

# Leveraging Language for Accelerated Learning of Tool Manipulation

Allen Z. Ren<sup>1</sup>, Bharat Govil<sup>2</sup>, Tsung-Yen Yang<sup>2</sup>, Karthik Narasimhan<sup>2\*</sup>, Anirudha Majumdar<sup>1\*</sup>

<sup>1</sup>Department of Mechanical and Aerospace Engineering, <sup>2</sup>Department of Computer Science  
Princeton University

{allen.ren, bgovil, ty3, karthikn, ani.majumdar}@princeton.edu

**Abstract:** Robust and generalized tool manipulation requires an understanding of the properties and affordances of different tools. We investigate whether linguistic information about a tool (*e.g.*, its geometry, common uses) can help control policies adapt faster to new tools for a given task. We obtain diverse descriptions of various tools in natural language and use pre-trained language models to generate their feature representations. We then perform language-conditioned meta-learning to learn policies that can efficiently adapt to new tools given their corresponding text descriptions. Our results demonstrate that combining linguistic information and meta-learning significantly accelerates tool learning in several manipulation tasks including pushing, lifting, sweeping, and hammering.<sup>2</sup>

**Keywords:** Language for Robotics, Tool Manipulation, Meta-learning

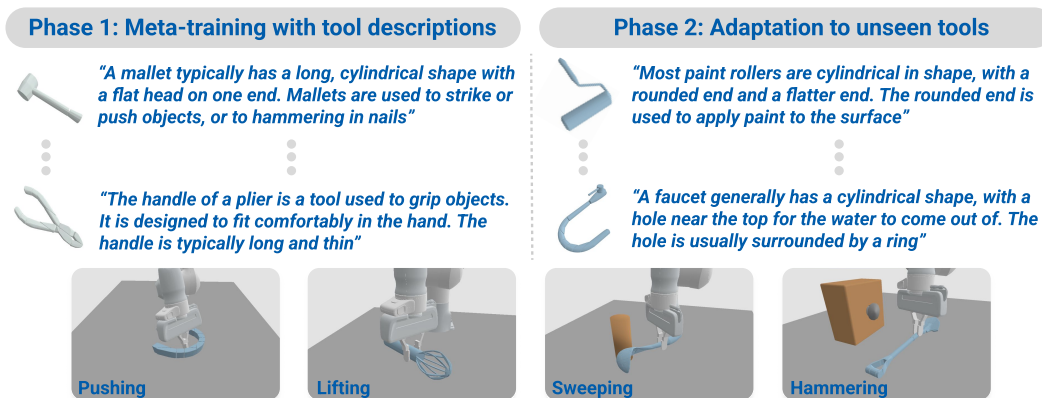


Figure 1: Rich, semantic knowledge from language descriptions, such as geometric features and common use of the tools, can help policies adapt faster to unseen tools (right) in pushing, lifting, sweeping, and hammering tasks (bottom) after meta-learning on training tools (left).

## 1 Introduction

The ability to quickly learn how to use a broad range of new tools is considered one of the hallmarks of human intelligence [1]. Humans are able to attain a degree of proficiency with a new tool (*e.g.*, a new hand tool or a virtual tool such as a joystick or remote) within just a few minutes of interaction [2]. This rapid learning relies on the ability to understand or discover the *affordances* [3] of a new tool, *i.e.*, the ability to perform a certain action with the tool in a given environment. For example, a hammer affords the opportunity to grasp it, use it to hammer a nail, or use it to push another object. In contrast, a spatula affords the opportunity to flip a pancake or sweep food ingredients into a bowl. Endowing robots with the ability to quickly discover and exploit affordances of a new tool in order to learn how to perform a given task has been a long-standing grand challenge in robotics [4].

\*Equal contribution in advising

<sup>2</sup>Video showing the four manipulation tasks: <https://shorturl.at/dmCEW>

While humans routinely rely on geometric priors and visual observations for tool manipulation, they can also ‘read’ text manuals or linguistic descriptions to understand affordances of new tools and quickly adapt to using them. In this work, we investigate whether *language* can help robots accelerate the process of learning to use a new tool. In particular, consider the following descriptions of two types of tools:

*The shape of tongs is typically that of a V, with two long, thin handles that come to a point at the top, and a gripping area in the middle.*

*A spatula is a kitchen utensil for flipping over food while cooking. The head of the spatula is usually rectangular or oval in shape. The handle of the spatula is usually long and thin.*

Our key intuition is that such natural language descriptions of tools contain information about the affordances of the tools, how to exploit these affordances for a given task, and how perceptual features of tools (*e.g.*, their visual appearance and geometry) relate to their affordances. Moreover, language can help capture the *shared structure* of tools and their affordances. Thus, if one has previously learned to use a set of tools (with corresponding language descriptions), a description of a new tool can help to exploit this prior knowledge in order to accelerate learning.

**Statement of Contributions.** Our primary contribution is to propose **ATLA** — Accelerated Learning of Tool Manipulation with **L**anguage — a meta-learning framework leveraging large language models (LLMs) to accelerate learning of tool manipulation skills. The overall approach is illustrated in Fig. 1. We propose to use LLMs in two distinct ways: to generate the language descriptions for tools and to obtain the corresponding feature representations. At meta-training time, the meta-learner updates a base-learner that quickly fine-tunes a manipulation policy; this fine-tuning process is conditioned on the LLM representations corresponding to the language descriptions of each tool. Specifically, we propose a simple gradient-based meta-learning setup based on Reptile [5] that performs off-policy updates. At test time, the base-learner adapts to a new tool using its language descriptions and interactions with it. To our knowledge, our approach is the first to utilize LLMs to accelerate learning of new tools. We demonstrate the benefits of using language in a diverse set of tool-use tasks including pushing, lifting, sweeping, and hammering.

