

# SUFIA: Language-Guided Augmented Dexterity for Robotic Surgical Assistants

Masoud Moghani<sup>1</sup>, Lars Doorenbos<sup>2</sup>, William Chung-Ho Panitch<sup>3</sup>, Sean Huver<sup>4</sup>, Mahdi Azizian<sup>4</sup>, Ken Goldberg<sup>3</sup>, Animesh Garg<sup>1,4,5</sup>

**Abstract**—In this work, we present SUFIA, the first framework for natural language-guided augmented dexterity for robotic surgical assistants. SUFIA incorporates the strong reasoning capabilities of large language models (LLMs) with perception modules to implement high-level planning and low-level control of a robot for surgical sub-task execution. This enables a learning-free approach to surgical augmented dexterity without any in-context examples or motion primitives. SUFIA uses a human-in-the-loop paradigm by restoring control to the surgeon in the case of insufficient information, mitigating unexpected errors for mission-critical tasks. We evaluate SUFIA on four surgical sub-tasks in a simulation environment and two sub-tasks on a physical surgical robotic platform in the lab, demonstrating its ability to perform common surgical sub-tasks through supervised autonomous operation under challenging physical and workspace conditions.

Project website: [orbit-surgical.github.io/sufia](https://orbit-surgical.github.io/sufia)

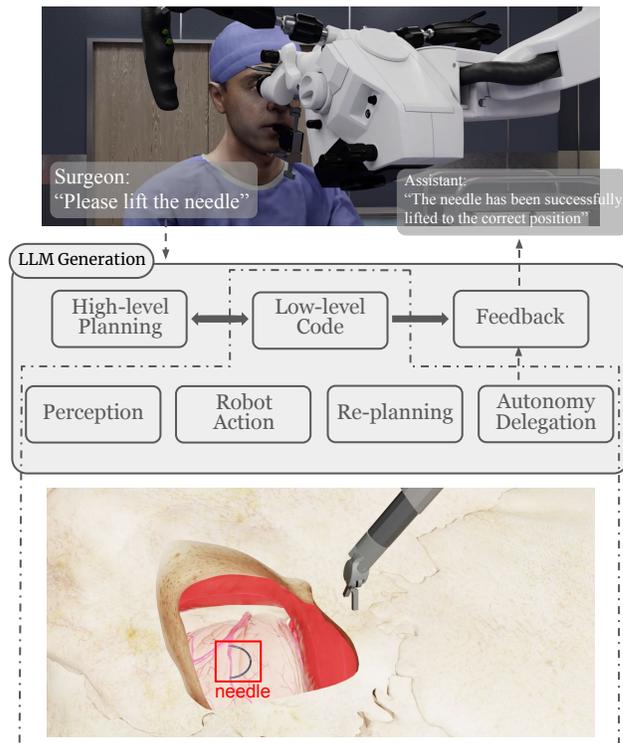
## I. INTRODUCTION

Recently, one prominent trend in surgery has been the increasing adoption of robotic surgical assistants (RSAs) in operating rooms. These RSAs are often controlled via local or remote teleoperation through a console by a trained human surgeon using hand controllers or other input peripherals, thereby enabling the surgeon to perform tasks with enhanced precision, dexterity, and control during an operation [1]. The teleoperated surgical procedures often involve tedious, repetitive, or time-consuming sub-tasks. Augmented dexterity in surgery holds the potential to simplify the surgical workflow, reduce surgeon fatigue, and improve patient outcomes [2], [3].

Learning-based approaches such as reinforcement and imitation learning learn policies to solve specific surgical sub-tasks [4], [5]. However, complex, long-horizon surgical sub-tasks are often computationally expensive, require extensive domain knowledge and reward engineering, and involve time-consuming dataset curation. Furthermore, the lack of generalizability limits the utility of learning-based models in safety-critical applications where unseen, in-domain variations are prevalent. As a result, most surgical robotic platforms still lack any level of autonomous capabilities [6].

In recent years, Large Language Models (LLMs) have received considerable attention for their ability to respond naturally to textual prompts and have been integrated into various domains, including the field of robotics and autonomous agents [7]. Language and vision models have demonstrated considerable promise in long-horizon robot planning and

<sup>1</sup>University of Toronto, <sup>2</sup>University of Bern, <sup>3</sup>University of California, Berkeley, <sup>4</sup>NVIDIA, <sup>5</sup>Georgia Institute of Technology  
[moghani@cs.toronto.edu](mailto:moghani@cs.toronto.edu), [animesh.garg@gatech.edu](mailto:animesh.garg@gatech.edu)



**Fig. 1: An overview of SUFIA automating the lifting of a suture needle from a surgical site.** SUFIA receives commands from a surgeon in natural language and converts them to high-level planning and low-level control code. If a task requires object interaction, SUFIA queries a perception module for object state information and generates low-level trajectories and robot actions accordingly. SUFIA can assist a surgeon with open-ended tasks, such as moving the robot in a desired motion to help complete a surgical task. In times of inefficient information, SUFIA delegates full control back to the surgeon.

control [8], [9], [10]. While these efforts still require pre-trained skills and motion primitives, they have demonstrated the potential of unified many-modality models for addressing a variety of complex tasks involving improved generalization to novel objects and unseen tasks.

In surgical settings, LLMs have the additional potential to aid interaction between a human surgeon and a robot via natural language teleoperation. This empowers the surgeon with the ability to use both fine-grained manual control and autonomous natural language conversational control in commanding the RSA to perform a sub-task. This approach promises both more natural human-robot coordination and the potential for developing general-purpose models for autonomous surgery beyond the capability of current task-by-

task automation approaches.

