

Language Movement Primitives: Grounding Language Models in Robot Motion

Yinlong Dai¹, Benjamin A. Christie¹, Daniel J. Evans¹, Dylan P. Losey¹, and Simon Stepputtis²

Abstract—Enabling robots to perform novel manipulation tasks from natural language instructions remains a fundamental challenge in robotics, despite significant progress in generalized problem solving with foundational models. Large vision and language models (VLMs) are capable of processing high-dimensional input data for visual scene and language understanding, as well as decomposing tasks into a sequence of logical steps; however, they struggle to ground those steps in embodied robot motion. On the other hand, robotics foundation models output action commands, but require in-domain fine-tuning or experience before they are able to perform novel tasks successfully. At its core, there still remains the fundamental challenge of connecting abstract task reasoning with low-level motion control. To address this disconnect, we propose Language Movement Primitives (LMPs), a framework that grounds VLM reasoning in Dynamic Movement Primitive (DMP) parameterization. Our key insight is that DMPs provide a small number of interpretable parameters, and VLMs can set these parameters to specify diverse, continuous, and stable trajectories. Put another way: VLMs can reason over free-form natural language task descriptions, and semantically ground their desired motions into DMPs — bridging the gap between high-level task reasoning and low-level position and velocity control. Building on this combination of VLMs and DMPs, we formulate our LMP pipeline for zero-shot robot manipulation that effectively completes tabletop manipulation problems by generating a sequence of DMP motions. Across 20 real-world manipulation tasks, we show that LMP achieves 80% task success as compared to 31% for the best-performing baseline. See videos at our website: <https://collab.me.vt.edu/lmp>

I. INTRODUCTION

Robots should be able to perform new tasks given simple user instructions. For example, imagine that a human tells their robot to “clean the plate.” Ideally, the robot reasons over this specification and the observed environment, breaks the task down into steps, and then cleans the plate. What makes this particularly challenging is that the human’s command implies motions that are never explained: e.g., picking up a cloth and moving it in a circular pattern. Consider Figure 1 where the robot realizes that it needs a cleaning implement, grasps a sponge, and then starts wiping. Unlike prior works in which the robot relies on human demonstrations [8, 12, 52] or multiple rounds of real-world experience [23, 56], we achieve *zero-shot* generalized intelligence by connecting the human’s language specification to the robot’s controlled motions.

Our work is related to recent efforts at the intersection of robotics and foundation models and view these efforts along

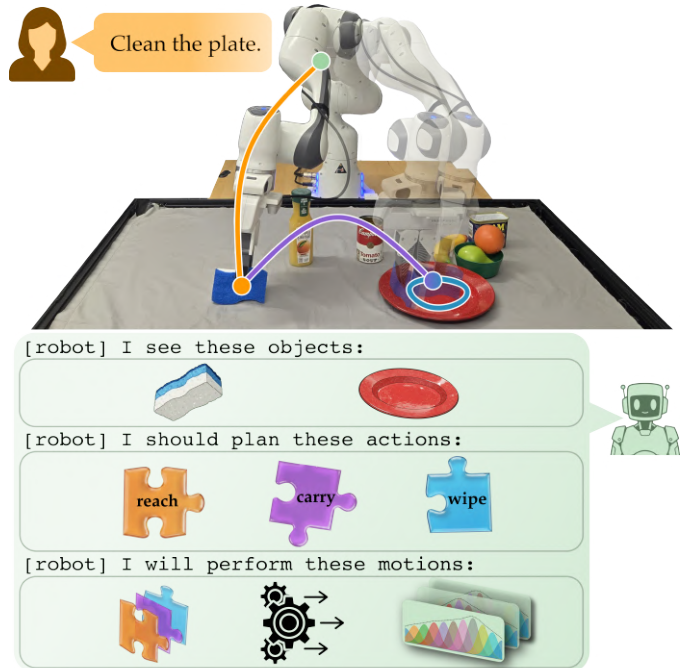


Fig. 1: Overview of Language Movement Primitives (LMP). Given a task description from the user, LMP first detects objects in the environment and composes a suitable sequence of high-level subtasks to achieve the overall goal. For each subtask, LMP then generates low-level parameters to define a Dynamic Movement Primitive (DMP). The robot tracks the continuous DMP trajectory, grounding its semantic reasoning for zero-shot robot manipulation.

a spectrum. At one extreme are methods that train models specifically on large-scale *robotic datasets*: these approaches obtain a policy that can directly convert language and vision into actions for the robot to execute [60, 29, 52, 4]. At the other extreme are methods that train on larger and more diverse *non-robotic datasets* [55, 32] with the goal of creating general-purpose foundation models, allowing these generalized vision and language models to decompose tasks and output actions at a symbolic level (e.g., “pick up the sponge”). However, both approaches have critical shortcomings. Robotics foundation models, while trained on extensive robot data, have limited common-sense reasoning capabilities and still require in-domain fine-tuning [4] or experience collection [21]. General purpose foundation models are not grounded in robotics, and thus their outputs are not immediately actionable — while the robot knows *what* to do, it does not know *how* to “pick up the sponge” with respect to the robot’s joint’s actuation.

Recent neuro-symbolic approaches try to sidestep this gap by separating high-level task reasoning from low-level robot control. In recent works like [53, 1, 35, 30] the robot combines

This work is supported in part by NSF Grant #2337884.

¹Collab, Dept. of Mechanical Engineering, Virginia Tech, Blacksburg, VA 24061. {daiyinlong, benc00, daniellevans, losey}@vt.edu

²TEA Lab, Dept. of Mechanical Engineering, Virginia Tech, Blacksburg, VA 24061. stepputtis@vt.edu

a high-level planner — which provides generalized reasoning — with a low-level controller — which converts symbolic plans into robotic motions. However, these approaches perform symbolic reasoning over *discrete* action primitives rather than the *continuous* motion parameters, constraining the system’s ability to perform complex movements. This means our fundamental bottleneck persists: to “clean the plate,” we must establish a robust connection between high-level task reasoning and low-level robotic control for precise position, velocity, and acceleration profiles [26, 2].

In this work we propose to bridge the gap and provide generalized robot policies by connecting vision and language models (VLMs) to Dynamic Movement Primitives (DMPs). We specifically consider table-top manipulation tasks. Given a user command (e.g., “clean the plate”, see Fig. 1), the VLM outputs task-based reasoning, breaking down the steps and even realizing what is missing from the human’s specification. DMPs offer a control-theoretic way to convert each of these steps into motions, parameterizing trajectories while guaranteeing the robot will reach its goal. Our insight is that:

DMPs physically ground VLMs because they embed diverse and expressive trajectories into a small number of parameters which the VLM can intuitively tune.