## 2 Related Work

**Tool Manipulation.** Tool manipulation [6, 7, 8, 9, 10] is one of the long-standing problems in robotics research. A major challenge is understanding the affordances of the tool in different tasks. Previous work has modeled and learned affordances from parameterized keypoints on the tools [7, 10], from human demonstrations [8], and from spatial-temporal parsed graphs of the tools [11]. Our work instead leverages natural language (*e.g.*, describing affordances of the tools in words) for generalization of affordances in tool manipulation and is compatible with previous approaches.

**Language-informed Control.** Natural language has been applied to enable efficient robotic learning through (1) generating primitive language instructions for producing control actions (*i.e.*, instruction following task) [12, 13, 14, 15, 16, 17, 18, 19], (2) learning language-informed reward functions for training control policies [20, 21, 22, 23, 24, 25], and (3) using language to correct or adapt the behavior of the robot [26, 27]. However, these works primarily translate natural language into action policies for a specific task with the text providing information on the desired actions that optimize returns (*e.g.*, “*push the door*”). This means that the text is tightly coupled with the task seen during training, making it difficult to generalize to a new distribution of tasks with different dynamics. In contrast, the text in our work only provides a high-level description of the property of each tool, encouraging the agent to extract useful information to generalize to a new task. Some prior work [28, 29, 30, 31] also uses language descriptions of environment dynamics to enable generalization of policies but does not leverage meta-learning.

**Meta-learning.** Our work uses the framework of meta-learning [32, 33, 34, 35, 36], in which the agent is trained with a distribution of tasks, and later adapts quickly to a previously unseen task. “Reptile” is proposed in [5] as a simple first-order, gradient-based meta-learning algorithm that learns an initialization of the neural network’s parameters for fast adaptation at test time. Recent papers [37, 38, 39, 40] have also explored providing additional context information (*e.g.*, the property of the task) to encode task-specific knowledge for a meta-learning agent. However, all these works directly provide the context information either through scalar signals or a learned task embedding, which require domain expertise or a pre-training stage. In this work, we assume that the agent is provided with a text description of the tool, which is more accessible and easier to collect.

### 3 Problem Formulation

We consider the following goal: given a new tool and corresponding language description(s), we aim to *quickly* learn a policy for using the tool to achieve a given task. We pose this problem in a meta-learning setting in which a policy is trained with a distribution of tools for a given task, and later adapts quickly to a previously unseen tool sampled from this distribution in the same task.

**Meta-training.** During meta-training, we assume access to a set  $\mathcal{T} = \{\tau_i\}_{i=1}^K$  of tools, where  $K$  is the number of available tools. For each tool  $\tau_i$ , we are also provided a set of corresponding language descriptions  $L_i = \{l_{ij}\}_{j=1}^{N_i}$ ,  $l_{ij} \in \mathcal{L}$ , where  $N_i$  is the number of available descriptions for tool  $\tau_i$  and  $\mathcal{L}$  is the set of possible textual descriptions. In addition, each  $l_{ij}$  can describe a different property of tools such as shape and common use. Given a particular robotic manipulator and a particular task (e.g., pushing, lifting, sweeping, or hammering), each tool  $\tau$  induces a partially-observable Markov decision process (POMDP):  $\langle \mathcal{S}_\tau, \mathcal{A}, \mathcal{O}, \mathcal{P}_\tau, R_\tau \rangle$ . Here,  $\mathcal{S}_\tau$  is the state of the entire environment (i.e., combined state of the robot, tool, and potentially other objects to be manipulated using the tool). The robot’s action space  $\mathcal{A}$  (e.g., corresponding to robot joint torques) and observation space (e.g., the space of RGB-D observations from a camera) are fixed across tools. The transition probabilities are given by  $\mathcal{P}_\tau : \mathcal{S}_\tau \times \mathcal{A} \times \mathcal{S}_\tau \rightarrow [0, 1]$ , and the reward function is  $R : \mathcal{S}_\tau \times \mathcal{A} \times \mathcal{S}_\tau \rightarrow [0, 1]$ . During meta-training, our goal is to learn a policy  $\pi_\theta : \mathcal{O} \times \mathcal{L} \rightarrow \mathcal{A}$  parameterized by  $\theta$  (e.g., weights of a neural network) that can be quickly fine-tuned at test time.

**Meta-testing.** At test time, we are provided a new tool  $\tau_\nu$  and corresponding language descriptions  $L_\nu = \{l_{\nu j}\}_{j=1}^{N_\nu}$ ,  $l_{\nu j} \in \mathcal{L}$ . We aim to let the meta-learned policy quickly adapt to this new tool in a fixed number of interactions with the tool in order to maximize the expected cumulative reward. This is a challenging task since the new tool can be quite different in terms of visual appearance and affordances as compared to previously seen tools in meta-training.

### 4 Approach

The key idea behind our approach is to collect and embed language information of the environment into meta-learning, allowing the policy to adapt faster and better to unseen environments.