In this work, we present SUFIA (Surgical First Interactive Autonomy Assistant), a framework for natural interaction between a human surgeon and a surgical robot to provide interactive surgical autonomy. As shown in Fig. 1, SUFIA takes in sub-task commands from a surgeon and outputs a high-level natural language task plan, as well as low-level Python code snippets for execution, if requested. A perception module grounds perceived surgical objects in the scene regardless of variations in shape, size, and pose and accounts for the characteristics of their often small, slender shapes. SUFIA also incorporates re-planning and human-in-the-loop control as safety measures. Our primary contributions are as follows:

- A general formulation for natural language interaction between a surgeon and a robot.
- A language-based control approach to facilitate surgical sub-task implementations.
- A systematic evaluation of the generalization of our approach to various surgical sub-tasks, showing its performance and robustness for challenging workspace conditions.

## II. RELATED WORK

### A. Large Language Models for Robotics

Large Language Models (LLMs) are state-of-the-art natural language processing systems built on the transformer architecture [11]. LLMs are pre-trained with self-supervised objectives on vast amounts of text corpora, enabling these models to exhibit impressive language understanding and generation capabilities and perform a wide range of tasks. They are typically further fine-tuned with labeled data and RLHF to create general-purpose assistants [12], [13] or more specialized models for use cases such as coding [14], [15] or report generation [16].

In robotics, LLMs have been recently employed to address the high-level planning aspect of robotic control [17], [18]. These models still require trajectory generators through cost or reward functions to compute the trajectory. Other works focused on leveraging LLMs to design reward functions [19], [20] to acquire complex skills via reinforcement learning. However, most of these research works perform well on predefined tasks and still require expensive training time to generalize. Building on a recent work [21] that revealed the potential of LLMs to directly reason trajectory paths for robot arms, SUFIA incorporates LLMs to directly control the gripper poses to perform surgical sub-tasks. This enables the surgeon to naturally interact with the robot by asking for a complete task (e.g., pick the needle, insert the soft tube) or an open-ended task (e.g., move the needle in semi-circular motion) to help complete a sub-task. Our work differs from [21] in that we do not rely on a separate object detector for validation, incorporate further safety mechanisms by delegation, and show results for surgical scenes, where we additionally study domain-relevant axes such as variations in needle shape.

### B. Surgical Augmented Dexterity

Augmented dexterity has been attempted for several sub-tasks with varying levels of autonomy [22], [23], [6] such as dexterous needle picking and handling [24], [3], suturing [2], [25], and tissue manipulation [26], [27], [28]. In particular, in contrast to full automation, an often-explored paradigm in surgical robotics is augmented dexterity [29], in which minimal surgical sub-tasks are automated under human supervision, enabling more precise actuation with less effort expended.

However, these works largely rely on access to expensive, task-specific surgical hardware and software. In order to enable wider exploration of robotic automation in surgery, prior work has often focused on reducing the hardware barrier by adapting traditional robotic arm geometries for medical sub-tasks [30], [31], or designing novel, lower-cost, multi-purpose medical robotic systems [32], [33], [34]. Additionally, learning robust perception and control models for surgical tasks often requires gathering very large and expensive in-vivo datasets to avoid safety-critical failure cases [35], [36]. In this work, we propose an alternate approach to this software barrier by relying on a general-purpose, natural language-guided framework for surgical augmented dexterity across multiple tasks.

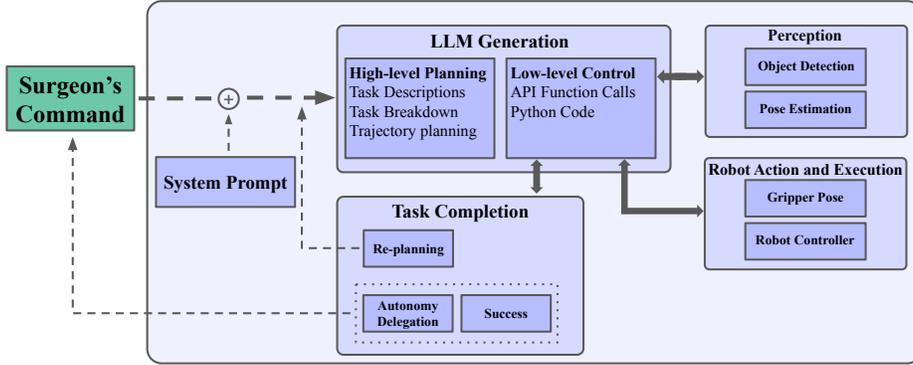
## III. PROBLEM FORMULATION

We focus on a novel approach to surgical augmented dexterity. In contrast to previous methods, we are investigating the potential of a generalist framework using large language models to address surgical augmented dexterity rather than training individual models for isolated tasks. We now briefly detail the assumptions with respect to the environment and available tools in our work. We do not provide any policies, trajectory optimizers, or in-context examples to the LLM [21]. Instead, we expect the LLM to reason over automating a benchmark simulated surgical sub-task with their internal knowledge and access to limited environment information through pre-defined function calls available in an API. All of our experiments are carried out on the da Vinci Research Kit (dVRK) robot platform [33]. The initial position and orientation of the dVRK grippers are available from the robot controller.

We assume access to a single RGB-D camera with a known intrinsic matrix, allowing for transformation between the camera’s perspective and the world coordinate space. With this, we design a perception module for the LLM to interact with and query object information. This module identifies and retrieves the pose information of objects present in the scene. For this, in simulation, we assume access to an instance segmentation model that, given an object name, outputs the segmentation maps of all instances of the queried object. In the physical experiments, we train a segmentation network based on the architecture from [37] for a needle segmentation model.

