

# SKILLWRAPPER: Autonomously Learning Interpretable Skill Abstractions with Foundation Models

Ziyi Yang<sup>1</sup>, Benned Hedegaard<sup>1</sup>, Ahmed Jaafar<sup>1</sup>, Skye Thompson<sup>1</sup>,  
Yichen Wei<sup>1</sup>, Everest Yang<sup>1</sup>, Haotian Fu<sup>1</sup>, Shreyas Sundara Raman<sup>1</sup>,  
Stefanie Tellex<sup>1</sup>, George Konidaris<sup>1</sup>, David Paulius<sup>1</sup>, Naman Shah<sup>1</sup>

**Abstract**—We envision a future where robots are equipped “out of the box” with a library of general-purpose skills. To effectively compose these skills into long-horizon plans, a robot must understand each skill’s preconditions and effects in a form that supports symbolic reasoning. Such representations should be human-interpretable so that robots may understand human commands and humans may understand robot capabilities. Unfortunately, existing approaches to skill abstraction learning often require extensive data collection or human intervention, and typically yield uninterpretable representations. We present SKILLWRAPPER, the first known *active learning* approach that leverages foundation models to learn human-interpretable abstractions of black-box robot skills, producing representations that are both *probabilistically complete* and *suitable* for planning. Given only RGB image observations before and after skill execution, our system actively collects data, invents symbolic predicates, and constructs PDDL-style operators to model the skills. We present preliminary simulation results demonstrating that the abstract representations learned by SKILLWRAPPER can be used to solve previously unseen, long-horizon tasks.

## I. INTRODUCTION

In the near future, robots will be deployed from the factory to the real world, equipped with a set of general-purpose skills to interact with their environment. However, these skills could be black boxes that were learned, engineered, or obtained in unknown ways. As a result, the conditions under which skills can be used in real-world settings may be unavailable or environment-dependent, potentially leading to failure when sequencing skills to solve unseen tasks. Herein lie two important problems: First, without understanding the conditions under which each skill can be successfully executed (i.e., *preconditions*) and the likely outcomes of execution (i.e., *effects*), a robot may fail to identify task plans that effectively use its skills. Second, skill preconditions and effects should be *interpretable* to everyday users, as they would allow users to understand the robot’s decision-making process, making it easier to specify goals or task constraints.

In this paper, we present SKILLWRAPPER, the first known approach that uses foundation models to autonomously characterize robot skills, emphasizing *human-interpretable* state abstractions while guaranteeing *probabilistic completeness* and *suitability* for planning. Our approach assumes a skill-type signature as input and learns a PDDL-style [11] symbolic model for each skill. Previous work has extensively explored learning symbolic representations of high-level skills [8, 16, 15, 5]. However, these works either assume access to privileged



“There are three items *Vase*, *TissueBox*, and *Bowl*, and three locations *Sofa*, *CoffeeTable*, and *DiningTable*. Their initial positions are shown as follows. The robot is near the *Sofa* initially, and everything is placed stably, and all items can fit in every location. The goal is to have all items on the *Sofa*.”

## Output Plan

```
GoTo3(Sofa,CoffeeTable),  
PickUp5(Vase,CoffeeTable),  
GoTo2(CoffeeTable, Sofa),  
DropAt2(Vase,Sofa),  
GoTo4(Sofa,DiningTable),  
PickUp1(Bowl,DiningTable),  
GoTo2(DiningTable,Sofa),  
DropAt2(Bowl,Sofa)
```

**Fig. 1:** An example multi-modal task specification using natural language and egocentric visual observations, followed by the corresponding plan found by planning using the PDDL-style operators learned by SKILLWRAPPER.