#### 4.1 Collecting Language Information Using Pre-Trained Large Language Models

A common use of language in robotics is to use it to provide an instruction to the robot (e.g., “pick up the green block on the table”). Such instructions are typically specified by humans manually through crowd-sourcing, which can be labor intensive. In our setting, we consider language as additional information about the environment (e.g., “the hammer has a long handle and large head at the top”). The language here is not used to describe the goal (e.g., what to do), but to provide information about properties of the environment (e.g., tool shape). This makes the text here *task-agnostic*, forcing the agent to learn generalizable policies. To obtain a diverse set of language descriptions, we are inspired by the recent advances in LLMs that are trained with vast amounts of online data and imbued with rich, semantic knowledge of different objects. We propose using LLMs to provide language descriptions of the tools in the form of *question answering*. Specifically, we provide the GPT-3 [41] model with the following template prompt through the OpenAI API:

“Describe the [feature] of [name] in a detailed and scientific response: ”

where “feature” is selected from one of [“shape”, “geometry”] or one of [“common use”, “purpose”] and “name” describes the tool (e.g., “a hammer”, “a pair of tongs”). We find that adding “*detailed and scientific*” to the prompt significantly improves the quality of the texts generated. For each tool, we generate 10 different paragraphs of descriptions for each of the four features, and then combine paragraphs in each of the four permutations of the features (“shape” and “common use”, etc). Each tool  $\tau_i$  is thus paired with a diverse set of 800 language descriptions  $L_i$  (see Appendix A1 for more examples). Each description  $l_{ij} \sim L_i$  is approximately 2-4 sentences long.

#### 4.2 Obtaining Feature Representations from Large Language Models

With the collected language descriptions, we now incorporate them into policy training. One common choice is to train a language module (e.g., long short-term memory (LSTM) [42]) from scratch

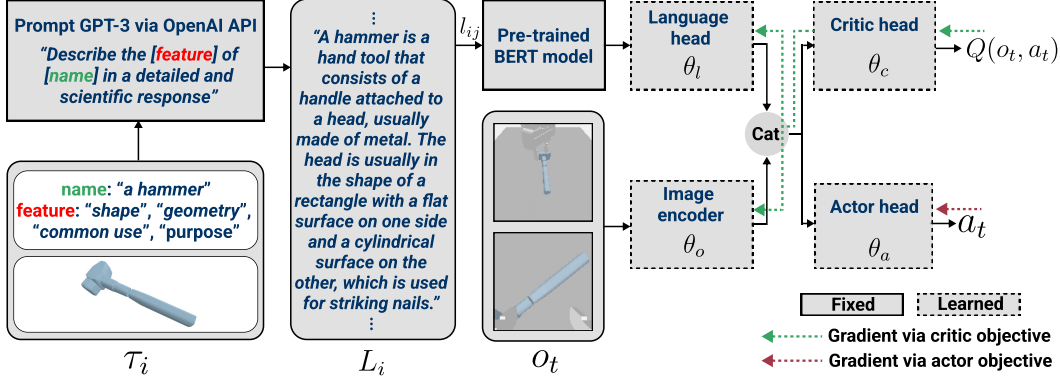


Figure 2: **Model Overview.** First, we prompt OpenAI GPT-3 to obtain a diverse set of language descriptions  $L_i$  of the tool  $\tau_i$ . Then for each episode collected, we sample a language description  $l_{ij}$  randomly from  $L_i$ , which is then fed into a pre-trained BERT model to obtain the representation. The language head further distills the language information. We concatenate the representations from the language head and the image encoder, and then the features are shared by the critic head and the actor head.

to embed features of the language input, which can take substantial time and effort to tune. Instead, we use a pre-trained LLM to distill the language descriptions into feature representations. Since LLMs are trained with vast amounts of data, they can better interpret and generalize to the diverse set of long descriptions. Since the GPT-3 model is not publicly available, we opt for the Google BERT-Base [43] model on HuggingFace, which has 110.1M parameters and outputs a 768-dimensional vector representation for each description input. T-SNE analysis shown in Fig. 3 demonstrates that, without any fine-tuning, the model already captures semantic meanings of the descriptions among tools (e.g., hammer and mallet are close to each other).

Fig. 2 shows the overall neural network architecture. We first prompt GPT-3 to obtain a set  $L_i$  of text descriptions for tool  $\tau_i$  via the procedure in Sec. 4.1. During meta-training, we randomly sample  $l_{ij}$  from  $L_i$  for each episode to ensure that the policy sees a diverse set of descriptions. We then freeze the BERT model during policy training. The output from BERT is fed into a single fully-connected layer with ReLU (language head,  $\theta_l$ ). The image observations (possibly from two camera angles—one from overhead and one from the wrist) are passed through convolutional layers (image encoder,  $\theta_o$ ), whose output is then concatenated with that from the language head. The actor head ( $\theta_a$ ) and critic head ( $\theta_c$ ) then output the action  $a_t$  and the corresponding value for the Q function  $Q(o_t, a_t)$ . See Appendix A3 for more details of the neural network setup.

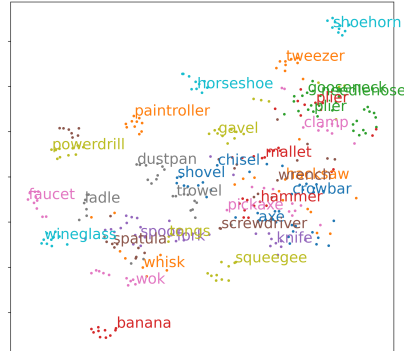


Figure 3: T-SNE results of BERT output for multiple language descriptions of each tool.