Our resulting system — which we refer to as **Language Movement Primitives (LMPs)** — integrates off-the-shelf VLMs into a full-stack pipeline for robotic task execution, combining the best of both worlds by leveraging the high-level reasoning capabilities of VLMs with the precise control of DMPs. During inference, LMP uses the human’s task description and an RGB-D image of the robot’s environment. By segmenting the input image, we generate a language description of the environment, including object types and locations. The VLM then reasons over this verbal state to create a high-level plan, and decides on the next subtask towards completing that plan. More specifically, the VLM chooses the parameters for the next DMP, and the robot executes this DMP in the environment. If a failure occurs, we incorporate a refinement loop where the human can provide corrective feedback (e.g., “use a towel instead”), and the robot iteratively refines its motion controller for subsequent attempts.

Overall, we see this paradigm as a step towards general-purpose robot arms that can be controlled by everyday humans without requiring additional demonstrations. We make the following contributions:

Formalizing LMPs. We frame Language Movement Primitives as an abstracted policy. The *state* is the image and generated text description of the environment, and the *action* is the weights of the next DMP. The human’s task specification is an implicit reward that the policy can leverage to choose actions, providing a semantically meaningful connection between high-level VLM planning and low-level DMP control.

Grounding Language to Control. We present a complete system that translates open-form user instructions and corrective feedback into fine-grained and stable low-level motion

controllers. Our user interface accepts commands specified in natural language, and also supports effortless and intuitive user feedback without requiring detailed knowledge of robot kinematics or low-level control specifications.

Comparing to Foundation Models. We compare LMPs to foundation models in robotics and other neuro-symbolic approaches across 20 real-world manipulation tasks. When using LMPs without any fine-tuning or refinement, our approach outperforms the best baseline by 38%. Performance particularly improves in scenarios that require trajectory shaping for obstacle avoidance or multi-stage tasks.

II. RELATED WORKS

Learning general-purpose robot policies remains an open challenge. Recent approaches range from small expert models tailored on specific robots and tasks [5, 39, 50] to large foundation models for cross-embodiment [4, 24]. Complementing these data-driven methods, neuro-symbolic approaches combine LLM reasoning with planning and motion control through symbolic representations [3, 25]. In the following subsections, we discuss these three approaches in more detail while positioning our proposed method alongside these approaches.

A. Small Expert Models for Robot Learning

Control-theoretic approaches such as [6] and [47] learn behaviors for particular tasks and robots with only a few examples, but struggle beyond their training distribution. Deep learning has demonstrated strong results for robot policy and multi-task learning by modeling complex action distributions [8, 48], whether through imitation [54, 57] or reinforcement learning [17, 20]. While effective, demonstration-centric paradigms impose a substantial burden on the user to provide suitable examples [43], while reinforcement learning methods require large-scale training due to the complexity of real-world environments [38]. Hybrid approaches that initialize policies with imitation or reinforcement learning and subsequently refine them through experience [45, 14, 19] partially address this issue; however, generalization to novel scenario remains limited. World models offer an alternative by enabling robots to plan with scenario roll-outs [18, 58], but remain computationally expensive and scale poorly to high-dimensional real-world environments.

A core challenge across these approaches — particularly in household and tabletop manipulation — is the high variety of tasks that need to be learned and the associated difficulty of obtaining sufficient training data [28, 9]. Open X-Embodiment [41] addresses this data scarcity challenge by aggregating datasets from different real-world robot tasks, allowing for positive transfer across embodiments within similar task domains. In contrast, LMP leverages the *reasoning capabilities of large foundation models* to generate task-specific robot controllers through a semantically interpretable parameter space, *removing the need for demonstrations*.

B. Robotics Foundation Models

Internet-scale foundation models have enabled progress in robot task planning [11], multimodal perception [29], and zero-shot adaptation to novel tasks and environments [35, 51]. Specifically, Vision-Language-Action (VLA) models [29, 4] leverage vision and language pipelines to ground high-dimensional data in a shared embedding space, conditioning downstream action generation [50]. RT-2 [60] popularized this VLA paradigm by co-fine-tuning vision-language models on robotic trajectory data, enabling emergent capabilities such as chain-of-thought reasoning and semantic understanding that transfers from internet-scale pre-training. Other methods prompt language models to output code that generates trajectories [35], while related works leverage LLMs as evaluators or reward models to provide feedback for policy learning [44, 1]. Despite these advances, a fundamental gap remains in *grounding high-level semantic reasoning into low-level continuous control while generalizing across diverse tasks*. LMP addresses this gap by providing a semantically meaningful parameter space such that large-scale foundation models can effectively generate low-level controllers while performing common-sense task reasoning across the control parameters.

C. Hybrid Models

Neuro-symbolic and hybrid systems have shown success by leveraging symbolic reasoning and planning as an integral part of the robot’s decision making and motion generation process [53, 11, 50]. Recent approaches incorporate neural vision and language-based components into symbolic pipelines, enabling robots to reason over unstructured sensory inputs while maintaining structured task representations [3, 36]. For example, [7] employ LLMs as neuro-symbolic task planners compatible with standard planning approaches, applying their generative capabilities to overcome limitations of traditional symbolic planners in dynamic human-robot collaboration scenarios. Leveraging LLMs enables robots to perform complex task decompositions, semantic grounding, and symbolic planning, while relying on learned controllers for execution of atomic task primitives [1, 33]. However, *existing hybrid approaches typically rely on hand-designed symbolic abstractions or a discrete library of learned controllers that are tightly coupled to specific robots and tasks*.

Beyond task planning, structured motion representation through control-theoretic approaches, such as DMP, has been used to model robot behavior [47, 46], offering a foundation with convergence guarantees while modeling motion through learned forcing functions. Recent extensions have combined DMPs with reinforcement learning [34], demonstrating their versatility for trajectory representation. Our LMP framework is inspired by these directions, leveraging a large language model for high-level reasoning, generating a set of parameters for a low-level DMP motion controller, bridging the gap between language understanding, common-sense reasoning, and fine-grained motion control.

III. PROBLEM STATEMENT

We consider settings where a robot arm is performing tabletop manipulation tasks. The human operator provides a natural language description of the task τ they want the robot to perform, and the robot converts this description into task execution (e.g., embodied motion). The robot does not have access to any in-domain demonstrations.

Observation. Let $o \in \mathcal{O}$ be the robot’s observation of its environment. In our experiments o contains the robot’s joint position (measured by onboard encoders) and an RGB-D image of the environment (taken using an external depth camera). This observation captures where the robot is, and also perceives the different objects on the table.

Policy. The robot converts the user’s task description $\tau \in \mathcal{T}$ and the environment observation o into actions:

$$a \sim \pi(o \mid \tau, o) \quad (1)$$

Here π is a *high-level* policy: the actions produced by this policy are not low-level position or joint torque commands. Instead, we treat the actions as motion primitives (with details below), and rely on the high-level policy π to select the next motion primitive that will help complete the given task τ .

Action. We define the robot’s actions $a \in \mathcal{A}$ as parameterized motions. Related neuro-symbolic works have instantiated actions as absolute positions [30], learned skills [11], or a library of general-purpose motions [15]. We instead use Dynamic Movement Primitives (DMPs) as the *actions* that connect vision and language models to embodied behaviors.