## IV. SUFIA

We propose SUFIA, a framework for natural interaction between surgeons and robots. SUFIA uses a human-in-the-



**Fig. 2: SUFIA architecture and workflows.** SUFIA enables a surgeon to naturally interact with the robot by either asking for a complete sub-task (e.g. “pick up the needle and hand it over to the other arm”) or generating a trajectory to help with performing a task (e.g. “move the needle 1 cm to the left”). SUFIA uses limited environmental knowledge in natural language (i.e. System Prompt) and scene understanding from a perception module to directly generate high-level plans and low-level sequences of gripper poses to interact with small-scale surgical objects. If SUFIA encounters difficulty in querying for an object or executing a necessary step to solve the surgical sub-task, it hands the control back to the surgeon for teleoperation.

loop approach, allowing either complete sub-task autonomy or assistance in open-ended tasks to help surgeons achieve their desired goals. The architectural framework, workflow, and primary elements of SUFIA are shown in Fig. 2. The following sections elaborate on the specifics of SUFIA.

#### A. LLM Generation and Planning

Crucial to the effectiveness of any LLM-based system is the design of the prompt, as only changing the prompt format can already lead to large differences in performance [38]. We build upon [21] who developed a single task agnostic prompt for performing low-level robot control for object grasping. We adapt it for surgical augmented dexterity with a four-part prompt, which consists of a role description and three core parts: the first part contains the API library available to the LLM; the second provides limited environment information (e.g. the state of the robot(s) to control and the orientation of the coordinate system); the third provides general instructions on how the LLM should generate the code, including the format of the desired output; the fourth describes prompt optimizations such as doing step-by-step reasoning [39].

#### B. API Library

The LLM has access to a library of functions that are available through an API. This API is documented in the main prompt, where for each function, its signature is given along with a brief description of its functionality [21], [40]. The API library mainly manages interaction with the robot control and perception modules. The modular approach with the API provides SUFIA the flexibility to adapt its respective modules independently, enabling integration into new embodiments and environments, such as switching from simulation to physical experiments.

Specifically, the API library includes robot control functionalities to execute a trajectory, rotate or open/close the gripper of the specified robot arm, and return the control back to the surgeon. Furthermore, perception functions detect the world poses of objects within the environment and can validate whether an object is at the expected position.

#### C. Perception

While LLMs lack the capability to ground physical worlds [41], they can still reason over the required steps to interact with objects and plan for task execution. To do so, we design a perception module that enables the processing of observations of the environment obtained from a single RGB-D camera to provide the object states to the LLM generator. This workflow is enabled by the API function *detect\_object*, through which the LLM queries and interacts with the perception module to retrieve object information.

*detect\_object* takes as input the name of the object to detect. After obtaining a segmentation of the named object and projecting it to world coordinates with the camera intrinsic matrix, we compute the 3D bounding cube and obtain the location and orientation. Moreover, for circular objects such as needles, we fit RANSAC to provide the object parameters and compute candidates for the location and orientation of a suitable interaction point.

#### D. Safety

A critical issue in surgical robotics is the reliability and safety of the robot control. To this end, we implement two components tailored to improve these aspects:

- 1) **Re-planning.** The original plan could become inappropriate due to, for example, mistakes in the planning or unforeseen circumstances, such as the gripper losing grip on the needle, in which case a new plan has to be devised. We encourage SUFIA to repeatedly use the *verify\_object* function to check whether the observed position of a given object matches the position expected by the framework. If the object being manipulated is not in the expected place, SUFIA re-plans the steps to complete the desired task given the updated knowledge of the environment.
- 2) **Human-in-the-loop approach.** In some cases, the perception module cannot find the desired object. In this case, rather than continuing blindly, SUFIA proceeds by handing control back to the surgeon for teleoperation with the API function *transfer\_control*. SUFIA is also instructed to call this function when it does not know how to solve a certain (sub-)task rather than operating on

insufficient information. Fig. 6 illustrates an instance when the system is unable to execute a command properly and returns the control to the surgeon to adjust the environment or provide further instruction.

Together, these two components enhance the safety and reliability of the assistant, which is crucial in the domain of surgical robotics. Note that the surgeon can also directly instruct SUFIA to re-take control of the robot, ensuring a smooth interplay between surgeon and robot.

## V. EXPERIMENTAL RESULTS

To empirically measure the efficacy of SUFIA, we perform experiments both in ORBIT-Surgical, a high-fidelity surgical simulation framework, and on a dVRK platform in the lab.

### A. Experimental Setup

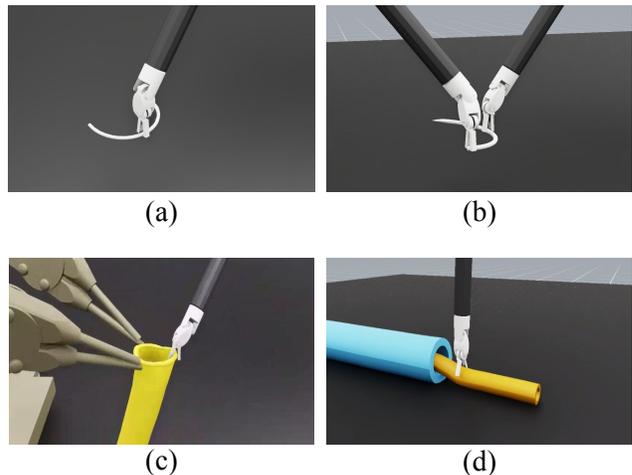
We conduct our simulation experiments in ORBIT-Surgical [42], which accurately imitates joint articulation and low-level controllers of the real dVRK platform, supports contact-rich physical interactions between rigid and deformable objects, and provides high-fidelity rendering. Furthermore, ORBIT-Surgical provides an interface for tele-operation, which enables the user to work together with SUFIA to solve a sub-task if needed. We use a camera sensor in NVIDIA Omniverse to acquire  $512 \times 512$  rendered RGB-D images and ground-truth semantic segmentation masks. In section V-E, we will discuss the utility of general segmentation models to adapt the workflow to objects with various configuration in simulation.