### 4.3 Meta-training and Testing Language-Conditioned Manipulation Policies

We hypothesize that additional language information of the tools promotes generalization. However, zero-shot transfer to unseen tools can be difficult given the distinct geometries and affordances. Thus we perform meta-training for explicitly training the policy to adapt to distinct tools within the same task. Our algorithm (shown in Algorithm 1) is based on Reptile [5], a simple first-order gradient-based meta-learning algorithm, but with an additional sampling strategy to prevent overfitting. At each iteration of the meta-training, one tool  $\tau$  is sampled from the training set  $\mathcal{T}$ . At the base level (line 5), we run the current policy with the language description  $l_{ij}$  sampled from  $L_i$ . We then add the collected experiences in a base replay buffer  $\beta_{\text{base}}$ , and perform  $B$  iterations of off-policy updates using Soft Actor Critic (SAC) [44] in order to obtain the final policy parameters  $\theta' = [\theta'_l, \theta'_o, \theta'_a, \theta'_c]$ . Then at the meta level (line 12), the network is updated with a gradient step towards  $\theta'$ :

$$\theta_{\text{new}} \leftarrow \theta + \alpha(\theta' - \theta), \quad (1)$$

where  $\alpha$  is the meta-learning rate and  $\theta$  is the collection of the old policy parameters  $\theta_l, \theta_o, \theta_a, \theta_c$ . We also highlight the following remarks including differences to the Reptile algorithm:

- To reduce variance, the actor and critic share the parameters of language head  $\theta_l$  and the image encoder  $\theta_o$ , and both modules are updated only with the critic objective.
- To prevent overfitting to language descriptions of one tool during training, experiences collected from all tools are saved in a meta replay buffer  $\beta_{\text{meta}}$  (line 14), and  $\beta_{\text{base}}$  for each tool is initialized with random samples from  $\beta_{\text{meta}}$  (line 3). During policy update at base training, 30% of the experiences are sampled from  $\beta_{\text{meta}}$ . During test time, we do not use any experiences from  $\beta_{\text{meta}}$ . We demonstrate the effectiveness in Sec. 6.
- In applications of Reptile in supervised learning, the meta update is often averaged over  $N$  sampled environments:  $\theta_{\text{new}} \leftarrow \theta + \alpha \sum_i^N (\theta'_i - \theta)$ . However, we find that  $N = 1$  trains faster and also matches the test objective of adapting to a single tool.  $N > 1$  does not offer better performance at test time. See Appendix A6 for sensitivity analysis.
- Performing only a single meta update after  $B = 5$  iterations of base update can be inefficient. Instead of using a large learning rate  $\alpha$  which causes unstable training, we perform  $M = 2$  meta updates for each adaptation to one tool, but collect experiences to  $\beta_{\text{base}}$  only at the first update.

Algorithm 2 shows the procedure of adaption at test time. First, the model is provided with a target test tool  $\tau_\nu$  and a set of language descriptions. With the adaptation budget  $B_\nu$ , we run the policy with the language description  $l_j$  sampled from  $L_\nu$ . The collected experiences are stored in the buffer  $\beta_{\text{base}}$  and used to update the policy parameters.

---

**Algorithm 1** ATLA: Meta-training,  $N = 1$

---

**Require:**  $\mathcal{T} = \{\tau_i\}_{i=1}^K$ : training set of tools;  $\{L_i\}_{i=1}^K$ : sets of language descriptions;  $\theta_l, \theta_o, \theta_a, \theta_c$ : policy modules;  $\beta_{\text{meta}}$ : meta replay buffer.

```

1: while meta-training do
2:   Sample  $\tau$  from  $T$  # meta level
3:   Reset  $\beta_{\text{base}}$  with samples from  $\beta_{\text{meta}}$ 
4:   for  $m = 1$  to  $M$  do
5:     for  $b = 1$  to  $B$  do # base level
6:       if  $m = 1$  then
7:         Collect episodes each with  $l_{ij} \sim L_i$ ; add to  $\beta_{\text{base}}$ 
8:       end if
9:       Sample from  $\beta_{\text{base}}$  and update  $\theta_o, \theta_l, \theta_c$  with the critic objective
10:      Sample from  $\beta_{\text{base}}$  and update  $\theta_a$  with the actor objective
11:     end for
12:     Meta update  $\theta_o, \theta_l, \theta_c, \theta_a$  with Eq. 1
13:   end for
14:   Add  $\beta_{\text{base}}$  to  $\beta_{\text{meta}}$ 
15: end while

```

---

**Algorithm 2** Adaption at test time

---

**Require:**  $\tau_\nu$ : test tool;  $L_\nu$ : set of language descriptions;  $\theta_l, \theta_o, \theta_a, \theta_c$ : policy modules;  $\beta_{\text{base}} \leftarrow \emptyset$ : base replay buffer

```

1: for  $b = 1$  to  $B_\nu$  do
2:   Collect episodes each with  $l_j \sim L_\nu$ ; add to  $\beta_{\text{base}}$ 
3:   Sample from  $\beta_{\text{base}}$  and update  $\theta_o, \theta_l, \theta_c$  with the critic objective
4:   Sample from  $\beta_{\text{base}}$  and update  $\theta_a$  with the actor objective
5: end for

```

---

## 5 Experiment Setup

Through different tool manipulation tasks in simulation, we aim to investigate the following questions: (1) Does language information help achieve better adaptation to new tools? (2) Does meta-learning improve adaptation to new tools? (3) How does the choice of pre-trained LLMs affect policy training? (4) Does language information help the policy utilize tools’ affordances effectively?

