# An LLM-powered Natural-to-Robotic Language Translation Framework with Correctness Guarantees

1st ZhenDong Chen
*School of Computer Science and Engineering*
*Sun Yat-Sen University*
GuangZhou, China
chenzhd29@mail2.sysu.edu.cn

2nd ZhanShang Nie
*School of Computer Science and Engineering*
*Sun Yat-Sen University*
GuangZhou, China
niezhsh@mail2.sysu.edu.cn

3rd ShiXing Wan
*School of Computer Science and Engineering*
*Sun Yat-Sen University*
GuangZhou, China
wanshx3@mail2.sysu.edu.cn

4th JunYi Li
*School of Computer Science and Engineering*
*Sun Yat-Sen University*
GuangZhou, China
lijy727@mail2.sysu.edu.cn

5th YongTian Cheng
*School of Computer Science and Engineering*
*Sun Yat-Sen University*
GuangZhou, China
chengyt8@mail2.sysu.edu.cn

6th Shuai Zhao*
*School of Computer Science and Engineering*
*Sun Yat-Sen University*
GuangZhou, China
zhaosh56@mail.sysu.edu.cn

*Abstract*—The Large Language Models (LLM) are increasingly being deployed in robotics to generate robot control programs for specific user tasks, enabling embodied intelligence. Existing methods primarily focus on LLM training and prompt design that utilize LLMs to generate executable programs directly from user tasks in natural language. However, due to the inconsistency of the LLMs and the high complexity of the tasks, such best-effort approaches often lead to tremendous programming errors in the generated code, which significantly undermines the effectiveness especially when the light-weight LLMs are applied. This paper introduces a natural-robotic language translation framework that (i) provides correctness verification for generated control programs and (ii) enhances the performance of LLMs in program generation via feedback-based fine-tuning for the programs. To achieve this, a Robot Skill Language (RSL) is proposed to abstract away from the intricate details of the control programs, bridging the natural language tasks with the underlying robot skills. Then, the RSL compiler and debugger are constructed to verify RSL programs generated by the LLM and provide error feedback to the LLM for refining the outputs until being verified by the compiler. This provides correctness guarantees for the LLM-generated programs before being offloaded to the robots for execution, significantly enhancing the effectiveness of LLM-powered robotic applications. Experiments demonstrate NRTrans outperforms the existing method under a range of LLMs and tasks, and achieves a high success rate for light-weight LLMs.

## I. INTRODUCTION

With the growing demand for intelligent and autonomous robots in a wide range of applications (*e.g.*, smart factories [1], healthcare [2], and householding [3]), conventional

robotic systems face significant challenges, which can only complete specific tasks in predefined scenarios. To bridge this gap, LLMs (*e.g.*, OpenAI's GPT [4], Meta's Llama [5], and Google's Gemma [6]) are deployed in robotics, leveraging their semantic comprehension and contextual reasoning to generate robotic control programs that fulfill the given tasks.

Most of the existing LLM-powered control program generation methods for robotics can be broadly categorized into the following three fundamental paradigms [7]–[9], as shown in Fig. 1. In Fig. 1(a), the pre-trained LLMs are fine-tuned as an end-to-end solution that directly generates low-level control programs for underlying hardware (*e.g.*, motor, sensor) of a robot given specific user tasks. This approach [7], [10] simplifies the workflow of the generation process; however, it requires both enormous datasets and computational resources to tune the LLMs, which are often unavailable in resource-constrained scenarios. To address this, prompt engineering is proposed that utilizes LLMs to transform user tasks into robot actions without model training. As shown in Fig. 1(b), the pre-trained LLMs with tailor-designed prompts are applied to decompose user tasks into a series of pre-defined sub-tasks (actions), which are mapped to specific low-level hardware control programs for execution. In Fig. 1(c), the LLMs are applied to generate the high-level programs (*e.g.*, in Python) for robotics, which are interpreted into the low-level programs for controlling the hardware.

As discussed above, the LLM fine-tuning approach suffers from notable limitations in applicability due to the need for massive computational resources [11]. Considering resource-
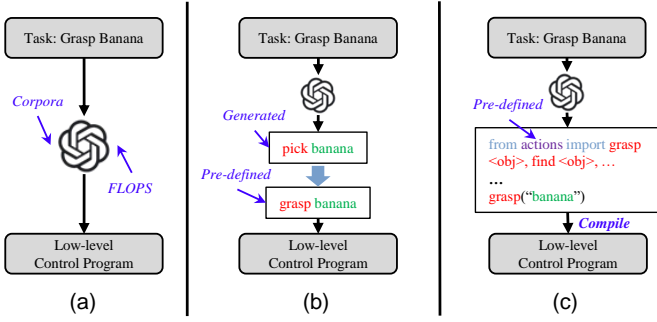
*Corresponding author: Shuai Zhao.

Fig. 1: Three paradigms of LLM-powered control program generation: (a) LLMs fine-tuned for low-level control program generation, (b) Task decomposition by LLMs into sub-tasks mapped to the pre-defined robot actions for low-level control programs, and (c) LLMs generate high-level programs compiled into the low-level control programs.

constrained scenarios, prompt engineering for control program generation is preferred due to the elimination of model training. However, with this approach, the generated control program is often prone to programming errors due to the inherent inconsistency of the LLMs and the high complexity of user tasks [12]. In addition, the complexity of tasks highlights the challenge of deploying light-weight LLMs on resource-constrained devices, as their performance often falls short of the required standards.