Dynamic Movement Primitives (DMP). DMPs parameterize a smooth motion between two points. Let ξ be a single degree of freedom in the robot’s end-effector space; e.g., ξ could be the robot’s position along the x axis. DMPs parameterize the trajectory $\xi(t)$ between the start state $\xi(0)$ and the goal state g as a second-order linear dynamical system [47, 31]:

$$T^2 \ddot{\xi} = \alpha(\beta(g - \xi) - T\dot{\xi}) + f(z) \quad (2)$$

Here ξ is position, $\dot{\xi}$ is velocity, and $\ddot{\xi}$ is acceleration. The temporal scaling parameter T modulates the execution speed of the trajectory, while α and β are positive constants that define a spring-damper system (regulating how quickly the trajectory is attracted to the goal position g).

We can shape the trajectory of the DMP using the nonlinear forcing function $f(z)$. In our experiments we instantiate this forcing function as a normalized weighted sum of Gaussian basis functions $\{\psi_i\}_{i=1}^N$ parameterized by the phase variable \tilde{t} and modulated by the canonical decay variable $z \in [0, 1]$:

$$f(z) = \frac{z(t) \sum_{i=1}^N w_i \psi_i(\tilde{t})}{\sum_{i=1}^N \psi_i(\tilde{t})}, \quad \tilde{t} = \min\left(\frac{t}{T}, 1\right) \quad (3)$$

Each Gaussian basis function is given by:

$$\psi_i(\tilde{t}) = \exp(-h(\tilde{t} - c_i)^2) \quad (4)$$

where the basis centers c_i are uniformly distributed over $[0, 1]$ and $h = 1/\sigma^2$ controls the basis width.

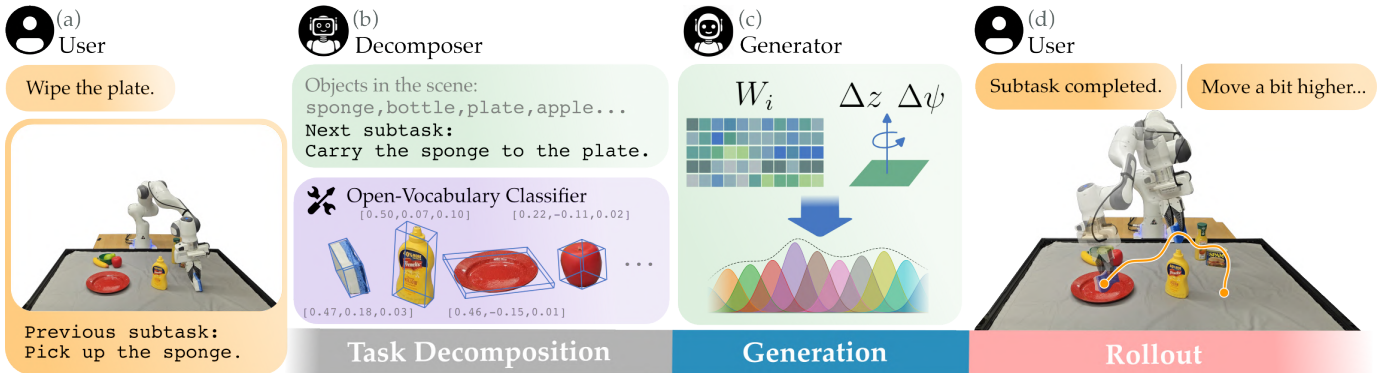


Fig. 2: LMP pipeline for a single subtask rollout. (a) The robot begins with a user-provided task description. The robot then collects an image capturing the current environment state, and remembers any previously performed subtask(s). (b) The decomposer $\pi_{\mathcal{D}}$ identifies scene objects and outputs a subtask for the next DMP to complete. An open-vocabulary classifier and depth sensing are used to estimate 3D object locations. The scene description and proposed subtask are then forwarded to the DMP weight generator $\pi_{\mathcal{G}}$. (c) The generator predicts DMP weights and auxiliary parameters to define the low-level reference trajectory. (d) The robot tracks the continuous trajectory generated from the predicted DMP parameters. Optionally, the user may observe the robot and provide natural-language feedback about any mistakes. If the user gives this refinement r , then the robot resets the rollout and the process repeats from (b).

Within our work the weights w are particularly important. Element w_i determines the relative impact of the i -th basis function; by increasing or decreasing these weights, we modulate the shape of trajectory $\xi(t)$ along a continuous spectrum. Indeed, we can specify the shape of the entire motion using weights w and goal position g . When using DMPs, the robot is guaranteed to converge to the goal as the canonical variable $z \rightarrow 0$. We update z as a time-dependent decay function of the normalized time \tilde{t} :

$$z(t) = \exp(-\gamma \tilde{t}^3) \quad (5)$$

where $\gamma > 0$ controls z 's rate of decay. In practice, z remains close to one until the end of the trajectory: when $z \rightarrow 1$, the forcing function can alter the motion profile. As $z \rightarrow 0$, the nonlinear terms are ignored and the robot converges to g .

Overall, DMPs serve as motion primitives that capture a wide range of trajectories through a small number of weights w and goals g . We know that trajectory $\xi(t)$ will converge to the goal g , and directly adjust the trajectory shape using weights w . DMPs enable rich task-specific motion shaping without sacrificing stability. Within the context of our high-level policy in Equation (1), the set of DMP weights and goals becomes the robot's action space \mathcal{A} .

Subtasks. At the start of the interaction the human provides a task description τ . The robot then completes this task in a sequence of steps. At each step the robot observes o , and queries the high-level policy π for the weights and goals of the next DMP. The robot arm rolls out the DMP and interacts with the physical environment. This process then repeats until the task is complete (outputting a null command). We refer to the steps as *subtasks* τ_i , and denote the robot's observation and action at each subtask as o_i and a_i , respectively.

Feedback. Ideally, the sequence of actions $[a_i]_0^K$ should solve the original task τ . But we recognize that no policy is perfect and that the generated weights w may not complete the task. To this end, we incorporate additional natural language feedback at the end of each subtask, allowing the robot to re-

attempt the task by refining the previously generated weights given verbal feedback. Alternatively, the human may have *new* instructions they want to provide. Our problem setting formulates this feedback as a natural language refinements r . We emphasize that this refinement is not necessary, but our framework incorporates it effectively when provided.

IV. LANGUAGE MOVEMENT PRIMITIVES (LMPs)

To convert high-level language prompts τ into low-level control trajectories we propose Language Movement Primitives (LMPs). The central idea of LMPs is a combination of VLMs and DMPs. The VLM serves as the high-level policy from Equation (1): this model takes in segmented state observations, breaks the task down into steps, and decides on the next subtask. As an action, the VLM outputs the trajectory parameters a_i for the physical robot to execute. DMPs then map these parameters into a controlled robot motion: ensuring the robot will reach the given goal and providing a continuous reference trajectory. We hypothesize that this combination will be effective because DMPs capture a diverse set of complex motions through a small number of parameters, and that these parameters are semantically meaningful (e.g., a given weight might correspond to increased motion in the x axis).