**Tasks.** Four different tool manipulation tasks are implemented (see bottom of Fig. 1): (1) pushing: pushing the tool to a fixed location on the table; (2) lifting: reaching and lifting the tool up from the table to some target height; (3) sweeping: using the tool to sweep a cylinder to a fixed location on the table; (4) hammering: using the tool to hammer a peg further into a hole in a block. Solving these tasks benefits from an understanding of the geometric affordances of the tools such as the grasp location. See App. A2 for more details of the task setup including the reward functions.

**Robot.** We build custom simulation environments with a 7-DOF Franka Panda arm in the PyBullet simulator [45]. We use RGB cameras with  $128 \times 128$  image outputs, placed at different off-arm locations and at the arm wrist depending on the needs of the tasks. For all tasks, we use 4-DOF cartesian velocity (3D translation and yaw) as the action output from the policy. The arm joints are then commanded with a jacobian-based velocity controller at 5Hz. The policy does not command the gripper; instead, we use the heuristic that once the gripper is below some height, the gripper closes to grasp the object. If the grasp fails, the gripper re-opens if it rises above the threshold.

**Tools.** We collect a total of 36 objects (See App. A1 for the full list) from open sources. Most of the objects are common tools such as a hammer and an axe. Some of them are less used as tools but have distinct geometry and affordances, such as a banana whose inner curvature may help push other objects. We split the objects into a training set of 27 and a test set of 9 — we try to separate objects with similar geometry or affordances (*e.g.*, hammer and mallet) into different sets.

**Baselines.** We compare ATLA (ours) with the following (Fig. 4): **(a) AT-TinyLA** (ours): ATLA with a smaller BERT encoder (BERT-Tiny [43] with 4.4 million parameters and 128-dimensional output). **(b) AT**: ATLA without language information. **(c) AT-XL**: ATLA without language information but larger networks for  $\theta_a$  and  $\theta_c$  (matching the number of parameters of ATLA). **(d) SAC-LA**: vanilla multi-environment training with SAC and language information but without meta-training objective. **(e) SAC**: SAC-LA but without language information. SAC-LA and SAC training follow Algorithm 3 in Appendix A4; they do not perform inner adaptation ( $B = 0$ ) at training time, and the gradient update is averaged over experiences sampled from any environment ( $N = \infty$ ). All meta learning baselines use  $B = 5$  and  $N = 1$ . See Appendix A6 for sensitivity analyses on  $N$  and  $B$ .

**Metric.** For all experiments, we save the model checkpoint with the highest running-average reward on the training dataset. After training, for each test tool we load the checkpoint and run a fixed number of iterations of adaptation. In Fig. 4, we report the highest reward at adaptation, averaged over 3 seeds for each test tool. See Appendix A6 for reward in numbers for each tool and task.

## 6 Results

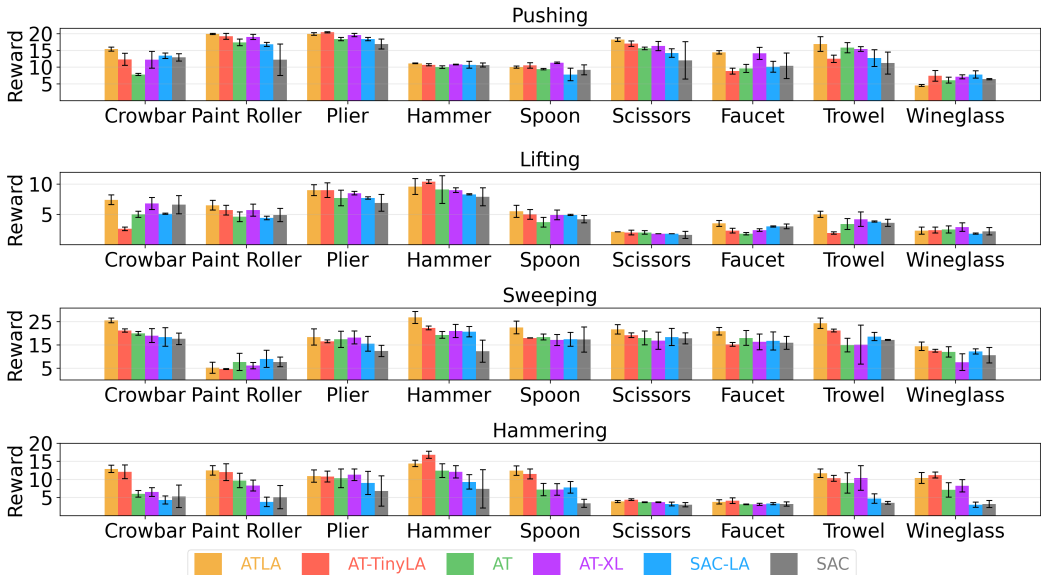


Figure 4: Post-adaptation reward in mean and standard deviation over 3 seeds across 4 tasks and 9 test tools.