**Contributions.** To tackle these issues, this paper proposes NRTrans, an LLM-powered natural-to-robotic language translation framework that translates user tasks into executable robot control programs with correctness guarantees. To achieve this, we first develop a high-level Robot Skill Language (RSL) that reduces the complexity of program generation for LLMs, effectively bridging user tasks with robot skills by removing the robot control program details. This allows the LLM to directly generate RSL program without considering details of robot control programs. To verify the correctness of the generated RSL program, an RSL compiler is constructed that verifies the correctness for generated RSL programs, providing correctness guarantees for robot control programs compiled from RSL programs. In addition, to address the detected programming errors, a RSL debugger is developed that provides a feedback-based fine-tuning mechanism to refine the program. This effectively improves the success rate for generating compiler-verified control programs, especially for the light-weight LLMs on resource-constrained robots.

The experimental results demonstrate that NRTrans outperforms existing methods [13] by 53.6% on average in terms of the success rate of control program generation. In particular, NRTrans increases the success rate by 91.6% without model training by enabling the feedback-based program fine-tuning, which effectively detects and corrects errors in the generated RSL program before sending for execution. In addition, experiments demonstrate the effectiveness of NRTrans for light-weight LLMs, attaining a success rate of 92% with an LLM of 2B parameters [6].

## II. THE STATE-OF-ART AND LIMITATION

As trained on extensive corpora, pre-trained LLMs demonstrate remarkable capabilities in code generation for high-level languages (*e.g.*, C, Python). To accommodate LLMs to generate low-level control programs in robotic applications, LLMs, which are fine-tuned with datasets collected from specific robots and scenarios, have been emerging, including Google's RT-1 [7], Google's PaLM-E [10], and BIGAI's LEO [14]. The LLMs process user input, and directly generate low-level control programs to complete user tasks. However, this suffers from the following limitation due to the constrained resources.

> **Limitation 1: Resource Constraint**
>
> The fine-tuning for LLMs [7], [10], [14] lacks practicality due to their reliance on scenario-specific corpora and extensive computational resources [11].

In order to reduce the reliance on datasets and computation resources for applicability, a portion of existing work utilizes prompt engineering to enable LLMs in robotic applications. The work [15] validates that well-designed prompts enable LLMs to decompose high-level tasks into sub-tasks executable by robots. Based on this, existing methods [14]–[19] decompose user tasks into sub-tasks that are mapped to pre-defined robot actions for low-level control program generation. For instance, Language Planner [15] generates mid-level action plans in natural language, mapping them to pre-defined actions. This leads to the following limitation.

> **Limitation 2: Capability Requisition**
>
> Task decomposition [17]–[19] requires superior capabilities of LLMs to generate reasonable and coherent sub-tasks, making it challenging for the light-weight LLMs (*e.g.*, Gemma2-2b and Gemma2-9b) [15].

As discussed above, LLMs excel at code generation for high-level languages, which is leveraged by some existing work [8], [9], [13], [20]–[24]. The above methods utilize LLMs to receive user tasks and generate high-level language programs, which are then compiled into low-level robot control programs. ProgPrompt [13] uses a Pythonic representation to enhance the program generation of LLMs. This method relies on high-level language generation through LLMs and leads to the following limitation.

> **Limitation 3: LLMs Inconsistency**
>
> High-level language generation [13], [25], [26] are prone to programming errors due to the inconsistency of LLMs, significantly limiting the effectiveness due to the lack of correctness guarantees [12].

In addition, most existing methods are tailored for specific robots [7] and scenarios (*e.g.*, simulating household Virtual-Home [27]), or depend on LLMs with superior capabilities, while struggling with the inconsistency of LLMs [12]. These limitations make them challenging to apply in real-world robotic applications, in which light-weight LLMs on resource-constrained robots generate control programs that face the
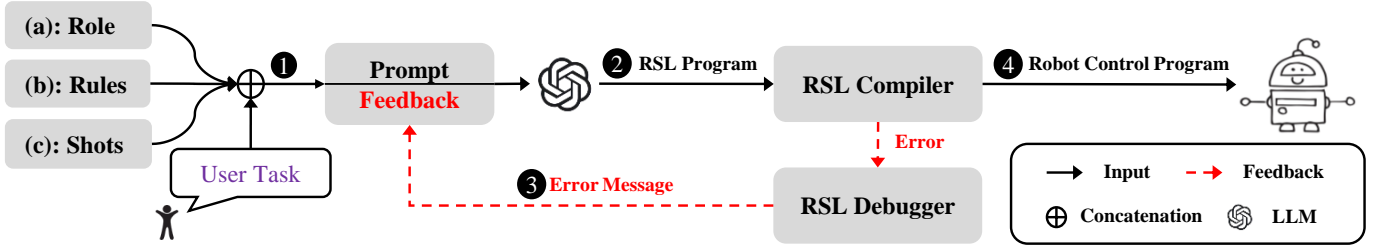
Fig. 2: The overall structure and workflow of NRTrans framework.

inconsistency of LLMs. To address this, we propose, NR-Trans, an LLM-powered natural-to-robotic language translation framework that offers correctness guarantees for programs and improves the success rate for light-weight LLMs.

## III. FRAMEWORK OVERVIEW AND REQUIREMENTS

To clarify how the NRTrans works, we provide an overview that translates user tasks to robot control programs, shown as four stages in Fig. 2. As the first work that provides correctness guarantees for robot control program generation, we focus on the design and construction of NRTrans while considering the following preliminaries and assumptions. This provides the groundwork for future extensions that include environmental interactions, real-world deployment, multi-robot collaboration, etc. Below, we present the requirements for NRTrans and discuss the solutions in Sec. IV.