In this section we present our LMP formalism. Our paradigm starts by generating a text summary of the state from image and depth observations (Section IV-A). Once the VLM is given this language information, it then *decomposes* the overall task into the next subtask the robot should perform (Section IV-B). Finally, the VLM *generates* the goals and weights for low-level DMPs that will achieve this subtask, and the robot executes the DMPs in its physical environment (Section IV-C). Across our experiments this pipeline was sufficient for completing most tasks on the first attempt. However, if the robot fails, we also incorporate a refinement step where the user can give natural language corrections, and the robot uses those suggestions to improve its next attempt (Section IV-D). Our overall method is summarized in Figure 1, Figure 2, and Algorithm 1.

Algorithm 1 Language Movement Primitives (LMPs)

Input: task prompt τ and DMP grounding prompt s_G

- 1: $\Pi \leftarrow \emptyset$ \triangleright Set of completed subtasks
- 2: **repeat**
- 3: $o \leftarrow \text{environment_state}$
- 4: $\theta \leftarrow \pi_{\text{class}}(o \mid o)$ \triangleright Extract object poses and labels
- 5: $\varphi_i \leftarrow \pi_{\mathcal{D}}(o \mid \tau, o, \theta, \Pi)$ \triangleright Decompose for next subtask
- 6: $(W_i, \Delta z, \Delta \psi) \leftarrow \pi_{\mathcal{G}}(o \mid \varphi_i, s_G, o, \theta)$
 \triangleright Generate DMP weights and goal offsets
- 7: DMP($W_i, \Delta z, \Delta \psi$) \triangleright Execute DMP in environment
- 8: $r \leftarrow \text{judge}()$ \triangleright Collect any feedback
- 9: **if** $r \neq \emptyset$ **then**
- 10: update_prompts(r) \triangleright Add r to existing prompts
- 11: reset(φ_i) \triangleright Reset to start of subtask i
- 12: $\Pi \leftarrow \Pi \cup \varphi_i$
- 13: **until** $\varphi_i = \text{done}$

A. From Observations to State Descriptions

The robot’s observations o include the robot’s joint position and RGB-D images of the environment. We directly provide the RGB environment image as an input to the VLM policy. However, we also translate our entire observation o into a templated natural language description of the state for the VLM to reason over. For example, “there is an [object label] located at [this position] and oriented with [this orientation].” We achieve this automated segmentation by using an open-vocabulary classifier to identify objects in the environment:

$$\theta = [(l, p)_i]_0^N \sim \pi_{\text{class}}(o \mid o) \quad (6)$$

Here θ is the set of all segmented object information. This includes l , the textual object labels, and p , their 3D position and orientation in the robot’s coordinate frame. The policy π_{class} converts the observations into segmented data including the object’s label and global position in 3D space. In our experiments we use Gemini-Robotics ER [51] and LangSAM [42, 37], but other approaches for object detection and localization are also suitable [27, 40, 16]. Note that most off-the-shelf classifiers output pixel coordinates. Using the camera intrinsics, depth measurements, and the calibrated camera-to-robot extrinsics, we back-project pixels into 3D and transform them into the robot frame. Once θ is obtained, we can then use the information in θ to automatically populate a text description of the environment. Overall, the VLM inputs both the RGB image and the text description from θ .

B. From State Descriptions to Decomposed Subtasks

Now that the VLM has the task and state, we move towards outputting DMP weights for robot motion. We divide this overall policy π into two parts: *decomposition* $\pi_{\mathcal{D}}$ and *generation* $\pi_{\mathcal{G}}$. In the decomposition, the VLM outputs a language description of the next subtask τ_i . Each subtask is a step towards the completed behavior. Consider our motivating example of cleaning the plate: this task τ breaks down into three subtasks

including grasping the sponge τ_1 , carrying it to the plate τ_2 , and wiping the plate τ_3 . We want each subtask to correspond to something executable by our DMPs. We therefore provide a *template* for the decomposition policy $\pi_{\mathcal{D}}$ to complete at each high-level timestep. The subtasks within this template can be of two forms. First is ACTION(object), where the robot acts on a single object (e.g., “grasp the sponge”). Second is ACTION(object) TO(object), which is used for actions with multiple objects (e.g., “carry the sponge to the plate”). These templates ensure that the decompositions are anchored in subtasks related to scene objects and that subtasks are only involving one primary object that is manipulated, ensuring the right difficulty for each DMP.

With this framework in mind, the decomposition policy $\pi_{\mathcal{D}}$ is a foundational vision and language model of the form:

$$\varphi_i \sim \pi_{\mathcal{D}}(o \mid \tau, o, \theta, [\varphi_k]_0^{i-1}) \quad (7)$$

Here $[\varphi_k]_0^{i-1}$ represents the sequence of previously proposed subtask templates. At the start of the interaction this sequence is empty, and the sequence iteratively grows each time the decomposition policy is queried. Looking at Equation (7), τ is the high level task, and o and θ are the RGB image and natural language state description from Section IV-A. The completed subtask template φ_i output by this model describes the motion that the next DMP should complete.

C. From Subtasks to Generating DMPs

Now that we know the next desired motion, we ground this motion in the robot’s physical environment with DMPs. More specifically, given our scene description and next subtask, we sample from the *generator policy* $\pi_{\mathcal{G}}$ to produce DMP weights W_i and offset parameters Δz and $\Delta \psi$:

$$(W_i, \Delta z, \Delta \psi) \sim \pi_{\mathcal{G}}(o \mid \varphi_i, s_G, o, \theta) \quad (8)$$

The generator VLM produces weights $W_i \in \mathbb{R}^{M \times B}$ where M is the number of controlled degrees of freedom and B denotes the number of basis functions used in the DMP formulation. As a reminder, the basis functions are applied in Equation (3) to shape the DMP trajectory. The details of the grounding prompt s_G will be outlined later in this subsection.