Physical experiments are performed on a da Vinci Research Kit (dVRK) [33] robot surgical assistant, using an Allied Vision Prosilica GC 1290 stereo camera pair for visual input. These cameras are capable of producing paired stereo frames at a resolution of  $1280 \times 960$  at 33 fps. Real-world depth images are then subsequently obtained by passing image pairs through RAFT-Stereo RVC [43], a state-of-the-art network for predicting image correspondences using optical flow, and then using the camera’s calculated intrinsic matrix to retrieve depth from these point discrepancies. We find that this approach provides better empirical results than traditional depth cameras in our use case due to the small, reflective objects and short focal lengths involved in the surgical setting. To emulate the real-world conditions encountered in a surgical setting, our workspace consists of a 3-D Med suturing tissue phantom on a red background. The phantom is then wrapped in blue cloth to imitate the use of a surgical cover during operation. These physical experiments introduce additional challenges, including a more challenging perception task, estimation and control noise, and more complex physics.

Throughout this section, we use GPT-4 Turbo [12] unless stated otherwise.

### B. Tasks and Evaluation Metrics

We demonstrate the generalizability of SUFIA by evaluating it across four distinct simulated surgical sub-tasks derived from ORBIT-Surgical, as shown in Fig. 3. We additionally select two of the subtasks (Needle Lift and Needle Handover) for evaluation using our physical setup. Each sub-task poses



**Fig. 3: Surgical sub-tasks.** (a) *Needle Lift*: lift a suture needle to a desired height, (b) *Needle handover*: pick and handover a suture needle, (c) *Vessel Dilation*: grip the vessel rim and dilate by pulling, (d) *Shunt Insertion*: insert a soft tube into larger vessel phantom. Best viewed in color.

unique challenges to show the robustness of the proposed workflow, described as follows:

*Needle Lift* – In this task, the needle (N1 in Fig. 5) is initialized at a random position and orientation within the reach area from a single dVRK arm. The task is successful if the robot grasps and lifts the needle to a specified height above the table.

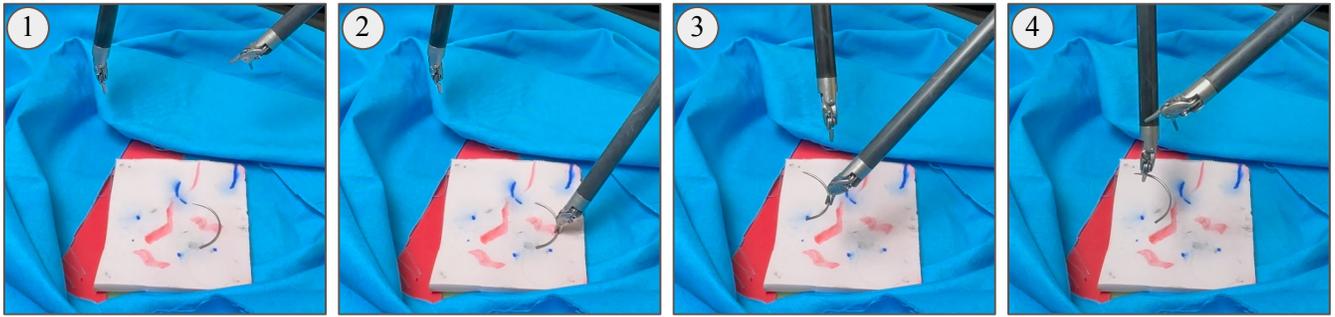
*Needle Handover* – This task involves transferring a needle using a dual-arm dVRK setup. The needle is initially positioned randomly. The arm closest to the needle first grasps and lifts it to a specified handover location. Subsequently, the second arm reaches for the needle, grasps it, and takes it to a desired position. The task is successful if the needle is effectively transferred from the initial to the second arm.

*Vessel Dilation* – In this task, a spring clamp assembly holds a soft vessel phantom from two points. The dVRK arm is required to grip the vessel rim from a third point facing the robot and dilate the vessel by pulling backward. A successful trial is defined if the robot fully dilates the vessel.

*Shunt Insertion* – This task requires using a dVRK gripper to insert a shunt into a vessel phantom. The arm grasps the shunt from the middle, lifts it slightly, and then inserts it into a vessel phantom. The task is considered successful if, upon release by the grippers, the shunt remains inside the vessel phantom.

### C. SUFIA Evaluation

We now discuss the effectiveness of SUFIA on solving the proposed surgical sub-tasks. SUFIA utilizes a perception module to localize the objects and proposes a sequence of sub-trajectories to perform the required task. We present the success rate for each sub-task for 10 trials in TABLE I. Overall, SUFIA is able to solve all proposed surgical sub-tasks requiring precise grasping of small surgical objects in simulation. Each task poses a unique challenge for automation, including object-gripper alignment and executing many steps to achieve successful results. In the *Vessel Dilation*



**Fig. 4: Physical Needle Handover task.** (1) Starting workspace configuration. The needle is placed in a fixed position within the workspace, and the gripper positions are randomized. In this stage, the SUFIA LLM planner queries for and identifies the pose of the suture needle, determines which robot arm is closest to it, and plans a trajectory for that robot arm to reach the suture needle. (2) The closest robot arm approaches and grasps the suture needle. (3) The suture needle is lifted to a neutral handover position. At this stage, the SUFIA LLM planner detects the pose of the suture needle at the handover position and plans a trajectory for the second robot arm to approach the needle. (4) The second robot arm descends and grasps the needle, then the first robot arm releases the needle after the second robot arm has grasped it. We provide task videos at [orbit-surgical.github.io/sufia](https://orbit-surgical.github.io/sufia)

Experiment	Success Rate	Planning Steps	Failure Modes	
			(P)	(E)
<b>Sim Experiments</b>				
Needle Lift	100 %	6	0	0
Needle Handover	90 %	14 - 16	1	0
Vessel Dilation	60 %	6 - 8	3	1
Shunt Insertion	70 %	8 - 9	3	0
<b>Physical Experiments</b>				
Needle Lift	100 %	6	0	0
Needle Handover	50 %	14 - 18	2	3

**TABLE I: Evaluation** Success rate and planning steps required for surgical sub-tasks automation (10 trials for each experiment). Failure modes: (P) denotes planning and (E) denotes execution failures. Sim experiments are carried out in ORBIT-Surgical, a high-fidelity surgical simulation framework. Physical Experiments are performed on a dVRK surgical platform.

task, all planning failures were due to not rotating the grippers to grasp the vessel’s rim correctly. In the *Shunt Insertion* task, the planning failures were from incorrect lift height calculations before insertion.