### A. Preliminaries and Assumptions

To facilitate the construction, we establish the following preliminaries and assumptions for the NRTrans:

P1. We adopt a uniform prompt format sourced from OpenAI's API and extend to a range of LLMs with existing inference tools (*e.g.*, llama.cpp [28]).

P2. Following prior work [29], we assume that robots are built on the Robot Operating System (ROS) [30] and equipped with Python-based interfaces for robot control programming. These Python interfaces are encapsulated as discrete functions without self-decision-making.

P3. Advanced language mechanisms (*e.g.*, conditional and loop statements) and robotic capabilities (*e.g.*, dynamic environment monitoring) are deferred to future work, allowing us to prioritize the design and validation of the proposed translation framework for most tasks, thus providing a proof-of-concept for compiler-verified language translation by LLMs.

### B. Framework Overview

For robot control program generation, resolving the inconsistency of LLMs is critical to guarantee consistent outputs for underlying robots. To address this, we propose NRTrans, an LLM-powered framework for natural-to-robotic language translation with correctness guarantees. Considering the deployment of the NRTrans framework, we propose a feedback-based fine-tuning approach to minimize the requisitions for superior capabilities of LLMs. This method can achieve a high success rate when light-weight LLMs are applied. As illustrated in Fig. 2, NRTrans employs a structured four-stage phase to establish an entire workflow.

*Stage ❶: Prompt Construction and RSL Generation.* In our work, the prompt is initialized with the system message, shots (optional), and user task. The system message $\mathcal{M}$ defines the role of LLMs, expected generation specifications, and target language rules to the LLM, effectively guiding the generation of RSL programs utilized in Stage ❷, as shown in Fig. 2. The shots $\mathcal{S}$ provide output templates comprising selected user tasks and their corresponding RSL programs, constraining the LLM to generate outputs consistent with the desired format. Given a user task $\tau$ expressed in natural language, the prompt $\mathcal{P}$ is constructed by directly concatenating the system message, shots, and user task in sequence. The resulting prompt $\mathcal{P}$ serves as the query to the LLM, which generates the RSL program $\chi$ for the user task $\tau$ without human intervention.

*Stage ❷: RSL Compilation and Validation.* The compiler typically receives a high-level language program as input and then compiles it into a lower-level language, detecting programming errors if the program violates the language rules. For the generated RSL program, the RSL compiler translates the RSL program $\chi$ into a robot control program $\mathcal{R}$ and identifies hidden errors $\mathcal{E}$ within the program. The compiler ensures the correctness of the RSL program and compiles it into an executable program for the underlying robot in Stage ❹, thereby providing compiler-verified correctness guarantees for the generated robot control programs. However, if the RSL compiler detects errors $\mathcal{E}$, the NRTrans framework enters a feedback-based fine-tuning iteration for RSL program debugging in Stage ❸.

*Stage ❸: Feedback Composition and RSL Fine-Tuning.* The debugger is critical in identifying and resolving errors within RSL programs, improving the success rate of program generation without model training. If errors $\mathcal{E}$ are detected in Stage ❷, the error messages (e.g., a missing semicolon) generated by the RSL debugger are incorporated into the prompt $\mathcal{P}$ to form the feedback $\mathcal{F}$ for LLM re-evaluation. As shown in Fig. 2, the feedback $\mathcal{F}$ is utilized by LLMs to fine-tune the RSL program $\chi$ for solving programming errors. In our work, Stage ❷ and Stage ❸ operate in a closed fine-tuning loop to improve the success rate of program generation.

*Stage ❹: Robot Control Program Execution.* If the RSL program $\chi$ successfully passes correctness verification in Stage ❷ or exits the fine-tuning iteration between Stage ❷ and Stage ❸, the compiled robot control program is deployed to the

underlying robot for task completion. The control program is automatically translated into the low-level control program to control robot hardware (*e.g.*, motor, sensor), enabling the robot to perform the corresponding actions.

### C. Requirements for NRTrans

Compared with NRTrans, existing work utilizes formatted error messages with intricate function invocation from compilers aimed to fix high-level language errors [25], [26] and employed feedback for robot planning to complete tasks [31], facing runtime risks during execution. To the best of our knowledge, our work is the first to achieve a high success rate of task completion with light-weight LLMs, addressing LLM inconsistencies and providing correctness guarantees for practical applications. In our work, we independently designed and implemented the NRTrans framework. The requirements for building NRTrans are summarized as follows:

R1. To decompose user tasks into robot-executable skills, RSL must represent the capabilities of robots, offering intuitive and actionable representations for LLMs.

R2. To enable feedback-based fine-tuning with acceptable response times, the RSL debugger should distill program errors into concise, intuitive messages, omitting unnecessary details and providing clear guidance for LLMs.

R3. To support diverse robots and scenarios, NRTrans must quickly adapt to new capabilities and environments, ensuring compatibility and scalability of NRTrans.

```
1   # Role for Regex
2   "Generate regular expressions for user tasks with the following rules.\n"
3   # Regex Rules
4   "Anchors (`\^`, `\$`) match the entire input.",
    "Character classes (`[a-zA-Z0-9]`) define valid ranges.",
    "Quantifiers (`*`, `+`, `\{n,m\}`) control repetition.",
    "Escaped characters (`.` or `*`) match literally."
5   # Shots
6   "Task: A capital letter with 2-4 digits, ending with a period: \n"
    "```
        Regex: ^[A-Z][0-9]{2,4}\.$
    ```"
```

Fig. 3: A prompt for Regex Generation.