To improve tractability and reduce unnecessary complexity, we restrict the generator to the minimal set of operational dimensions required for task execution. Our experiments focus on top-down robot manipulation tasks where the end-effector is positioned above the target object. We accordingly consider the motion of the robot end-effector in Cartesian space. Under this formulation, the weight matrix W_i is composed of the following weight vectors:

$$\mathbf{w}_i^{(x)}, \mathbf{w}_i^{(y)}, \mathbf{w}_i^{(z)}, \mathbf{w}_i^{(\theta_z)}, \mathbf{w}_i^{(g)} \in \mathbb{R}^B$$

where $\mathbf{w}_i^{(x)}$, $\mathbf{w}_i^{(y)}$, $\mathbf{w}_i^{(z)}$ correspond to translational motion along the Cartesian axes, $\mathbf{w}_i^{(\theta_z)}$ represents the end-effector rotation about the z -axis, and $\mathbf{w}_i^{(g)}$ controls the gripper state. Note that the gripper command is represented as a continuous variable rather than a binary signal, allowing it to be controlled

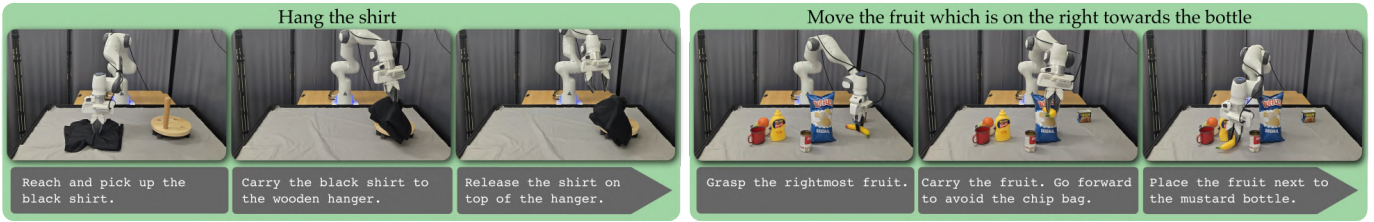


Fig. 3: Our experiments evaluate LMP’s performance on tabletop-manipulation tasks, converting natural-language task descriptions into robot controllers. In our tests, we evaluate 20 diverse household tasks requiring semantic task understanding, awareness of obstacles, and spatial reasoning.

in the same manner as the other degrees of freedom. In practice, we often just need the gripper to open or close at a specific timestep along the trajectory. We encode this knowledge in the basis functions of the gripper DMP by replacing the Gaussian in Equation (4) with a step function. More generally, designers can modify the DMP hyperparameters (e.g., types and number of basis functions) to capture their domain knowledge.

The weights W_i determine the continuous trajectory *shape*. But we also need to set the *goal* g that the DMP converges towards. This goal is specified in two parts by offsets Δz and $\Delta\psi$. We first extract the location of the center of the object from θ in Section IV-A. We then move towards a point aligned with the object’s xy position, but offset in both height and orientation. The generator policy in Equation (8) outputs a target height along the z -axis (Δz) and a target orientation around the z -axis ($\Delta\psi$).

Overall, the generator π_G provides the parameters that define DMPs in each Cartesian axis. Our experiments show that we can use off-the-shelf vision and language models as π_G , provided that we give these models some context on how DMPs operate. We input this grounding through the task-invariant prompt s_G . This prompt (available on our [website](#)) encourages the generator to internalize the relationship between the physical workspace and the DMP weight space, i.e., how variations in DMP weights induce corresponding changes in Cartesian end-effector motion. For example, the prompt explains that “increasing the weight of parameters in the x -axis will cause the trajectory to move in that direction.” Equipped with s_G , the generator policy is able to translate subtask descriptions into low-level DMP parameters; hence, the subtask φ_i only needs to be in natural, human-understandable language.

After decomposition and generation, we execute the DMP with weights W_i and goal pose g :

$$\mathbf{g} = [p_k^x \quad p_k^y \quad p_k^z + \Delta z \quad p_k^\theta \quad p_k^\phi \quad p_k^\psi + \Delta\psi] \quad (9)$$

Here p_k is the position of the target object. The position of all objects is originally estimated in θ from Section IV-A, and the target object is assigned by subtask φ_i . The robot rolls out this motion in the environment. More precisely, the robot uses the DMP as a reference trajectory, and outputs low-level control commands to track the reference.

D. Using Feedback to Refine the Motion

Ideally, the motion produced in Section IV-C is correct and the robot completes the desired subtask. But in practice this

may not be the case; we therefore enable iterative, step-by-step feedback through a *subtask judge*. The subtask judge observes how the environment changes over the course of subtask φ_i , and provides qualitative feedback in text form to the decomposer and the generator. The judge outputs a qualitative feedback statement r that describes the error and/or how the robot should improve (e.g., “too high, you missed the sponge”). To incorporate this feedback and refine the robot’s motion, we add the refinement r to the robot’s existing prompts for both the decomposer and the generator. The environment is then reset to the configuration it was in *before* the DMP was executed, and the decomposer and generator reason about appropriate corrections. Through this in-context learning mechanism, we are able to iteratively improve both the task decomposition and generated weights, leading to safer and more efficient interactions.

V. EXPERIMENTS

To demonstrate the effectiveness of LMP, we compare it to state-of-the-art baselines and conduct extensive ablations in a controlled tabletop manipulation setting involving household objects. We evaluate LMP on 20 household tasks (Figure 3), assessing overall task success rate compared to $\pi_{0.5}$ [22] and TrajGen [30]. Furthermore, we ablate our method to evaluate the impact of free-form natural language feedback on motion generation as well as the impact of grounding subtask identification through task decomposition. Videos of our experiments are available here: <https://collab.me.vt.edu/lmp>

Setup. We evaluate these methods on a 7-DoF Franka Emika Panda robot arm with a UMI gripper (see Figure 3). We use a mounted Orbbec Femto Mega camera for RGB-D imagery. For $\pi_{0.5}$, we instead use two Intel Realsense D435 cameras: one mounted on the robot end-effector and one statically mounted with an isometric view of the workspace, covering the same field of view. We use Gemini Robotics-ER 1.5 [51] to obtain object labels $\{l_i\}_{i=1}^M$. These labels are then passed to LangSAM [37, 42], which produces fine-grained object segmentations. These segmentation masks are subsequently projected onto the depth image to recover the corresponding 3D object bounding boxes. We use these bounding boxes to determine the approximate yaw of each object relative to the end-effector orientation. Leveraging its strong capabilities in physical scene understanding and multimodal reasoning for robotic tasks, we also employ Gemini Robotics-ER as our task decomposer. We use GPT-5.2 [49] as the generator policy.

TABLE I: Tasks and Performances.

Reported in success rate (%). Numbers in parenthesis (o) represent the number of feedbacks per task, averaged across trials. Columns 4-6 refer to our judge-free (J/F), decomposer free (D/F), and judge- & decomposer-free (J/F D/F) ablations. Task Categories: ● Semantic understanding ● Obstacle awareness ● Spatial reasoning

