-clear-instructions, accessed: 2024-09-16.

[19] LangChain Developers, "LangChain Chat Models," https://python.langchain.com/v0.1/docs/integrations/chat/, accessed: 2024-10-01.

[20] I. Asimov, "Runaround," https://web.williams.edu/Mathematics/sjmiller/public_html/105Sp10/handouts/Runaround.html, accessed: 2024-09-16.

## BIOGRAPHY

**Rob Royce** received his degree in Computer Science and Engineering from the University of California, Los Angeles (UCLA), where he was admitted as a Regents Scholar. He is currently a Data Scientist in the Artificial Intelligence & Analytics Group at NASA's Jet Propulsion Laboratory (JPL) and serves as the Systems Integration Lead for various robotics projects, including NAISR/NeBula and EELS. Rob is a leading practitioner and subject matter expert in Generative AI (GenAI) at JPL, and is also registered as a NASA AI Advisor. He has worked on numerous innovative GenAI projects, such as utilizing retrieval-augmented generation to analyze vast datasets and creating applications that enable users to interact with their documents through conversational interfaces. His extensive experience in software engineering, artificial intelligence, machine learning, and robotics has culminated in the development of ROSA, the Robot Operating System Agent.

**Marcel Kaufmann** is a Data Scientist in the Future Technology Exploration and Infusion Group at NASA's Jet Propulsion Laboratory (JPL). His research focuses on enabling innovative technologies in multi-robot systems, human-robot interaction, remote sensing, and data visualization. He holds a Ph.D. (2024) in Computer Engineering from Polytechnique Montreal, Canada, and an M.Sc. (2016) and a B.Sc. (2015) in Photonics and Computer Vision from the University of Applied Sciences Darmstadt, Germany. He previously developed an autonomy assistant for multi-robot operations called Copilot MIKE as part of the

*DARPA SubT Challenge. He is also a former Vanier Canada Graduate Scholar (2019) and an alumnus of the International Space University's Space Studies Program (2017), hosted by the Cork Institute of Technology in Ireland.*



**Jonathan Becktor** *is a Robotics Technologist in the Maritime and Multiagent Systems Group at NASA's Jet Propulsion Laboratory (JPL). His research centers on developing advanced autonomy and perception systems for multiagent robotics, with key roles in the NAISR/NeBula program and the DARPA LINC challenge. Jonathan's work focuses on robust object detection, semantic mapping, and integrating machine learning techniques for autonomous systems. He holds a Ph.D. (2023) in Computer Science and Electrical Engineering from the Technical University of Denmark. Prior to JPL, he worked as a machine learning engineer and conducted research in perception technologies for robotics.*



**Sangwoo Moon** *is a Postdoctoral Fellow in the Mobility and Robotic Systems at NASA Jet Propulsion Laboratory (JPL). His primary focus is on developing target belief model representation from multi-modal sensor data and leveraging learning-based decision making for proactive exploration and target localization using vision-language-action models (VLAM). His research interests encompass uncertainty-aware and information-driven robotic mission and path planning, capturing information measure in stochastic environments, and communication-aware informative planning for multi-robot systems. He received a Ph.D. from the Ann and H.J.Smead Aerospace Engineering Sciences at the the University of Colorado Boulder (2021).*



**Kalind Carpenter** *Kalind is a Robotics Engineer in the robotic Advanced Robotic Systems group at JPL. He focuses on robotic access, manipulation and sampling technologies. His latest work includes the sample tube gripper lead for the Sample Recovery Helicopter. Kalind is the inventor and was the Principal Investigator of the Exobiology Extant Life Surveyor (EELS), an adaptable mobility capability aimed to traverse through the plume vent crevasses on Enceladus to reach the ocean below the ice. Previous work includes co-inventor of PUFFER and principle investigator (PI) of an autonomous swarm of robots with ground penetrating radar to map the changes in sub glacial topology in Greenland and Antarctica created by ocean currents that directly. This work was the precursor to the current CADRE mission.*



**Kai Pak** *is a Data Scientist in the Artificial Intelligence & Advanced Analytics Group at NASA's Jet Propulsion Laboratory (JPL) and co-lead for the Lab's exploration and experimentation with Large Language Models for scientific, engineering, and institutional applications. His research focuses on developing computer vision capabilities for satellite and other remote sensing imagery for Earth and planetary science applications. He holds an M.S. in Applied Physics from Columbia University and an M.S. in Computer Science from the University of Illinois Urbana-Champaign (UIUC).*



**Amanda Towler** *is the Group Supervisor for the Future Technology Exploration and Infusion Group at JPL, where she co-leads the lab's exploration of Large Language Models and GenAI capabilities. Prior to joining JPL, Amanda was an Associate Researcher with the Applied Research Laboratory for Intelligence & Security (ARLIS). Previously, she was the Co-Founder and Principal Investigator of Hyperion Gray LLC, a woman-owned small business providing technology R&D, rapid prototyping, and independent security research services to the US Government.*



**Rohan Thakker** *is a Robotics Technologist in the Robotics and Mobility Group at NASA's Jet Propulsion Laboratory since 2017 after finishing his graduate studies in robotics at Carnegie Mellon University. His research focuses on autonomous decision-making under uncertainty in real-world environments across the entire robotics stack. He is the Autonomy Lead of the EELS Snake Robot team which developed and deployed the NEO autonomy stack for exploring unknown ice worlds at the Athabasca Glacier in Canada. He has led autonomy sub-system teams across three DARPA programs including team costar that finished first in DARPA Subterranean Challenge's Tunnel Circuit.*



**Shehryar Khattak** *is a Robotics Technologist in the Perception Systems Group at NASA's Jet Propulsion Laboratory (JPL). His work focuses on enabling resilient robot autonomy in complex environments through multi-sensor information fusion. Currently, he is the Principal Investigator (PI) for the Multi-robot Autonomous Intelligent Search and Rescue task at JPL. He previously served as the perception lead for JPL's team in the DARPA RACER project and for Team CERBERUS, winners of the DARPA Subterranean Challenge. Before joining JPL, Shehryar was a postdoctoral researcher at ETH Zurich. He earned his Ph.D. (2019) and M.S. (2017) in Computer Science from the University of Nevada, Reno Additionally, he holds an M.S. in Aerospace Engineering from KAIST (2012) and a B.S. in Mechanical Engineering from GIKI (2009).*