## IV. RATIONALE AND DETAIL FOR NRTRANS

As discussed in Sec. III-C, the design and implementation of the NRTrans need to satisfy requirements that affect the effectiveness and efficiency of translation from user tasks to robot control programs. In this section, we present design rationales and implementation details about each component of NRTrans to illustrate how to meet the requirements, providing thorough insights to support further research and development. In addition, characteristics of NRTrans are discussed to address limitations mentioned in Sec. II.

### A. Prompt Setting

As described in Sec. III-B, prompts are constructed with the system message and (optional) shots, which provide task-related foundational information and generation templates for

LLMs. To demonstrate the system message and shots, we use **Regular Expression** as a concrete example to illustrate how LLMs can translate user tasks in Fig. 3.

As demonstrated above, shots provide representative examples to enhance the performance of LLMs for subsequent user tasks. In our work, one shot is provided per RSL keyword; however, shots are not mandated in NRTrans, as error feedback can fine-tune generated programs. Experiments with zero-shot prompts in Sec. V-D validate this. As shown in Fig. 4, the system message, the optional shots, and the task are concatenated to form the prompt for LLMs.

```
1   # Role for RSL
2   "Break down a natural language task into subtasks and create a
    program with custom code and clear comments for each subtask.\n"
3   # RSL Rules
4   "FORWARD: 'forward'", "NUMBER: '-'?[0-9]+ ('.' [0-9]+)?; ", ...
    "forwardCommand: FORWARD NUMBER;", ...
5   # Shots
6   "move to position (0, 1): \n"
    "``` \n
        // navigate to definite position (0, 1) \n
        goto 0, 1; \n
    ```"
7   # User Task
8   "grasp the bottle on the table: \n"
```

Fig. 4: An illustrative example of the input prompt ((a) ⊕ (b) and optional (c) in Fig. 2).

### B. RSL Design and Compilation

The RSL language with its compiler forms the foundation for translation from user tasks to robot control programs with correctness guarantees in NRTrans. In this section, we illustrate the rationales of RSL design and how to implement RSL compiler. Guided by R1, RSL abstracts robot capabilities to achieve translation from RSL programs to robot control programs, while RSL compiler provides correctness guarantees within RSL program compilation.

*RSL.* The keywords of RSL are derived from robot capabilities, where each keyword is associated with one or more interfaces. These keywords abstract away from the details of robot control programs, highlighting their semantics of robot skills for LLMs to facilitate RSL program generation. As shown in Tab. I, the keyword meaning is intuitive, designed to express the robot's functionalities in natural language. The additional elements in RSL, such as the identifiers, numbers, and comment style, follow the conventions in C programming.

> **Characteristic 1: RSL Simplicity**
>
> RSL encapsulates the intuitive semantics of robot skills, enabling LLMs to generate RSL programs for user tasks without model fine-tuning, addressing Limitation 1.

The syntax of RSL is command-based, where each RSL statement consists of a keyword and optional parameters, and ends with a semicolon. Tab. I outlines the syntax rules of RSL, showing the usage and parameters for each keyword. We enforce that the numeric parameters must be positive to ensure

TABLE I: RSL Lexical and Syntax Design

| Robot Capability | Description | Keyword | Syntax |
|---|---|---|---|
| Forward | Moves the robot forward by a specified distance. | FORWARD | FORWARD NUMBER; |
| Backward | Moves the robot backward by a specified distance. | BACKWARD | BACKWARD NUMBER; |
| Turn Left | Rotates the robot left by a specified angle. | TURNLEFT | TURNLEFT NUMBER; |
| Turn Right | Rotates the robot right by a specified angle. | TURNRIGHT | TURNRIGHT NUMBER; |
| Look Up | Tilts the camera or sensor upward by a specified angle. | LOOKUP | LOOKUP NUMBER; |
| Look Down | Tilts the camera or sensor downward by a specified angle. | LOOKDOWN | LOOKDOWN NUMBER; |
| Look Left | Rotates the camera or sensor left by a specified angle. | LOOKLEFT | LOOKLEFT NUMBER; |
| Look Right | Rotates the camera or sensor right by a specified angle. | LOOKRIGHT | LOOKRIGHT NUMBER; |
| Perceive | Builds a SLAM map of the environment. | PERCEIVE | PERCEIVE; |
| Approach | Identifies and navigates to a specified object. | APPROACH | APPROACH OBJECT; |
| GoTo | Navigates to a given coordinate. | GOTO | GOTO NUMBER, NUMBER; |
| Grasp | Picks up a specified object with the robotic arm. | GRASP | GRASP OBJECT; |
| Others | Includes numbers, objects, and delimiters following C syntax rules. | - | - |

correct robot actions (*e.g.*, the forward instruction would not cause the robot to move backward). Identifiers represent objects manipulated by the robot with actions represented by specific keywords.