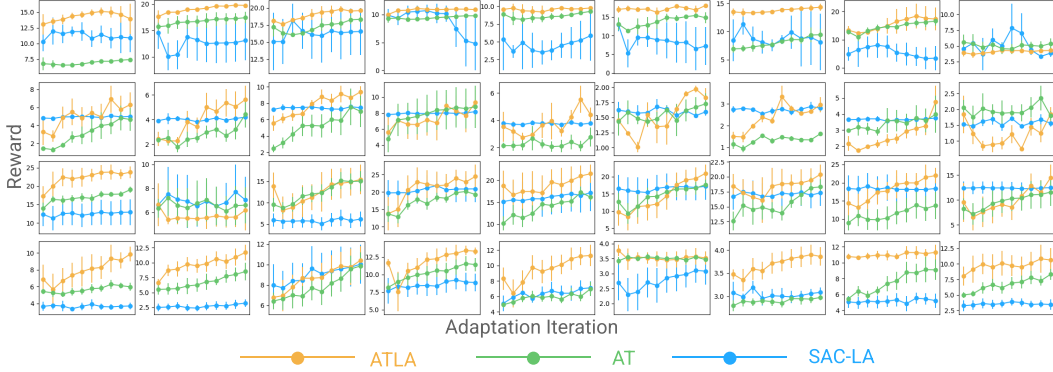


Figure 5: Curves of adaptation iteration vs. reward of the test tools (columns) in the four tasks (rows).

**Q1: Does language information help achieve better adaptation to new tools?** Fig. 4 shows that ATLA and AT-TinyLA perform better than AT and AT-XL among most tools in the 4 tasks. The differences are more significant in sweeping and hammering, which are more difficult and language information can better inform the affordances of the tools. Among the 9 tools, ATLA and AT-TinyLA always perform better with crowbar, plier, hammer, scissors, faucet, and trowel. ATLA and AT-TinyLA do not perform better mostly when low reward is achieved for all baselines for that tool, such as wineglass in pushing and lifting and paint roller in sweeping. Fig. 5 also shows that language helps faster learning in ATLA compared to AT in most cases, with the agent achieving higher rewards with fewer episodes of adaptation.

We also find that different tools learn better for different tasks. For example, hammer is better in sweeping than plier probably due to its long bar, but worse in pushing also due to the small inertia along the long bar causing instability during pushing. As ATLA performs better among most tools, language information can provide useful information about the tool affordances in different tasks.

**Q2: Does meta-learning improve adaptation to new tools?** Fig. 4 shows that with or without language information (ATLA / AT-TinyLA vs. SAC-LA, or AT / AT-XL vs. SAC), meta-learning improves final performance after adaptation. Without meta-learning, SAC-LA shows smaller improvement over SAC (*e.g.*, plier, hammer, scissors, and trowel). This demonstrates that language information particularly helps when combined with meta-learning. Fig. 5 also compares the adaptation curves between ATLA and SAC-LA: those of SAC-LA tend to stagnate or fluctuate while those of ATLA tend to rise steadily. This indicates meta-learning trains the policy to better adapt to new tools after training. In Appendix A6 we also perform sensitivity analysis on the number of inner adaptation during meta-training,  $B$ , which highlights the effectiveness of performing a few steps of inner adaptation.

**Ablation: meta replay buffer.** We investigate the effect of re-using experiences from other tools (saved in  $\beta_{meta}$  during adaptation). For this, we run ATLA without  $\beta_{meta}$  on the pushing task. Across the 9 test tools, the post-adaptation reward is mostly lower (−20%, −10%, −13%, −18%, +5%, −25%, −18%, −18%, −5%) compared to ATLA with  $\beta_{meta}$ . Note that the effect is more prominent when the reward difference between ATLA and AT is larger (*e.g.*, 20% with crowbar and 25% with scissors), indicating that language information is more effective if  $\beta_{meta}$  is applied. We also find using  $\beta_{meta}$  accelerates the meta-learning process (Fig. 6) — demonstrating that sharing experiences among tools makes training more efficient.

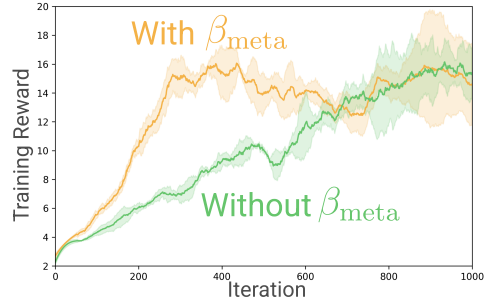


Figure 6: Using meta replay buffer also accelerates training.

**Q3: How does the choice of pre-trained LLMs affect policy training?** Fig. 4 shows that ATLA usually attains higher post-adaptation reward than AT-TinyLA, which uses a smaller pre-trained BERT model, indicating that the policies benefit from the richer representation of the language descriptions that the bigger BERT model offers.

**Q4: Does language information help the policy utilize tools’ affordances effectively?** The results above have shown that language descriptions of the geometric features and common use of the tools help policies adapt to new tools for a given task. Fig. 7(a) visualizes the effect in the example of using a crowbar for sweeping. Language descriptions of a crowbar often contain phrases including “long and thin bar”, “curved”, “hook”, “used to leverage”, and “used to pry open things”. With the descriptions, ATLA (orange curve in Fig. 7(a)) enables the policy to adapt quickly to this tool unseen during meta-training — the policy learns to use the curved hook to better steer the cylinder towards the target. As a comparison, we replace the descriptions with only the sentence “A crowbar is a long and thin bar,” and the policy (green curve in Fig. 7(a)) does not adapt as well.