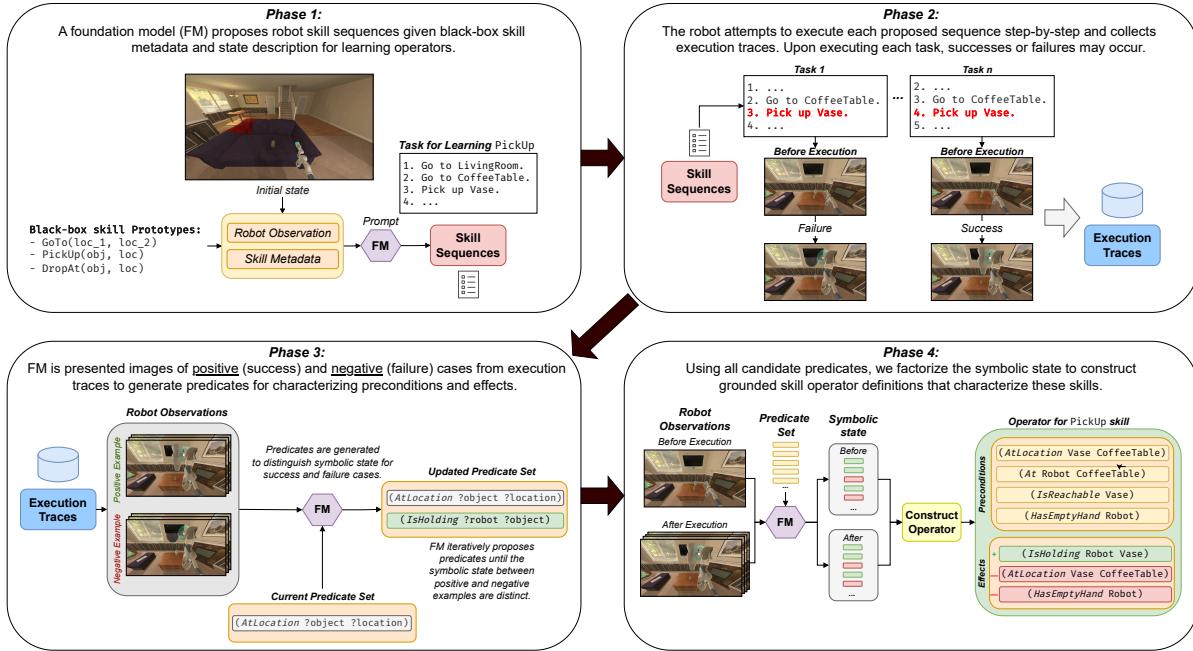
information (e.g., object poses [15] or extensive human feedback [5]) or fail to produce human-interpretable representations [8]. Although prior work has explored extracting symbols and language directly from demonstrations [4], existing approaches require significant manual effort to define features and train specialized classifiers for representation learning.

Hence, to facilitate learning skill abstractions from raw robot observations while reducing effort from human experts, we utilize foundation models, such as large language models (LLMs) and vision-language models (VLMs). Several works have exploited language models for robot decision-making and planning [1, 2, 9, 17, 14]. In contrast to these approaches, our method generates planning operators compatible by design with task planners, thus benefiting from efficient domain-independent heuristics [6] and correctness guarantees. This paper briefly introduces SKILLWRAPPER while highlighting key insights from our preliminary experiments in simulation.

## II. METHOD

Briefly, SKILLWRAPPER (Figure 2) learns an abstract model for planning with a library of black-box skills by (1) actively proposing and (2) executing exploratory skill sequences to collect execution traces, (3) inventing predicates by contrasting

<sup>1</sup>Brown University, Providence, RI, USA.



**Fig. 2:** Overview of SKILLWRAPPER: in (1), given a description of the robot’s environment and metadata about its skills, a foundation model (FM) proposes skill sequences useful for representation learning. In (2), the robot attempts to execute the proposed sequences, collecting the initial and final state of each action as images stored in a database. In (3), this database is presented to the FM as *contrastive pairs* (i.e., success and failure images as positive and negative examples) from which the FM will invent predicates to describe the symbolic state across all skills. Finally, in (4), the FM is used to infer the abstract states corresponding to the states before and after each successful execution trace. The resulting abstract transitions are used to construct planning operators.

pairs of successful and failed skill executions, and (4) using these predicates to generate PDDL-style operators compatible by design with off-the-shelf classical AI planners.

**Skill Sequence Proposal:** Our system first queries a foundation model to generate skill sequences intended to explore the symbolic state space in a directed manner.

**Predicate Invention:** SKILLWRAPPER uses foundation models to generate interpretable predicates, along with their semantic meanings as English sentences. Predicate invention aims to generate predicates that distinguish state features responsible for successful or unsuccessful skill executions.

**Operator Learning by Clustering:** Using the collected dataset of skill execution traces, SKILLWRAPPER evaluates the truth value of each predicate at every traced state, inducing a dataset of abstract state transitions. The operator learning algorithm identifies the effects and preconditions of the potentially multiple subgoal options corresponding to each skill [8].

**Planning with Learned Operators:** Having “wrapped” the black-box skills in corresponding learned operators, SKILLWRAPPER can solve task planning problems conveniently specified using natural language and images (Figure 1). To convert the multi-modal task specification into an initial PDDL state, the system queries a foundation model to classify whether each predicate holds given the current state description, in a way similar to existing work [10].

### III. EXPERIMENTS

We demonstrate the capabilities of the SKILLWRAPPER system using preliminary simulation experiments in the ManipulaThor [7, 3] environment. We prompt GPT-4o [12] with egocentric observations to evaluate the truth value of each

predicate. For predicate invention and skill sequence proposal, we utilize o1-preview [13]. We provide the simulated robot with three high-level actions: `PickUp(obj, loc)`, `DropAt(obj, loc)`, and `GoTo(loc1, loc2)`. These high-level actions function as black-box skills by executing a deterministic sequence of low-level motions.