Task Name	●	●	●	LMP	TrajGen [30]	$\pi_{0.5}$ [22]	LMP J/F	LMP D/F	LMP J/F D/F
pick the chip bag on the left of the table			✓	100 (0.2)	40	20	80	40 (2.2)	20
pick the rightmost can			✓	100 (-)	20	40	100	20 (2.4)	20
pick the fruit in the middle			✓	60 (1.2)	60	40	60	80 (0.6)	60
pick the chip bag which is to the right of the can			✓	100 (-)	60	20	100	60 (0.6)	60
move the fruit which is on the right towards the bottle			✓	100 (-)	60	0	100	80 (0.6)	60
move the banana near the pear			✓	100 (-)	100	60	100	80 (0.6)	40
move the banana near the pear (obstacles included)		✓	✓	80 (1.0)	20	60	40	0 (3)	0
move the can to the center of the table			✓	80 (0.6)	20	20	80	80 (0.6)	40
move the lonely object to the others	✓		✓	80 (0.6)	0	60	80	100 (0.2)	80
place the apple in the bowl				60 (1.2)	60	40	60	40 (0.6)	0
place the apple in the bowl (obstacles included)		✓		60 (1.6)	0	20	20	60 (1)	0
pick the apple from the bowl and place it on the table				80 (0.6)	60	40	80	20 (1.2)	20
wipe the plate	✓			80 (0.6)	0	0	80	0 (3)	0
drop the ball into the cup				60 (1.4)	0	20	40	20 (1.2)	0
drop the ball into the cup (obstacles included)		✓		80 (2.4)	0	0	20	20 (1.2)	0
insert the bread into the toaster				80 (0.6)	0	0	80	20 (2.2)	0
pick up the bowl				80 (0.6)	20	40	80	100 (0.4)	60
hang the shirt	✓			40 (1.8)	0	40	40	40 (0.6)	20
put the banana on the plate			✓	80 (0.6)	80	60	80	60 (1.2)	60
put the banana on the plate (obstacles included)		✓	✓	100 (0.4)	0	40	60	40 (3)	0
Overall Performance				80	30	31	69	46	27

Feedback, when needed, is provided by a single person to ensure consistency across the evaluation.

Tasks. For our experiments, we choose a set of 20 tasks, inspired by Kwon et al. [30], which contain a diverse set manipulation challenges including spatial reasoning, semantic scene understanding, and obstacle awareness (see Figure 3).

In contrast to Kwon et al. [30], we introduce task variants that require the robot to reason about object interactions and potential collisions along its motion, which are not explicitly mentioned in the task instructions, but are implicitly relevant. To further evaluate our method’s ability to interpret semantically complex tasks, we provide generalized task descriptions. For example, we use instructions such as “Wipe the plate” instead of “Wipe the plate with a sponge.” In these cases, we rely on the decomposer to infer the implicit steps and complete the full sequence of actions.

We expect language-enabled motion generators to exhibit three core capabilities: (1) semantic understanding, inferring intent from under-specified instructions; (2) obstacle awareness, recognizing configurations the robot must avoid; and (3) spatial reasoning, understanding spatial relations between objects in the environment. We label tasks according to these characteristics to analyze method capabilities.

Metrics. We employ two performance metrics across tasks:

- 1) Success rate: For each task, the success rate is evaluated over 5 independent trials. In each trial, the positions and orientations of the target objects, as well as the obstacles and distractors, are randomly initialized. We only consider a task successful if all constituent subtasks

are completed successfully and no collisions occur.

- 2) We analyze the failure modes of the model with respect to five common categories.

A. Baseline Comparison

We compare our approach against two representative state-of-the-art baselines that use pretrained language models to guide robot policies from natural language instructions. Our selection of baselines reflects two conventions by which language models are incorporated into robotics: as common sense reasoning agents and language-grounded action generators.

TrajGen [30]. Kwon et al. [30] leverages the common sense reasoning capabilities of LLMs to generate executable motion representations. Specifically, TrajGen employs an LLM to synthesize code that generates dense waypoint trajectories based in the task instruction, which are subsequently executed by the robot. To ensure a consistent evaluation, we provide TrajGen with the same information from the object detection pipeline and employ GTP-5.2 as its internal LLM.

Pi-0.5 [22]. Vision-Language-Action models such as $\pi_{0.5}$ leverage pretrained language models to effectively ground language instructions in action and visual perception. This allows VLAs to generalize across tasks by encoding high-level semantic knowledge acquired during large-scale pretraining [60]. We train $\pi_{0.5}$ using an aggregated demonstration dataset across all 20 tasks, with five demonstrations per task. We follow the officially released Libero fine-tuning scheme [22] and fine-tune from the standard checkpoint for 20000 steps.

Baseline Results. We report the performance of LMP when compared to the baselines in Table I. We observe that LMP consistently outperforms TrajGen and $\pi_{0.5}$ (col. 1-3) with an overall success rate of 80% as compared to 30% and 31% for TrajGen and $\pi_{0.5}$ respectively. We attribute this improvement to decoupling high-level semantic reasoning from low-level motion generation: the semantically interpretable parameter space allows general-purpose foundation models to specify task intent without requiring knowledge of robot dynamics, while the low-level controller handles trajectory execution.

We observe that TrajGen particularly struggles with tasks that require nonlinear motions, such as “wipe the plate”, which requires the robot to move in a circular pattern. TrajGen tends to output linear motion, whereas LMP is capable of generating smooth splines. Additionally, in tasks that require semantic understanding, TrajGen often neglects task-relevant objects that are not directly specified in the task description.

We also found that $\pi_{0.5}$ experiences causal confusion [10]. For example, if an orange is always near an apple during data collection, $\pi_{0.5}$ will fail to develop a strong semantic distinction between these two objects. We also observe that when the task-relevant entities do not provide a strong visual signal, such as referencing the “middle of the table” as opposed to “the plate” as a target position, it is more challenging for $\pi_{0.5}$ to complete the task. This observation is common when fine-tuning a vision-language backbone for low-level control, as it diminishes its overall reasoning capability [59, 13].

B. Ablation Study

Two key components of our method are the ability to incorporate feedback and the utilization of a task decomposition module. To assess the impact of these two components, we conducted extensive ablation studies (see Table I, col 4-6).

Judge-Free Ablation (J/F). We test our method without the judge feedback outlined in Section IV-D. This ablation (LMP J/F) focuses the success of generating DMP weights on the first try, removing the ability to refine the motion through feedback.

Decomposer-Free Ablation (D/F). We also evaluate our method without the task decomposer described in Section IV-B (LMP D/F). In this setting, the generator VLM is prompted to infer a sequence of DMPs, goals, and offsets, decomposing the task implicitly. This ablation aims to demonstrate the importance of a dedicated decomposer to anchor the DMP generation in appropriate subtask complexity.

Judge- & Decomposer-Free Ablation (J/F D/F). We evaluate a further ablation (LMP J/F D/F) of *Decomposer-Free* by removing expert feedback from the generator policy. This ablation intends to show that expert feedback is especially critical when the generator lacks structured subtask grounding.