We observed that the performance of SUFIA was relatively robust to the more complex physics and observation spaces of the physical environment, with 0 and 2 planning failures encountered during the *Needle Lift* and *Needle Handover* experiments, respectively. This aligns closely with the framework’s performance in simulation. However, we found that hysteresis and encoder mismatch within the cable-driven dVRK resulted in variation between the commanded and actual gripper positions. Although SUFIA was often able to recover from the failures induced by this mismatch through its re-planning behavior, the lack of explicit servoing can result in dropping the needle during more complicated handovers.

#### D. Task Prompt Analysis

Simple tasks such as *Needle Lift* require a simple prompt to function properly. The surgeon can specify a position to transfer the needle to or allow the LLM to

determine a specific lift height above the table.

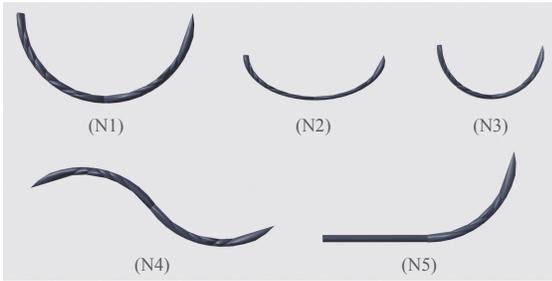
More sophisticated prompts are needed for tasks that require several steps for successful completion. In the *Needle Handover* task, the surgeon can provide additional notes for SUFIA to consider (e.g., "please note that for a handover, each robot should grasp the needle from the side closest to it."). The sequence in which the robot arms grasp and hand over to each other, as well as the location of the handover, can either be specified directly or left for the SUFIA to decide based on the distance to the needle or other environmental states.

The SUFIA planner may suggest unnecessary steps that may not be required for task completion and may potentially elongate task execution time. For instance, in the *Vessel Dilation* task, the vanilla prompt for dilating a vessel can sometimes lead to an additional step of "Lift the vessel slightly by moving the end-effector upwards to provide clearance from the table." The surgeon can provide additional information about the fact that the clamps are holding the vessel vertically to eliminate the suggestion of lifting steps in dilating the vessel. Similarly, in the *Shunt Insertion* task, additional information such as "please lift the small tube by a specific amount off of the table and horizontally insert it" helps to achieve better planning and execution.

Vision language models (VLMs) can also be incorporated in SUFIA to enhance the general visual understanding of the LLM planner. For instance, in the *Vessel Dilation* task, GPT4-Vision [12] can provide the planner with environmental context regarding the orientation of the vessel phantom. In this example, the VLM response can complement the user prompt: I see a vertical yellow tube on the right side of the image. It appears to be standing upright on one of its ends on a flat surface. While useful for providing general visual context, similar to [44], we find GPT4-Vision unreliable as a standalone perception module for detecting (small) objects’ spatial states and omit it for the remainder of our experiments.

The prompts used for the tasks are as follows:

*Needle Lift* – "Pick up the needle and lift it."



**Fig. 5: Needle variations in simulation.** We consider five instances of simulated suture needles (N1 - N5) with various sizes and shapes to conduct the generalizability experiment in ORBIT-Surgical.

*Needle Handover* – "Pick up the needle with the arm closest to it, move it directly to the handover location between the two arms, and keep holding the needle. Grasp the right side of the needle with the other robot arm, then right after that, release the needle from the first robot and stay put."

*Vessel Dilation* – "Grasp the vessel from its leftmost side with robot 0 and pull it backward to the left by 5 millimeters while holding on to it to dilate. When grasping the vessel, grasp it 15 millimeters below the left point."

*Shunt Insertion* – "Lift the small shunt from the middle and insert it into the left opening of the large tube. Approach the large tube from the left. Only lift the tube by 8 millimeters and move horizontally to insert."

#### E. Perception Adaptation and Domain Variation

While we envision domain-specific perception models for applications of our framework in a given surgical environment (e.g., a needle segmentation network in our real-world physical experimentation), here we investigate whether a general-purpose segmentation model, LangSAM [45], can be utilized by the perception module in simulated surgical environments to enhance the generalizability of the SUFIA framework to various object shapes.

We found that due to the slender shape of dVRK arms, the use of LangSAM in the perception module sometimes returned both the dVRK arm and the suture needle when prompted to find the "needle." However, descriptive adjectives (e.g. "round" or "small white needle") enable LangSAM to correctly identify and segment the suture needle well enough for SUFIA to generate the required steps to grasp and lift it. GPT-4 was also able to reason over the sizes of the segmented objects and properly determine the object of interest; here is an example of the LLM planner: "there are two objects detected as "white needle," but only one of them has dimensions that match a needle (Width: 0.011, Length: 0.032, Height: 0.002). The second object's dimensions are too large to be the needle we are interested in. Therefore, we will focus on the first object with the correct dimensions."