In total, we collected data from five skill sequences consisting of 80 image-based states and 40 transitions, of which 24 skill executions were successful. Given the dataset of environment transitions and the abstract state space induced by the invented predicates, SKILLWRAPPER learned ten operators to model the three high-level skills. Once the system has learned the human-interpretable predicates and operators, we use a foundation model to formulate an unseen task planning problem. We present an example multi-modal task specification using natural language and images at the top of Figure 1. Given the learned abstract transition model and the task planning problem inferred from the above task specification, the task planner returned the plan shown at the bottom of Figure 1.

### IV. CONCLUSION

This paper formulates the problem of actively learning interpretable, abstract representations of black-box skills. We propose SKILLWRAPPER as a solution that exploits the multi-modal reasoning capabilities of foundation models. We demonstrate our approach using a proof-of-concept example in a simulated mobile manipulation setting. In ongoing and future work, we plan to compare our approach to alternative methods for active relational abstraction learning.

## ACKNOWLEDGEMENT

The work is supported by the ONR under the grant N00014-22-1-2592. We thank Everest Yang for his help in preparing in experimental setup.

## REFERENCES

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do As I Can, Not As I Say: Grounding Language in Robotic Affordances. In *Proceedings of The 6th Conference on Robot Learning*, volume 205 of *Proceedings of Machine Learning Research*, pages 287–318. PMLR, 14–18 Dec 2022.
- [2] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. PaLM-E: An Embodied Multimodal Language Model. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 8469–8488. PMLR, 23–29 Jul 2023.
- [3] Kiana Ehsani, Winson Han, Alvaro Herrasti, Eli VanderBilt, Luca Weihs, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. ManipulaTHOR: A Framework for Visual Object Manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4497–4506, 2021.
- [4] Nakul Gopalan, Eric Rosen, GD Konidaris, and Stefanie Tellex. Simultaneously Learning Transferable Symbols and Language Groundings from Perceptual Data for Instruction Following. In *Robotics: Science and Systems (RSS)*, 2020.
- [5] Muzhi Han, Yifeng Zhu, Song-Chun Zhu, Ying Nian Wu, and Yuke Zhu. InterPreT: Interactive Predicate Learning from Language Feedback for Generalizable Task Planning. In *Robotics: Science and Systems (RSS)*, 2024.
- [6] Malte Helmert. The Fast Downward Planning System. *Journal of Artificial Intelligence Research*, 26:191–246, 2006.
- [7] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv*, 2017.
- [8] George Konidaris, Leslie Pack Kaelbling, and Tomas Lozano-Pérez. From Skills to Symbols: Learning Symbolic Representations for Abstract High-Level Planning. *Journal of Artificial Intelligence Research*, 61:215–289, 2018.
- [9] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500, 2023.
- [10] Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. LLM+P: Empowering Large Language Models with Optimal Planning Proficiency. *arXiv preprint arXiv:2304.11477*, 2023.
- [11] D. McDermott, M. Ghallab, A. Howe, C. Knoblock, A. Ram, M. Veloso, D. Weld, and D. Wilkins. PDDL – The Planning Domain Definition Language. Technical report, CVC TR-98-003/DCS TR-1165, Yale Center for Computational Vision and Control, 1998.
- [12] OpenAI. Hello GPT-4o, 2024. URL <https://openai.com/index/hello-gpt-4o/>. Accessed: May 29, 2025.
- [13] OpenAI. Introducing OpenAI o1-preview, 2024. URL <https://openai.com/index/introducing-openai-o1-preview/>. Accessed: May 29, 2025.
- [14] Shreyas Sundara Raman, Vanya Cohen, Ifrah Idrees, Eric Rosen, Ray Mooney, Stefanie Tellex, and David Paulius. CAPE: Corrective Actions from Precondition Errors using Large Language Models. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 14070–14077, 2024.
- [15] Naman Shah, Jayesh Nagpal, Pulkit Verma, and Siddharth Srivastava. From Reals to Logic and Back: Inventing Symbolic Vocabularies, Actions and Models for Planning from Raw Data. *arXiv preprint arXiv:2402.11871*, 2024.
- [16] Tom Silver, Rohan Chitnis, Nishanth Kumar, Willie McClinton, Tomás Lozano-Pérez, Leslie Kaelbling, and Joshua B. Tenenbaum. Predicate Invention for Bilevel Planning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(10):12120–12129, Jun. 2023.
- [17] Jimmy Wu, Rika Antonova, Adam Kan, Marion Lepert, Andy Zeng, Shuran Song, Jeannette Bohg, Szymon Rusinkiewicz, and Thomas Funkhouser. TidyBot: personalized robot assistance with large language models. *Autonomous Robots*, 47(8):1087–1102, November 2023.