In addition, the conciseness of the RSL reduces the generation length for LLMs compared with conventional high-level language, which is advantageous for resource-constrained devices where token generation may be limited to one per second. However, RSL programs generated by LLMs are mostly accompanied by illustrative content. Considering the scope of this paper, the solution for RSL program generation and detailed analysis are deferred to future work.
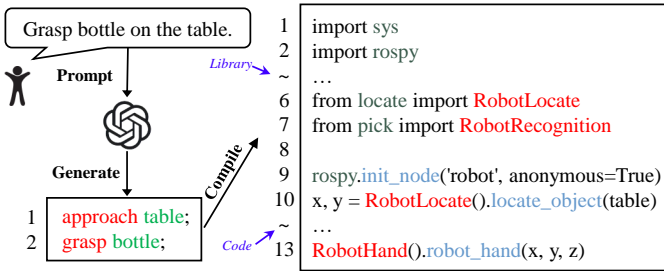


Fig. 5: A translation from a user task to a robotic program.

**RSL Compiler.** The RSL compiler verifies the syntax correctness of RSL programs, binds RSL rules to robot control interfaces, and generates executable robot control programs. In our work, lexical and syntactic rules of RSL are defined in regular expressions (see Tab. I), which are compatible with existing tools (*e.g.*, ANTLR [32]) that automatically generate the lexer and parser for RSL. The generated lexer and parser analyze RSL programs that are separated as individual tokens and parsed as an abstract syntax tree (AST) using LL(1) syntax analysis. A custom code generator, implemented with only a few lines of code (LOC), processes the AST in a depth-first manner, translating sub-trees (*i.e.*, RSL statements) into robot control programs. The automatic generation for lexer and parser, and the compact code generator support R3.

As illustrated in Fig. 5, the LLM generates the correspond-ing RSL program given a user task. The RSL program serves as input for RSL compiler, then is compiled to an executable robot control program. In our work, RSL compiler compiles RSL program into a Python program, automatically importing necessary libraries and sequentially invoking required func-tions to complete tasks.

> **Characteristic 2: Correctness Guarantees**
>
> The RSL compiler provides correctness verification for gen-erated RSL programs to solve the inconsistency of LLMs, effectively addressing Limitation 3.

### C. RSL Debugger

During the RSL compilation, programming errors are de-tected by RSL compiler for RSL programs, constructing associated error messages for LLMs by RSL debugger. Error messages enable LLMs to debug RSL programs, improving the success rate for RSL program generation.

**Feedback-based Tuning.** Based on the error messages, we propose a feedback-based fine-tuning method that iteratively refines generated RSL programs. If errors are spotted by RSL compiler, error messages will be produced by RSL debugger. As shown in Fig. 6, messages are concatenated with prompts to form error feedback for LLMs, fixing existing errors through regeneration of the RSL programs. The fine-tuning iteration finishes until the program passes the compiler verification. As illustrated in Sec. V-D, the feedback-based tuning effectively addresses the lack of semantic information for ambiguous tasks and zero-shot scenarios, improving the applicability of the NRTrans across different scenarios.

> **Characteristic 3: Capacity Enhancement**
>
> NRTrans employs a feedback-based fine-tuning method that enables light-weight LLMs to achieve a high success rate on resource-constrained devices, addressing Limitation 2.

**Error Messages.** Traditional error messages produced by compilers often contain sophisticated technical details, pre-venting an untrained LLM from understanding and fixing the errors. To overcome this, we design semantic-intuitive error

messages based on the RSL lexical and syntactic rules, with the line number and exact token highlighted. The lexical error messages are classified into five categories for every lexeme in RSL, including the keyword, identifier, number, character, and comment. Four types of syntactic error messages are constructed for syntactic errors in RSL. These error messages are presented in natural language with intuitive semantics for LLMs, as demonstrated in Tab. II. Based on the above intuitive error messages, NRTrans improves the success rate with feedback-based fine-tuning, satisfying the R2.

```
1   # Prompt
~   " ... "
6   " ... "
7   # User Task
8   "grasp the bottle on the table: \n"
9   # Error Message
10  "fixing and avoiding the compiler errors: 'keywords should be lowercase' \n"
```

Fig. 6: An example of the feedback for fine-tuning.

TABLE II: Error Message Design for RSL

| Type | Error Message | Example |
|------|---------------|---------|
| Keyword | Keywords should be lowercase. | APPROACH table; |
| Identifier | The identifier is illegal. | approach 3apple; |
| Number | The number is illegal. | forward 123.23.45; |
| Character | The $ is an illegal character. | $forward 1; |
| Comment | This comment has errors. | / This is a comment |
| Command | The command (keyword) is illegal. | move 1.5; |
| Parameter | Parameter types of the command are invalid. | forward table; |
| Quantity | The number of parameters is illegal. | goto 1; |
| Semicolon | The statement must end with a semicolon. | approach table |

### D. Robot Capsulation

In NRTrans, any requirements for decision-making or higher-level reasoning are delegated to the LLM, making robots receive definite parameters to execute specific actions. Based on definite robot skills, lexical and syntax rules can be easily constructed as keywords representing robot actions, with required parameters. As discussed above, rules of RSL in regular expressions and the code generator are implemented within hundreds of LOC, making NRTrans easily customized or extended to accommodate various robotic platforms.

## V. EVALUATION

To evaluate NRTrans, we compare it with ProgPrompt [13], which utilizes a Pythonic prompt for LLMs to generate robot control programs, in terms of the success rate, and accuracy of generated programs for NRTrans. Current task decomposition methods are not included due to the lack of open-source implementations and the incompatibility with the evaluation objectives for NRTrans [16], [17].