Ablation Results. We find that removing either the judge or the decomposer leads to a reduced performance. J/F achieves a success rate of 69%, while removing the decomposer D/F degrade performance to 46%. Removing both results in a performance of 27%. The sources of failures differs across ablations, as shown in Figure 4. We find that both the

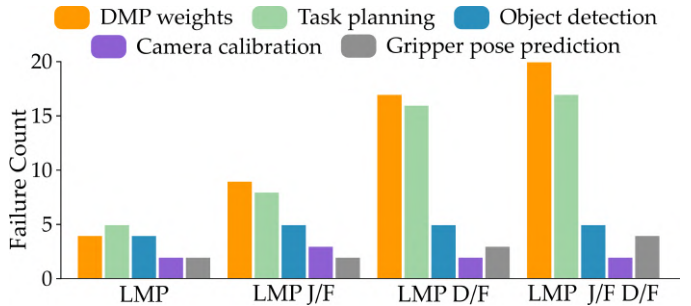


Fig. 4: We identify five dominant failure modes and analyze the impact of the judge as well as the task decomposition on these failure modes across task executions. Particularly the addition of a grounded task-decomposition substantially reduces the errors in task planning and subsequently more effective weight generation. Furthermore, we find that adding feedback has a strong influence on successful DMP weight generation.

judge and decomposer are necessary components to improve *DMP Weight Generation* and *Task Planning*. Removing the decomposer leads to a major degradation in overall task performance. Without it, LMP tends to manipulate task-irrelevant objects and reason about the task at an inappropriate level of abstraction: either too complex to be effectively captured by DMP parameterization, or too granular, at which point general-purpose foundation models lack the required knowledge to successfully reason over embodied robot dynamics. Furthermore, the judge particularly improves weight generation by allowing LMP to refine DMP parameters.

In summary, these ablations highlight the complementary roles of the judge and decomposer in grounding the foundation model at an appropriate level of abstraction.

VI. CONCLUSION

We presented a framework to bridge the gap between language understanding and motion control. Our core idea was to ground vision-language models into control-theoretic trajectories through dynamic movement primitives. We hypothesized that this combination is effective because DMPs have a relatively small number of semantically meaningful parameters that VLMs can inherently tune; hence, the VLM should be able to specify these weights and translate its high-level semantic reasoning into low-level robot behaviors. Our experiments across 20 tabletop manipulation tasks support this hypothesis: when using LMP the robot completed tasks on its first attempt 69% of the time, and 80% of the time when given up to 3 rounds of natural language corrections. This contrasts with VLM approaches that do not incorporate DMPs (30% success rate), and with robotics foundation models that need in-domain fine-tuning (31% success rate).

Limitations. While LMP provides a framework that bridges the gap between high-level semantic reasoning and low-level control, it requires the low-level controller to have a set of semantically interpretable parameters. Furthermore, due to the nature of the low-level controller, dynamically changing environments can currently not be handled and requires future work into dynamics modeling. Finally, while feedback from a human judge has proven useful, we will investigate the use of an autonomous judge, such as a VLM, in future work.

REFERENCES

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, et al. Do as I can and not as I say: Grounding language in robotic affordances. In *Conference on Robot Learning (CoRL)*, 2023.
- [2] Shuanghao Bai, Wenxuan Song, Jiayi Chen, Yuheng Ji, et al. Embodied robot manipulation in the era of foundation models: Planning and learning perspectives. *arXiv preprint arXiv:2512.22983*, 2025.
- [3] Sarthak Bhagat, Samuel Li, Joseph Campbell, Yaqi Xie, Katia Sycara, and Simon Stepputtis. Let me help you! Neuro-symbolic short-context action anticipation. *IEEE Robotics and Automation Letters*, 9(11):9749–9756, 2024.
- [4] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, et al. π_0 : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, et al. RT-1: Robotics transformer for real-world control at scale. In *Robotics: Science and Systems (RSS)*, 2023.
- [6] Joseph Campbell and Heni Ben Amor. Bayesian interaction primitives: A SLAM approach to human-robot interaction. In *Conference on Robot Learning*, pages 379–387, 2017.
- [7] Alessio Capitanelli and Fulvio Mastrogiovanni. A framework for neurosymbolic robot action planning using large language models. *Frontiers in Neurorobotics*, 18: 1342786, 2024.
- [8] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025.
- [9] Yinlong Dai, Andre Keyser, and Dylan P Losey. Prepare before you act: Learning from humans to rearrange initial states. *arXiv preprint arXiv:2509.18043*, 2025.
- [10] Yinlong Dai, Robert Ramirez Sanchez, Ryan Jeronimus, Shahabedin Sagheb, Cara M Nunez, Heramb Nemlekar, and Dylan P Losey. Civil: Causal and intuitive visual imitation learning. *arXiv preprint arXiv:2504.17959*, 2025.
- [11] Murtaza Dalal, Min Liu, Walter Talbott, Chen Chen, Deepak Pathak, Jian Zhang, and Ruslan Salakhutdinov. Local policies enable zero-shot long-horizon manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 13875–13882, 2025.
- [12] Sudeep Dasari, Oier Mees, Sebastian Zhao, Mohan Kumar Srirama, and Sergey Levine. The ingredients for robotic diffusion transformers. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 15617–15625, 2025.
- [13] Sombit Dey, Jan-Nico Zaech, Nikolay Nikolov, Luc Van Gool, and Danda Pani Paudel. Revla: Reverting visual domain limitation of robotic foundation models. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8679–8686, 2025.
- [14] Esraa Elelimy, David Szepesvari, Martha White, and Michael Bowling. Rethinking the foundations for continual reinforcement learning. *arXiv preprint arXiv:2504.08161*, 2025.
- [15] Jensen Gao, Bidipta Sarkar, Fei Xia, Ted Xiao, Jiajun Wu, Brian Ichter, Anirudha Majumdar, and Dorsa Sadigh. Physically grounded vision-language models for robotic manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 12462–12469, 2024.
- [16] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *International Conference on Learning Representations (ICLR)*, 2021.
- [17] Tuomas Haarnoja, Ben Moran, Guy Lever, Sandy H Huang, et al. Learning agile soccer skills for a bipedal robot with deep reinforcement learning. *Science Robotics*, 9(89):eadi8022, 2024.
- [18] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2024.
- [19] Jinwu Hu, Zihao Lian, Zhiquan Wen, Chenghao Li, Guohao Chen, Xutao Wen, Bin Xiao, and Mingkui Tan. Continual knowledge adaptation for reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2025.
- [20] Julian Ibarz, Jie Tan, Chelsea Finn, Mrinal Kalakrishnan, Peter Pastor, and Sergey Levine. How to train your robot with deep reinforcement learning: Lessons we have learned. *The International Journal of Robotics Research*, 40(4-5):698–721, 2021.
- [21] Physical Intelligence et al. $\pi_{0,6}^*$: a VLA that learns from experience. *arXiv preprint arXiv:2511.14759*, 2025.
- [22] Physical Intelligence et al. $\pi_{0,5}$: a vision-language-action model with open-world generalization. In *Conference on Robot Learning (CoRL)*, 2025.
- [23] Dmitry Kalashnikov et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning (CoRL)*, 2018.
- [24] Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Posner, and Yuke Zhu. Vision-language-action models for robotics: A review towards real-world applications. *IEEE Access*, 2025.
- [25] Leon Keller, Daniel Tanneberg, and Jan Peters. Neuro-symbolic imitation learning: Discovering symbolic abstractions for skill learning. *arXiv preprint arXiv:2503.21406*, 2025.
- [26] Muhammad Tayyab Khan and Ammar Waheed. Foundation model driven robotics: A comprehensive review. *arXiv preprint arXiv:2507.10087*, 2025.
- [27] Rahima Khanam and Muhammad Hussain. YOLOv11: An overview of the key architectural enhancements. *arXiv preprint arXiv:2410.17725*, 2024.