To study the generalizability of the perception module across various needle sizes and shapes, we conduct a study with five different needles, three needles in different sizes (N1

Perception Module	N1	N2	N3	N4	N5
Isaac Sim Camera	5 / 5	4 / 5	5 / 5	5 / 5	4 / 5
LangSAM	4 / 5	5 / 5	4 / 5	3 / 5	3 / 5

**TABLE II: Domain variation evaluation in simulation.** We report the success rates for lifting suture needles with varied sizes and irregular shapes (suture needles N1 - N5) over 5 trial runs with two variations of the perception module.

- N3; Fig. 5) and two irregular shapes (N4 and N5; Fig. 5). TABLE II shows the performance of SUFIA to lift various needles, which is robust to their shape and size.

#### F. LLMs Investigation

Here, we investigate the effect of different LLMs on the performance of the needle lift task. We use the same prompt and needle locations for all LLMs. As the error handling in SUFIA would, in principle, allow an LLM to keep trying endlessly until it generates code where no exceptions are raised, we limit the number of errors to five before terminating the program.

None of the open-source LLMs can perform the simple task of needle picking and have a hard time following the instructions in the prompt. All models struggle with understanding that `detect_object()` will print its result rather than return it as a variable in a Python script. When faced with errors, Mixtral [46] typically only outputs updated code snippets when asked to improve a code block rather than the whole code. CodeLlama [15] calls many undefined functions, such as `get_end_effector_pose()`, despite the end-effector pose being given in the prompt. Llama 2 [47] has a variety of mistakes related to understanding the steps in the task, such as forgetting to close the gripper or moving it down before lifting the needle.

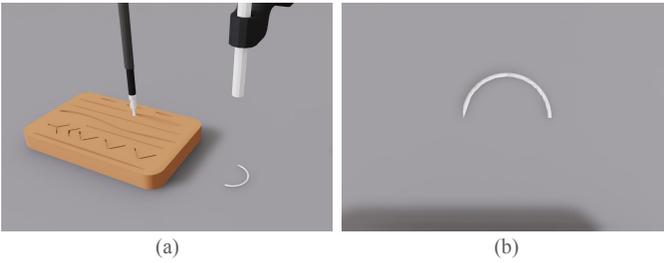
GPT3.5 Turbo similarly misunderstands `detect_object()`, often assigning its value to a variable called `needle_position`, despite the prompt stating the function does not return anything. Beyond that, GPT3.5 Turbo does consistently define a proper plan to lift the needle, but even when it calls `detect_object()` correctly, the information is not incorporated successfully.

All in all, in our experiments, only GPT-4 Turbo could follow all instructions and appropriately plan and execute the relatively simple task of lifting a suture needle.

#### G. Re-planning

To illustrate the benefits of our safety modules, we provide an example in the Needle Lift environment in Fig. 7. In the first row, SUFIA executes the plan it came up with to perform the task desired by the user, i.e., orienting its gripper with the needle, moving to a position where it can grab it, and picking it up. While picking it up, we move the needle to a different position. Because SUFIA validates the expected and observed position of the objects it manipulates, it correctly identifies the needle is not where it should be. Based on the newly observed state, SUFIA devises a new plan to proceed with the user instruction, finally lifting it to the desired height.

Surgeon: Please find the suturing site.



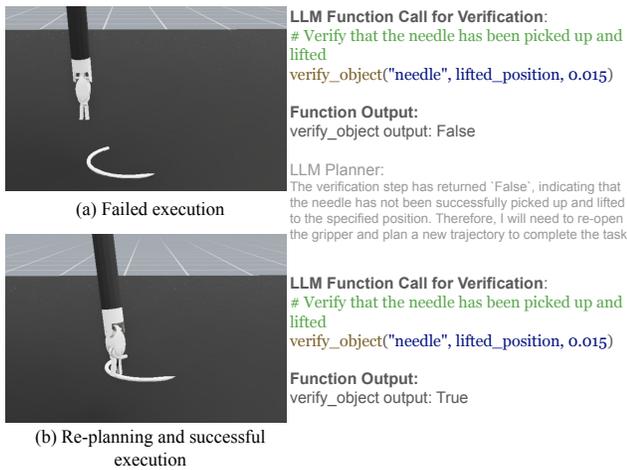
```
(c) SuFIA LLM Planner: To begin the task of finding the suturing site, I will first need to detect the location of the suturing site within the environment. I will call the 'detect_object' function with the parameter "suturing site" to obtain its position.
'''python
detect_object("suturing site")
'''

System execution output: Object not found!

SuFIA LLM Planner: Since the "suturing site" object was not found in the environment, I will defer control to the surgeon to provide further instructions or to adjust the environment so that the suturing site can be detected.
'''python
transfer_control()
'''

Control of the robot is transferred to the Surgeon.
```

**Fig. 6: Interactive human-in-the-loop approach.** (a) An overview of the environment showing the dVRK robotic arm and endoscope camera as well as a needle and a suturing pad in ORBIT-Surgical, (b) RGB image from the endoscope camera focused on the needle as input to the perception module, (c) delegation of control back to the surgeon due to the inability of the system to identify a suturing site. The uncertainty and occlusion in a surgical scene might lead to undesired object localization and retrieval. A human-in-the-loop approach enables a fail-safe paradigm for interactive robotic surgical assistants.



**Fig. 7: Re-planning.** (a) A failed execution as a result of not finding a proper grasping point on the needle. The LLM verification step indicates that the task was not executed successfully. (b) LLM planner proposes a new plan to re-identify and lift the needle by the robot.

### H. Limitations

The current best results are obtained with API calls to GPT-4 Turbo. Due to the generation speed of OpenAI’s API, SUFIA does not operate in real-time; for real-world needle handover trials, the API calls invoked during planning (including sensing and replanning) took an average of 25.5 seconds to complete, out of an average total trial time of 61.4 seconds. However, with open-source models constantly improving, running a quantized open-source LLM on-device will soon be a viable way to improve inference time greatly.

Moreover, while we incorporate two measures specifically designed to improve safety and reliability, deploying autonomous or semi-autonomous RSAs in real-world scenarios still has the potential to bring risks from unexpected circumstances an AI system might not be able to handle.