### A. Experimental Setup and Metrics

**Setup.** The NRTrans is deployed on an Intel i5-13400 CPU with an Nvidia RTX 4060Ti GPU. The Tiago [29] robot is

TABLE III: Examples of User Tasks for Evaluation.

| Group | Task |
|-------|------|
| **Simple Tasks** | Turn left 1 rad. |
| | Turn right 3 rads. |
| | Approach the door. |
| | Move forward 5 meters. |
| | Move backward 3 meters. |
| | Turn around a half circle. |
| **Ambiguous Tasks** | Go forward a little. |
| | Go forward further. |
| | Go forward a long distance. |
| | Turn around several circles. |
| **Multi-step Tasks** | Move forward 10 meters, then look around. |
| | Approach the workbench, then grasp a tool. |
| | Grasp the book, then move backward 2 meters. |
| | Move forward 5 meters, perceive, then turn right. |
| | Forward 2 meters, turn right 2 rads, and forward 1 meter. |
| | Move forward 2 meters, grasp the banana, turn left, then move forward 3 meters. |
| **Complex Tasks** | Grasp the bottle on the table and go back for me. |
| | Forward 4 meters, turn right, perceive environment. |
| | Repeat forward 1 meter and turn right 1.57 rads 4 times. |
| | Perceive the environment, approach the ball, and grasp it. |
| | Approach the table, look for a cup, and grasp it carefully. |
| | Approach the table, return to the origin in the coordinate. |
| | Move forward 3 meters, turn left, then move another 2 meters, then turn right. |
| | Repeat moving forward 1 meter and turning right 3 times for a triangular trajectory. |
| | Move forward 2 meters, turn right, look around; repeat this pattern until a full circle. |

applied for evaluation, equipped with ROS 1.0 and Python interfaces. The competing methods are evaluated under twenty-five user tasks, which are designed based on the skills of the Tiago robot. To assess the performance of NRTrans under different LLMs, Gemma2-2b [6], Gemma2-9b, Llama-70b [5], Gemini-1.5-Flash [33], and GPT-4o [34] are considered in the evaluation, covering a large scale of parameter size from 2 billion to 1.76 trillion (*i.e.*, GPT-4o) [35].

**Dataset.** Experimental tasks are categorized into four groups: (i) simple tasks, (ii) ambiguous tasks, (iii) multi-step tasks, and (iv) complex tasks, as shown in Tab. III. The tasks are designed with varying levels of complexity, demanding advanced capabilities from LLMs while addressing inconsistencies in their outputs. Using this task set, we aim to validate the effectiveness of NRTrans.

**Metric.** To quantify the performance of the NRTrans on the above tasks, three major metrics of the generated programs are considered: (i) **success rate**, verifying whether generated programs pass their compiler and conform to a uniform format that can be automatically executed; (ii) **accuracy**, assessing whether the program completes the requested task without the consideration of robot location and pose; and (iii) **pass**, representing the average number of iterations required for the LLM to generate the correct program, when feedback-based fine-tuning is enabled.
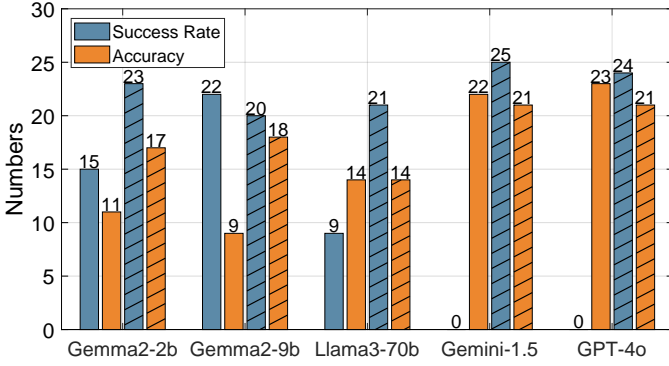
Fig. 7: Success Rate and Accuracy of Methods *(Bars without markers – ProgPrompt; Bars with markers - NRTrans;)*.

### B. Overall Performance Comparison under Varied LLMs

The NRTrans is evaluated against the ProgPrompt [13], which constructs a Python-structured prompt to import available objects and Tiago's interfaces for control program generation. The ProgPrompt can integrate with different LLMs, allowing us to evaluate and compare the effectiveness of both methods across a wide range of LLMs. Based on experimental tasks, Fig. 7 reports the number of executable programs and accurate programs generated by both methods under the considered LLMs, respectively.

As shown by the results, we first observe that the NRTrans constantly outperforms ProgPrompt, achieving an improvement of 53.6% in success rate and 9.6% in accuracy on the average case. Notably, NRTrans exceeded ProgPrompt by 30% in accuracy under two light-weight LLMs (*i.e.*, Gemma2-2b and Gemma2-9b). This observation demonstrates that the NRTrans effectively facilitates program generation of LLMs by abstracting away from the complex details of robot control programs, validating Characteristic 1. Interestingly, the ProgPrompt with the two largest and most powerful LLMs (*i.e.*, Gemini-1.5-Flash and GPT-4o) showed the lowest success rate but with high accuracy. This means that without effective constraints and guidance, such models often produce control programs in various forms (e.g., classes instead of functions), leading to execution failures without additional manual adjustments.

As for the NRTrans, it achieves at least 80% success rate and 56% accuracy across all LLMs. In general, the success rate of NRTrans is consistently higher than its accuracy for each LLM, translating the understanding of user tasks into correct control programs, and validating Characteristic 2. With Gemini-1.5-Flash and GPT-4o applied, NRTrans can produce correct programs in a single pass for nearly all user tasks. Notably, compared to GPT-4o, the light-weight LLMs (*i.e.*, Gemma2-2b and Gemma2-9b) demonstrate competitive performance on NRTrans by utilizing the feedback-based fine-tuning mechanism (*i.e.*, with more than one pass), validating Characteristic 3. This illustrates the effectiveness of the constructed RSL debugger, and more importantly, the applicability of NRTrans on resource-constrained devices.