- [28] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, et al. DROID: A large-scale in-the-wild robot manipulation dataset. In *Robotics: Science and Systems*, 2024.
- [29] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, et al. OpenVLA: An open-source vision-language-action model. In *Conference on Robot Learning (CoRL)*, 2024.
- [30] Teyun Kwon, Norman Di Palo, and Edward Johns. Language models as zero-shot trajectory generators. *IEEE Robotics and Automation Letters*, 9(7):6728–6735, 2024.
- [31] Ge Li, Zeqi Jin, Michael Volpp, Fabian Otto, Rudolf Lioutikov, and Gerhard Neumann. ProDMP: A unified perspective on dynamic and probabilistic movement primitives. *IEEE Robotics and Automation Letters*, 8(4): 2325–2332, 2023.
- [32] Qixiu Li, Yu Deng, Yaobo Liang, Lin Luo, et al. Scalable vision-language-action model pretraining for robotic manipulation with real-life human activity videos. *arXiv preprint arXiv:2510.21571*, 2025.
- [33] Samuel Li, Sarthak Bhagat, Joseph Campbell, Yaqi Xie, Woojun Kim, Katia Sycara, and Simon Stepputtis. ShapeGrasp: Zero-shot task-oriented grasping with large language models through geometric decomposition. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10527–10534, 2024.
- [34] Shangde Li, Wenjun Huang, Chenyang Miao, Kun Xu, Yidong Chen, Tianfu Sun, and Yunduan Cui. Efficient robot manipulation via reinforcement learning with dynamic movement primitives-based policy. *Applied Sciences*, 14(22), 2024.
- [35] 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. *arXiv preprint arXiv:2209.07753*, 2022.
- [36] Yichao Liang, Nishanth Kumar, Hao Tang, Tom Silver, et al. VisualPredicator: Learning abstract world models with neuro-symbolic predicates for robot planning. In *International Conference on Learning Representations (ICLR)*, 2025.
- [37] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, et al. Grounding DINO: Marrying DINO with grounded pre-training for open-set object detection. In *European Conference on Computer Vision (ECCV)*, pages 38–55, 2024.
- [38] Jianlan Luo, Zheyuan Hu, Charles Xu, You Liang Tan, et al. SERL: A software suite for sample-efficient robotic reinforcement learning. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 16961–16969, 2024.
- [39] Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data. In *Robotics: Science and Systems (RSS)*, 2020.
- [40] Matthias Minderer, Alexey Gritsenko, Austin Stone, et al. Simple open-vocabulary object detection with vision transformers. In *European Conference on Computer Vision (ECCV)*, pages 728–755, 2022.
- [41] Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, et al. Open X-embodiment: Robotic learning datasets and RT-X models: Open X-embodiment collaboration⁰. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903, 2024.
- [42] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, et al. SAM 2: Segment anything in images and videos. In *International Conference on Learning Representations (ICLR)*, 2024.
- [43] Harish Ravichandar, Athanasios S Polydoros, Sonia Chernova, and Aude Billard. Recent advances in robot learning from demonstration. *Annual Review of Control, Robotics, and Autonomous Systems*, 3(1):297–330, 2020.
- [44] Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner. Vision-language models are zero-shot reward models for reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2024.
- [45] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2011.
- [46] Matteo Saveriano, Fares J Abu-Dakka, Aljaž Kramberger, and Luka Peternel. Dynamic movement primitives in robotics: A tutorial survey. *The International Journal of Robotics Research*, 42(13):1133–1184, 2023.
- [47] Stefan Schaal. Dynamic movement primitives—a framework for motor control in humans and humanoid robotics. *Adaptive Motion of Animals and Machines*, pages 261–280, 2006.
- [48] Nur Muhammad Shafiullah, Zichen Cui, Ariuntuya Arty Altanzaya, and Lerrel Pinto. Behavior transformers: Cloning k modes with one stone. *Advances in Neural Information Processing Systems (NeurIPS)*, pages 22955–22968, 2022.
- [49] Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, et al. OpenAI GPT-5 system card. *arXiv preprint arXiv:2601.03267*, 2025.
- [50] Simon Stepputtis, Joseph Campbell, Mariano Phielipp, Stefan Lee, Chitta Baral, and Heni Ben Amor. Language-conditioned imitation learning for robot manipulation tasks. *Advances in Neural Information Processing Systems (NeurIPS)*, pages 13139–13150, 2020.
- [51] Gemini Robotics Team et al. Gemini robotics 1.5: Pushing the frontier of generalist robots with advanced embodied reasoning, thinking, and motion transfer. *arXiv preprint arXiv:2510.03342*, 2025.
- [52] Octo Model Team et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- [53] Alvaro Velasquez, Neel Bhatt, Ufuk Topcu, Zhangyang Wang, Katia Sycara, Simon Stepputtis, Sandeep Neema, and Gautam Vallabha. Neurosymbolic ai as an antithesis to scaling laws. *PNAS Nexus*, 4(5):pgaf117, 2025.
- [54] Fan Xie, Alexander Chowdhury, M De Paolis Kaluza, Linfeng Zhao, Lawson Wong, and Rose Yu. Deep imitation learning for bimanual robotic manipulation.

- In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2327–2337, 2020.
- [55] Seonghyeon Ye, Joel Jang, Byeongguk Jeon, SeJune Joo, et al. Latent action pretraining from videos. In *International Conference on Learning Representations (ICLR)*, 2025.
- [56] Andy Zeng, Shuran Song, Stefan Welker, Johnny Lee, Alberto Rodriguez, and Thomas Funkhouser. Learning synergies between pushing and grasping with self-supervised deep reinforcement learning. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4238–4245, 2018.
- [57] Tianhao Zhang, Zoe McCarthy, Owen Jow, Dennis Lee, Xi Chen, Ken Goldberg, and Pieter Abbeel. Deep imitation learning for complex manipulation tasks from virtual reality teleoperation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 5628–5635, 2018.
- [58] Gaoyue Zhou, Hengkai Pan, Yann LeCun, and Lerrel Pinto. DINO-WM: World models on pre-trained visual features enable zero-shot planning. In *International Conference on Machine Learning (ICML)*, 2024.
- [59] Zhongyi Zhou, Yichen Zhu, Minjie Zhu, Junjie Wen, Ning Liu, Zhiyuan Xu, Weibin Meng, Yaxin Peng, Chaomin Shen, Feifei Feng, et al. Chatvla: Unified multimodal understanding and robot control with vision-language-action model. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 5377–5395, 2025.
- [60] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pages 2165–2183, 2023.