## VI. CONCLUSION

We present SUFIA, a modular framework for natural surgeon-robot interaction. We show that our training-free approach, which uses pre-trained LLMs to provide low-level

control of surgical robots, can successfully interact with small surgical objects and execute surgeon commands for automating surgical sub-tasks. Safety is bolstered through re-planning capabilities and a human-in-the-loop approach. We evaluate the efficacy of SUFIA for common surgical sub-tasks in simulated and physical experiments in the lab and show that the proposed method succeeds across different sub-tasks with various difficulty levels. These results suggest that language-guided autonomy has the potential to enhance surgeon’s efficiency in surgical procedures.

In future work, we plan to test the viability of quantized open-source LLMs on-device to improve inference time. This will also address any privacy concerns stemming from transmitting highly sensitive medical information to off-site servers. Furthermore, we intend to explore the usefulness of fine-tuned large language and vision models in SUFIA.

## REFERENCES

- [1] M. Hwang, D. Seita, B. Thananjeyan, J. Ichnowski, S. Paradis, D. Fer, T. Low, and K. Goldberg, “Applying depth-sensing to automated surgical manipulation with a da vinci robot,” in *2020 international symposium on medical robotics (ISMR)*. IEEE, 2020, pp. 22–29.
- [2] S. Sen, A. Garg, D. V. Gealy, S. McKinley, Y. Jen, and K. Goldberg, “Automating multi-throw multilateral surgical suturing with a mechanical needle guide and sequential convex optimization,” in *2016 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2016, pp. 4178–4185.
- [3] H. Lin, B. Li, X. Chu, Q. Dou, Y. Liu, and K. W. S. Au, “End-to-end learning of deep visuomotor policy for needle picking,” in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 8487–8494.
- [4] P. M. Scheikl, B. Gyenes, R. Younis, C. Haas, G. Neumann, M. Wagner, and F. Mathis-Ullrich, “LapGym - An open source framework for reinforcement learning in robot-assisted laparoscopic surgery,” *Journal of Machine Learning Research*, vol. 24, no. 368, pp. 1–42, 2023.
- [5] J. Xu, B. Li, B. Lu, Y.-H. Liu, Q. Dou, and P.-A. Heng, “Surrol: An open-source reinforcement learning centered and dvrk compatible platform for surgical robot learning,” in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 1821–1828.
- [6] A. Attanasio, B. Scaglioni, E. De Momi, P. Fiorini, and P. Valdastri, “Autonomy in surgical robotics,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 4, pp. 651–679, 2021.
- [7] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin *et al.*, “A survey on large language model based autonomous agents,” *arXiv preprint arXiv:2308.11432*, 2023.
- [8] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter,