### C. Benefits of Semantic-intuitive Error Messages

To validate the effectiveness of error messages (See Tab. II), an experiment is conducted to evaluate the semantic-intuitive error messages constructed in Sec. IV-C. The NRTrans with the original error messages produced by ANTLR [32] is applied as the baseline method. The Gemma2-2b is applied for both methods. As shown in Tab. IV, NRTrans outperforms the baseline in both success rate and accuracy while requiring fewer passes on average. This indicates that semantic-intuitive error messages can aid the LLM in understanding the type and the location of errors, providing effective guidance for the LLM during the fine-tuning process.

TABLE IV: Effectiveness of Customized Error Messages

| Method | Success Rate | Accuracy | Pass |
|---------|--------------|----------|------|
| Baseline | 22 / 25 | 16 / 25 | 1.56 |
| NRTrans | 23 / 25 | 17 / 25 | 1.4 |

### D. Performance of NRTrans with Zero Shot

The prompt often requires modifications to accommodate diverse scenarios. However, it can be challenging for users to provide specific shots that adapt robot control programs into an expected format (*e.g.*, a complete function replacing code statements). To investigate the performance of NRTrans in such situations, an experiment is conducted on NRTrans with zero shot under the GPT-4o, which possesses the necessary ability to comprehend the semantics of user tasks compared to other LLMs such as Gemma and Llama.

As shown in Tab. V, NRTrans with zero shot shows a slightly lower success rate and a significant drop in accuracy, which is expected due to the absence of detailed semantic information. However, compared to the success rate in one pass (*i.e.*, without feedback-based fine-tuning), NRTrans with zero shot significantly increases the success rate by 91.6%. This experiment again verifies the effectiveness of the design of NRTrans, and demonstrates its applicability for scenarios where no shots can be provided.

TABLE V: Performance of NRTrans with Zero Shot

| Method | Success Rate | Success Rate (One Pass) | Accuracy | Pass |
|--------|--------------|-------------------------|----------|------|
| NRTrans with Zero Shot | 23 / 25 | 12 / 25 | 17 / 25 | 2 |
| NRTrans | 24 / 25 | 24 / 25 | 21 / 25 | 1.16 |

### E. Effectiveness of Feedback-Based Tuning

As shown in Fig. 8(a), experiments across LLMs reveal that NRTrans generates robot control programs for user tasks within an average of two passes. To evaluate the effectiveness of feedback-based fine-tuning, we analyzed user tasks completed at each pass in a zero-shot setting. As depicted in Fig. 8(b), most tasks are solved within three passes, accounting for up to 88% of the experimental tasks. These findings validate the efficiency and effectiveness of feedback-based fine-tuning, highlighting its applicability to real-world scenarios.
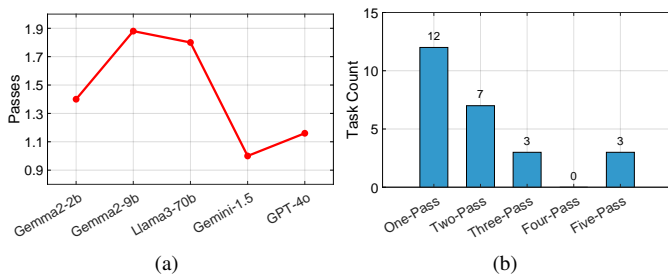
Fig. 8: (a) Average passes for feedback-based fine-tuning across LLMs, (b) Tasks completed per pass in zero-shot.

## VI. Conclusion

This paper presents NRTrans, an LLM-powered natural-to-robotic language translation framework. The framework introduces RSL with its compiler to facilitate robot control program generation and provide correctness guarantees for generated programs. Additionally, a RSL debugger incorporates a feedback-based fine-tuning method to improve the success rate of program generation, enabling light-weight LLMs on resource-constrained devices. The experiments validate the effectiveness of the proposed framework and highlight its potential for intuitive error messaging and feedback-based fine-tuning. In future work, we will (i) integrate environmental feedback for real-world interactions; (ii) extend RSL to support more complex user tasks, and (iii) improve the fine-tuning process for RSL program generation.

## VII. Acknowledgement

## References

[1] S. Harapanahalli, N. O. Mahony, G. V. Hernandez, S. Campbell *et al.*, "Autonomous navigation of mobile robots in factory environment," *Procedia Manufacturing*, vol. 38, pp. 1524–1531, 2019.

[2] M. Kyrarini, F. Lygerakis *et al.*, "A survey of robots in healthcare," *Technologies*, vol. 9, no. 1, p. 8, 2021.

[3] H. Sahin and L. Guvenc, "Household robotics: autonomous devices for vacuuming and lawn mowing [applications of control]," *IEEE Control Systems Magazine*, vol. 27, no. 2, pp. 20–96, 2007.

[4] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan *et al.*, "Language models are few-shot learners," 2020.

[5] H. Touvron, T. Lavril, G. Izacard *et al.*, "Llama: Open and efficient foundation language models," 2023. [Online]. Available: https://arxiv.org/abs/2302.13971

[6] G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S. Kale, J. Love *et al.*, "Gemma: Open models based on gemini research and technology," *arXiv preprint arXiv:2403.08295*, 2024.