One common feature among tools is the handle. Language descriptions of a trowel includes phrases like “flat, triangular blade”, “handle to be grasped”, and “used for scooping”. While ATLA learns to grasp at the handle (Fig. 7(b) top), when we remove “handle” from all the descriptions, the robot fails to grip firmly on the handle and loses the grip eventually (Fig. 7(b) bottom).

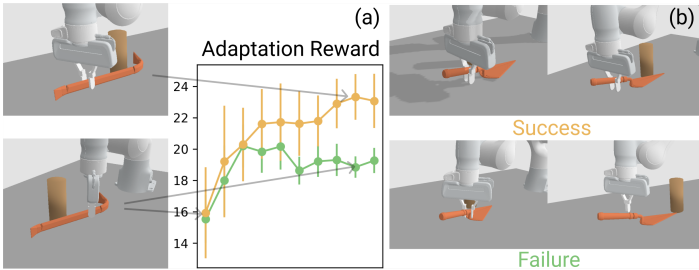


Figure 7: With language information, ATLA is able to adapt the policy to utilize the affordances of the tools – (a) curved hook on a crowbar; (b) handle on a trowel.

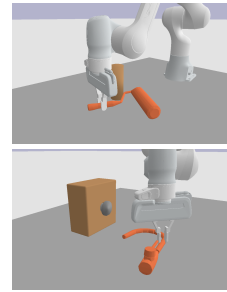


Figure 8: Policies may fail to correctly utilize affordances.

**Limitations.** In some cases, we observe that the policies can still fail to correctly utilize the tool affordances with language information. In the sweeping task, the paint roller is the only tool that ATLA fails to perform the best with. Fig. 8 (top) shows the grasp learned by the policy. It failed to use the bigger opening between the roller and the handle on the other side to sweep the cylinder. Fig. 8 (bottom) shows that in the hammering task, the policy fails to use the faucet, specifically, its relatively flat head to push towards the nail. However, these affordances can be sensitive to the initial pose of the tool and can be difficult to explore. This also highlights one of the limitations of our work. We use a relatively simple task policy setup for tool manipulation tasks, which is directly mapping image inputs to Cartesian velocity commands. This creates challenges in exploration and learning the skills even with language information. One remedy is to combine with approaches like keypoint-based methods [7] that inject additional domain knowledge into the policy. In addition, we use a relatively small dataset of tools, which may limit the potential of using language information. It would be particularly interesting to model revolute joints of the tools and perform more complex tasks such as picking up objects with a pair of tongs. Furthermore, our current evaluation does not consider real robot experiments, which is *not* the focus of this work.

## 7 Conclusion

In this work, we investigate using large language models (LLMs) to accelerate adaptation of policies to new tools in tool manipulation tasks. We use LLMs to both (1) generate diverse language descriptions of the tool geometry and common use, and (2) obtain vector representations of the descriptions. We then propose language-conditioned meta-learning that trains policies to quickly adapt to new tools. The results demonstrate that combining language information and meta-learning significantly improves the performance when adapting to unseen tools.



## Acknowledgments

The authors were partially supported by the Toyota Research Institute (TRI), the NSF CAREER Award [#2044149], the Office of Naval Research [N00014-21-1-2803, N00014-18-1-2873], and the School of Engineering and Applied Science at Princeton University through the generosity of William Addy '82. This article solely reflects the opinions and conclusions of its authors and not NSF, ONR, Princeton SEAS, TRI or any other Toyota entity.

## References

- [1] A. H. Taylor and R. D. Gray. Is there a link between the crafting of tools and the evolution of cognition? *Wiley Interdisciplinary Reviews: Cognitive Science*, 5(6):693–703, 2014.
- [2] K. R. Gibson. Evolution of human intelligence: The roles of brain size and mental construction. *Brain, Behavior and Evolution*, 59(1-2):10–20, 2002.
- [3] J. J. Gibson. The theory of affordances. *Hilldale, USA*, 1(2):67–82, 1977.
- [4] M. T. Mason. Toward robotic manipulation. *Annual Review of Control, Robotics, and Autonomous Systems*, 1(1), 2018.
- [5] A. Nichol, J. Achiam, and J. Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- [6] C. C. Kemp, A. Edsinger, and E. Torres-Jara. Challenges for robot manipulation in human environments. *IEEE Robotics & Automation Magazine (RAL)*, 14(1):20–29, 2007.
- [7] K. Fang, Y. Zhu, A. Garg, A. Kurenkov, V. Mehta, L. Fei-Fei, and S. Savarese. Learning task-oriented grasping for tool manipulation from simulated self-supervision. *The International Journal of Robotics Research (IJRR)*, 39(2-3):202–216, 2020.
- [8] A. Xie, F. Ebert, S. Levine, and C. Finn. Improvisation through physical understanding: Using novel objects as tools with visual foresight. *arXiv preprint arXiv:1904.05538*, 2019.
- [9] R. Holladay, T. Lozano-Pérez, and A. Rodriguez. Force-and-motion constrained planning for tool use. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019.
- [10] M. Qin, J. Brawer, and B. Scassellati. Rapidly learning generalizable and robot-agnostic tool-use skills for a wide range of tasks. *Frontiers in Robotics and AI*, 8, 2021.
- [11] Y. Zhu, Y. Zhao, and S. Zhu. Understanding tools: Task-oriented object modeling, learning and recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [12] T.-Y. Yang, A. Lan, and K. Narasimhan. Robust and interpretable grounding of spatial references with relation networks. In *Findings of the Association for Computational Linguistics: EMNLP*, 2020.
- [13] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, D. Ho, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, E. Jang, R. J. Ruano, K. Jeffrey, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, K.-H. Lee, S. Levine, Y. Lu, L. Luu, C. Parada, P. Pastor, J. Quiambao, K. Rao, J. Rettinghouse, D. Reyes, P. Sermanet, N. Sievers, C. Tan, A. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, S. Xu, and M. Yan. Do as I can and not as I say: Grounding language in robotic affordances. In *arXiv preprint arXiv:2204.01691*, 2022.
- [14] S. Tellex, T. Kollar, S. Dickerson, M. R. Walter, A. G. Banerjee, S. J. Teller, and N. Roy. Understanding natural language commands for robotic navigation and mobile manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2011.
- [15] C. Matuszek, E. Herbst, L. Zettlemoyer, and D. Fox. Learning to parse natural language commands to a robot control system. In *Proceedings of the International Symposium on Experimental Robotics*, 2013.