- A. Wahid, J. Tompson, Q. Vuong, T. Yu *et al.*, “Palm-e: An embodied multimodal language model,” *arXiv preprint arXiv:2303.03378*, 2023.
- [9] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, “RT-1: Robotics transformer for real-world control at scale,” *arXiv preprint arXiv:2212.06817*, 2022.
- [10] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Chormanski, T. Ding, D. Driess, A. Dubey, C. Finn *et al.*, “RT-2: Vision-language-action models transfer web knowledge to robotic control,” *arXiv preprint arXiv:2307.15818*, 2023.
- [11] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [12] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, “GPT-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [13] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *Advances in neural information processing systems*, vol. 36, 2024.
- [14] Z. Luo, C. Xu, P. Zhao, Q. Sun, X. Geng, W. Hu, C. Tao, J. Ma, Q. Lin, and D. Jiang, “Wizardcoder: Empowering code large language models with evol-instruct,” *arXiv preprint arXiv:2306.08568*, 2023.
- [15] B. Roziere, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin *et al.*, “Code llama: Open foundation models for code,” *arXiv preprint arXiv:2308.12950*, 2023.
- [16] Z. Wang, L. Liu, L. Wang, and L. Zhou, “R2GenGPT: Radiology report generation with frozen LLMs,” *Meta-Radiology*, vol. 1, no. 3, p. 100033, 2023.
- [17] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng, “Code as policies: Language model programs for embodied control,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9493–9500.
- [18] S. Vemprala, R. Bonatti, A. Buckner, and A. Kapoor, “ChatGPT for robotics: Design principles and model abilities,” *Microsoft Auton. Syst. Robot. Res.*, vol. 2, p. 20, 2023.
- [19] T. Xie, S. Zhao, C. H. Wu, Y. Liu, Q. Luo, V. Zhong, Y. Yang, and T. Yu, “Text2Reward: Automated dense reward function generation for reinforcement learning,” *arXiv preprint arXiv:2309.11489*, 2023.
- [20] Y. J. Ma, W. Liang, G. Wang, D.-A. Huang, O. Bastani, D. Jayaraman, Y. Zhu, L. Fan, and A. Anandkumar, “Eureka: Human-level reward design via coding large language models,” *arXiv preprint arXiv:2310.12931*, 2023.
- [21] T. Kwon, N. Di Palo, and E. Johns, “Language models as zero-shot trajectory generators,” *arXiv preprint arXiv:2310.11604*, 2023.
- [22] T. D. Nagy and T. Haidegger, “Autonomous surgical robotics at task and subtask levels,” in *Advanced Robotics and Intelligent Automation in Manufacturing*. IGI global, 2020, pp. 296–319.
- [23] F. Ficuciello, G. Tamburrini, A. Arezzo, L. Villani, and B. Siciliano, “Autonomy in surgical robots and its meaningful human control,” *Paladyn, Journal of Behavioral Robotics*, vol. 10, no. 1, pp. 30–43, 2019.
- [24] A. Wilcox, J. Kerr, B. Thananjeyan, J. Ichnowski, M. Hwang, S. Paradis, D. Fer, and K. Goldberg, “Learning to localize, grasp, and hand over unmodified surgical needles,” in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 9637–9643.
- [25] S. Krishnan, A. Garg, S. Patil, C. Lea, G. Hager, P. Abbeel, and K. Goldberg, “Transition state clustering: Unsupervised surgical trajectory segmentation for robot learning,” *International Journal of Robotics Research (IJRR)*, 2017.
- [26] C. Shin, P. W. Ferguson, S. A. Pedram, J. Ma, E. P. Putson, and J. Rosen, “Autonomous tissue manipulation via surgical robot using learning based model predictive control,” in *2019 International conference on robotics and automation (ICRA)*. IEEE, 2019, pp. 3875–3881.
- [27] N. D. Nguyen, T. Nguyen, S. Nahavandi, A. Bhatti, and G. Guest, “Manipulating soft tissues by deep reinforcement learning for autonomous robotic surgery,” in *2019 IEEE International Systems Conference (SysCon)*. IEEE, 2019, pp. 1–7.
- [28] Y. Li, F. Richter, J. Lu, E. K. Funk, R. K. Orosco, J. Zhu, and M. C. Yip, “Super: A surgical perception framework for endoscopic tissue manipulation with surgical robotics,” *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 2294–2301, 2020.
- [29] K. Goldberg. (2023) Augmented Dexterity: How Robots Can Enhance Surgeon Dexterity. [Online]. Available: <https://bit.ly/Augmented-Dexterity-S24>
- [30] A. Shademan, R. S. Decker, J. D. Opfermann, S. Leonard, A. Krieger, and P. C. Kim, “Supervised autonomous robotic soft tissue surgery,” *Science translational medicine*, vol. 8, no. 337, pp. 337ra64–337ra64, 2016.
- [31] H. Saeidi, J. D. Opfermann, M. Kam, S. Wei, S. Léonard, M. H. Hsieh, J. U. Kang, and A. Krieger, “Autonomous robotic laparoscopic surgery for intestinal anastomosis,” *Science robotics*, vol. 7, no. 62, p. eabj2908, 2022.
- [32] B. Hannaford, J. Rosen, D. W. Friedman, H. King, P. Roan, L. Cheng, D. Glozman, J. Ma, S. N. Kosari, and L. White, “Raven-ii: an open platform for surgical robotics research,” *IEEE Transactions on Biomedical Engineering*, vol. 60, no. 4, pp. 954–959, 2012.
- [33] P. Kazanzides, Z. Chen, A. Deguet, G. S. Fischer, R. H. Taylor, and S. P. DiMaio, “An open-source research kit for the da vinci® surgical system,” in *2014 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2014, pp. 6434–6439.
- [34] T. Ranzani, G. Gerboni, M. Cianchetti, and A. Menciassi, “A bioinspired soft manipulator for minimally invasive surgery,” *Bioinspiration & biomimetics*, vol. 10, no. 3, p. 035008, 2015.
- [35] K. B. Ozyoruk, G. I. Gokceler, T. L. Bobrow, G. Coskun, K. Incetan, Y. Almalioğlu, F. Mahmood, E. Curto, L. Perdigoto, M. Oliveira *et al.*, “Endoslam dataset and an unsupervised monocular visual odometry and depth estimation approach for endoscopic videos,” *Medical image analysis*, vol. 71, p. 102058, 2021.
- [36] M. S. Yasar and H. Alemzadeh, “Real-time context-aware detection of unsafe events in robot-assisted surgery,” in *2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, 2020, pp. 385–397.
- [37] K. Dharmarajan, W. Panitch, B. Shi, H. Huang, L. Y. Chen, M. Moghani, Q. Yu, K. Hari, T. Low, D. Fer *et al.*, “Robot-assisted vascular shunt insertion with the dvrc surgical robot,” *Journal of Medical Robotics Research*, vol. 8, no. 03n04, p. 2340006, 2023.
- [38] Z. Zhao, E. Wallace, S. Feng, D. Klein, and S. Singh, “Calibrate before use: Improving few-shot performance of language models,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 697–12 706.
- [39] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large language models are zero-shot reasoners,” *Advances in neural information processing systems*, vol. 35, pp. 22 199–22 213, 2022.
- [40] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg, “Progprompt: program generation for situated robot task planning using large language models,” *Autonomous Robots*, pp. 1–14, 2023.
- [41] R. Liu, J. Wei, S. S. Gu, T.-Y. Wu, S. Vosoughi, C. Cui, D. Zhou, and A. M. Dai, “Mind’s eye: Grounded language model reasoning through simulation,” *arXiv preprint arXiv:2210.05359*, 2022.
- [42] Q. Yu, M. Moghani, K. Dharmarajan, V. Schorp, W. C.-H. Panitch, J. Liu, K. Hari, H. Huang, M. Mittal, K. Goldberg, and A. Garg, “ORBIT-Surgical: An open-simulation framework for learning surgical augmented dexterity,” *arXiv preprint arXiv:2404.16027*, 2024.
- [43] H. Jiang, R. Xu, and W. Jiang, “An improved raftstereo trained with a mixed dataset for the robust vision challenge 2022,” *arXiv preprint arXiv:2210.12785*, 2022.
- [44] M. Skreta, Z. Zhou, J. L. Yuan, K. Darvish, A. Aspuru-Guzik, and A. Garg, “Replan: Robotic replanning with perception and language models,” *arXiv preprint arXiv:2401.04157*, 2024.
- [45] L. Medeiros. (2023) LangSam: Language Segment-Anything. [Online]. Available: <https://github.com/luca-medeiros/lang-segment-anything>
- [46] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l. Casas, E. B. Hanna, F. Bressand *et al.*, “Mixtral of experts,” *arXiv preprint arXiv:2401.04088*, 2024.
- [47] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.