[7] A. Brohan, N. Brown, J. Carbajal *et al.*, "Rt-1: Robotics transformer for real-world control at scale," *arXiv preprint arXiv:2212.06817*, 2022.

[8] K. Rana, J. Haviland, S. Garg, J. Abou-Chakra, I. Reid, and N. Suender-hauf, "Sayplan: Grounding large language models using 3d scene graphs for scalable task planning," *arXiv preprint arXiv:2307.06135*, 2023.

[9] 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.

[10] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch *et al.*, "Palm-e: An embodied multimodal language model," *arXiv preprint arXiv:2303.03378*, 2023.

[11] J. Yang, W. Tan, C. Jin, K. Yao, B. Liu, J. Fu *et al.*, "Transferring foundation models for generalizable robotic manipulation," 2024. [Online]. Available: https://arxiv.org/abs/2306.05716

[12] S. Vemprala, R. Bonatti, A. Bucker, and A. Kapoor, "Chatgpt for robotics: Design principles and model abilities," 2023. [Online]. Available: https://arxiv.org/abs/2306.17582

[13] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu *et al.*, "Progprompt: Generating situated robot task plans using large language models," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 11 523–11 530.

[14] J. Huang, S. Yong, X. Ma *et al.*, "An embodied generalist agent in 3d world," *arXiv preprint arXiv:2311.12871*, 2023.

[15] W. Huang, P. Abbeel *et al.*, "Language models as zero-shot planners: Extracting actionable knowledge for embodied agents," in *International Conference on Machine Learning*. PMLR, 2022, pp. 9118–9147.

[16] M. Ahn, A. Brohan, N. Brown, Y. Chebotar *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," *arXiv preprint arXiv:2204.01691*, 2022.

[17] R. Hazra, P. Z. Dos Martires, and L. De Raedt, "Saycanpay: Heuristic planning with large language models using learnable domain knowledge," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 18, 2024, pp. 20 123–20 133.

[18] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng *et al.*, "Inner monologue: Embodied reasoning through planning with language models," *arXiv preprint arXiv:2207.05608*, 2022.

[19] C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, and Y. Su, "Llm-planner: Few-shot grounded planning for embodied agents with large language models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2998–3009.

[20] D. Honerkamp, M. Buchner, F. Despinoy, T. Welschehold *et al.*, "Language-grounded dynamic scene graphs for interactive object search with mobile manipulation," *arXiv preprint arXiv:2403.08605*, 2024.

[21] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg, "Text2motion: From natural language instructions to feasible plans," *Autonomous Robots*, vol. 47, no. 8, pp. 1345–1365, 2023.

[22] D. Tanneberg, F. Ocker, S. Hasler, J. Deigmoeller, A. Belardinelli, C. Wang, H. Wersing, B. Sendhoff, and M. Gienger, "To help or not to help: Llm-based attentive support for human-robot group interactions," *arXiv preprint arXiv:2403.12533*, 2024.

[23] S. Wang, M. Han, Z. Jiao, Z. Zhang, Y. N. Wu, S.-C. Zhu, and H. Liu, "Llmˆ 3: Large language model-based task and motion planning with motion failure reasoning," *arXiv preprint arXiv:2403.11552*, 2024.

[24] J. Yang, Y. Dong, S. Liu, B. Li, Z. Wang, C. Jiang, H. Tan *et al.*, "Octopus: Embodied vision-language programmer from environmental feedback," *arXiv preprint arXiv:2310.08588*, 2023.

[25] M. Schäfer, S. Nadi, A. Eghbali, and F. Tip, "An empirical evaluation of using large language models for automated unit test generation," *IEEE Transactions on Software Engineering*, 2023.

[26] Z. Xu *et al.*, "Lecprompt: A prompt-based approach for logical error correction with codebert," *arXiv preprint arXiv:2410.08241*, 2024.

[27] X. Puig, K. Ra, M. Boben *et al.*, "Virtualhome: Simulating household activities via programs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8494–8502.

[28] G. Gerganov, "ggerganov/llama.cpp: Port of facebook's llama model in c/c++," https://github.com/ggerganov/llama.cpp, 2023.

[29] J. Pages, L. Marchionni, and F. Ferro, "Tiago: the modular robot that adapts to different research needs," in *International workshop on robot modularity, IROS*, vol. 290, 2016.

[30] M. Quigley, "Ros: an open-source robot operating system," in *IEEE International Conference on Robotics and Automation*, 2009. [Online]. Available: https://api.semanticscholar.org/CorpusID:6324125

[31] F. Joublin, A. Ceravola, P. Smirnov, F. Ocker, J. Deigmoeller, A. Belardinelli *et al.*, "Copal: corrective planning of robot actions with large language models," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 8664–8670.

[32] T. J. Parr and R. W. Quong, "Antlr: A predicated-ll (k) parser generator," *Software: Practice and Experience*, vol. 25, no. 7, pp. 789–810, 1995.

[33] G. Team, R. Anil, S. Borgeaud *et al.*, "Gemini: a family of highly capable multimodal models," *arXiv preprint arXiv:2312.11805*, 2023.

[34] J. Achiam, S. Adler, S. Agarwal, S. Ahmad, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[35] Wikipedia contributors, "Gpt-4 — Wikipedia, the free encyclopedia," 2024, [Online; accessed 20-September-2024]. [Online]. Available: https://en.wikipedia.org/w/index.php?title=GPT-4&oldid=1246495302