- [16] Y. Artzi and L. Zettlemoyer. Weakly supervised learning of semantic parsers for mapping instructions to actions. *Transactions of the Association for Computational Linguistics*, 1:49–62, 2013.
- [17] D. Chen and R. Mooney. Learning to interpret natural language navigation instructions from observations. In *Proceedings of AAAI Conference on Artificial Intelligence*, 2011.
- [18] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [19] S. R. Branavan, H. Chen, L. Zettlemoyer, and R. Barzilay. Reinforcement learning for mapping instructions to actions. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, 2009.
- [20] J. MacGlashan, M. Babes-Vroman, M. desJardins, M. L. Littman, S. Muresan, S. Squire, S. Tellex, D. Arumugam, and L. Yang. Grounding english commands to reward functions. In *Proceedings of Robotics: Science and Systems (RSS)*, 2015.
- [21] D. Bahdanau, F. Hill, J. Leike, E. Hughes, S. A. Hosseini, P. Kohli, and E. Grefenstette. Learning to understand goal specifications by modelling reward. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [22] P. Goyal, S. Niekum, and R. J. Mooney. Using natural language for reward shaping in reinforcement learning. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- [23] E. C. Williams, N. Gopalan, M. Rhee, and S. Tellex. Learning to parse natural language to grounded reward functions with weak supervision. In *Proceedings of the IEEE/RSJ International Conference on Robotics and Automation (ICRA)*, 2018.
- [24] T.-Y. Yang, M. Y. Hu, Y. Chow, P. J. Ramadge, and K. Narasimhan. Safe reinforcement learning with natural language constraints. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [25] T. R. Sumers, M. K. Ho, R. D. Hawkins, K. Narasimhan, and T. L. Griffiths. Learning rewards from linguistic feedback. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [26] J. D. Co-Reyes, A. Gupta, S. Sanjeev, N. Altieri, J. Andreas, J. DeNero, P. Abbeel, and S. Levine. Guiding policies with language via meta-learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [27] S. Karamcheti, D. Sadigh, and P. Liang. Learning adaptive language interfaces through decomposition. In *Proceedings of the First Workshop on Interactive and Executable Semantic Parsing*, 2020.
- [28] A. W. Hanjie, V. Y. Zhong, and K. Narasimhan. Grounding language to entities and dynamics for generalization in reinforcement learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [29] S. Branavan, D. Silver, and R. Barzilay. Learning to win by reading manuals in a Monte-Carlo framework. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 2011.
- [30] V. Zhong, T. Rocktäschel, and E. Grefenstette. Rtfm: Generalising to novel environment dynamics via reading. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [31] K. Narasimhan, R. Barzilay, and T. Jaakkola. Grounding language for transfer in deep reinforcement learning. *Journal of Artificial Intelligence Research*, 63:849–874, 2018.

- [32] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [33] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel. A simple neural attentive meta-learner. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [34] A. Rajeswaran, C. Finn, S. M. Kakade, and S. Levine. Meta-learning with implicit gradients. In *Advances in Neural Information Processing (NeurIPS)*, 2019.
- [35] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. B. Tenenbaum, H. Larochelle, and R. S. Zemel. Meta-learning for semi-supervised few-shot classification. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [36] C. Finn, K. Xu, and S. Levine. Probabilistic model-agnostic meta-learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [37] R. Dubey, E. Grant, M. Luo, K. Narasimhan, and T. Griffiths. Connecting context-specific adaptation in humans to meta-learning. *arXiv preprint arXiv:2011.13782*, 2020.
- [38] R. Vuorio, S.-H. Sun, H. Hu, and J. J. Lim. Multimodal model-agnostic meta-learning via task-aware modulation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [39] S. W. Yoon, J. Seo, and J. Moon. Tapnet: Neural network augmented with task-adaptive projection for few-shot learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.
- [40] L. Lan, Z. Li, X. Guan, and P. Wang. Meta reinforcement learning with task embedding and shared policy. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- [41] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [42] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [43] I. Turc, M.-W. Chang, K. Lee, and K. Toutanova. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962*, 2019.
- [44] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- [45] E. Coumans and Y. Bai. Pybullet, a Python module for physics simulation for games, robotics and machine learning. 2016.
- [46] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.