

MetaWorld: Skill Transfer and Composition in a Hierarchical World Model for Grounding High-Level Instructions

Anonymous CVPR submission

Paper ID ****

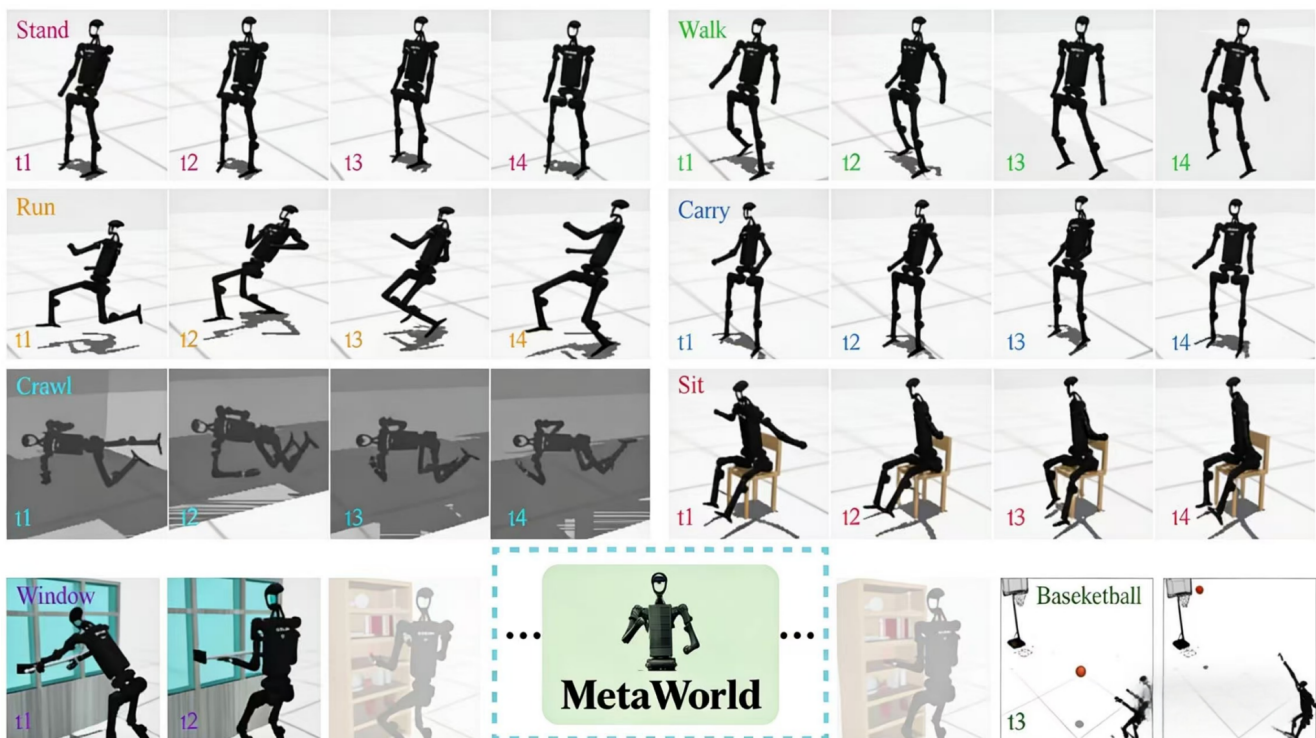


Figure 1. **MetaWorld** is a hierarchical world model framework that bridges the semantic-physical gap in robotic tasks. It leverages Vision-Language Models (VLMs) for high-level semantic task planning, integrates motion priors through a multi-expert policy library for action fusion, and achieves efficient low-level physical control via a latent dynamics model.

Abstract

001 Humanoid loco-manipulation requires coordinating
 002 whole-body joints for simultaneous locomotion and object
 003 interaction, yet it remains challenging due to the high-
 004 dimensional action space. Current methods face three pri-
 005 mary bottlenecks: (1) low sample efficiency in reinforc-
 006 ement learning, (2) poor generalization of imitation learn-
 007 ing in long-horizon tasks, and (3) physical inconsistency in
 008 VLM-based planning. To address these, we propose Meta-

World, a hierarchical world model that integrates seman- 009
 tic planning and physical control via expert policy transfer. 010
 Specifically, (i) we introduce a motion prior fusion mech- 011
 anism that leverages a pre-trained expert library within a 012
 compact latent dynamics model to accelerate online adap- 013
 tation; (ii) a hierarchical task decoupling strategy is em- 014
 ployed to decompose long-term logic from local execution, 015
 ensuring robust cross-scene transfer; and (iii) the VLM 016
 is utilized as a semantic interface to map high-level in- 017
 structions directly to pre-validated physical skills, bypass- 018

019 *ing the symbol grounding problem and ensuring dynamic*
020 *feasibility. Evaluated on challenging tasks in Humanoid-*
021 *Bench, MetaWorld significantly outperforms TD-MPC2 and*
022 *Dreamer V3 in both success rate and motion smoothness,*
023 *achieving a 139.1% increase in average reward. Our*
024 *code will be found at [https://anonymous.4open.](https://anonymous.4open.science/r/metaworld-2BF4/)*
025 *science/r/metaworld-2BF4/*

026 1. Introduction

027 Humanoid robots, as quintessential embodiments of em-
028 bodied intelligence, have long faced a fundamental chal-
029 lenge: executing semantically-driven loco-manipulation
030 tasks within unstructured and dynamic environments. At
031 the core of this challenge lies a significant “abstraction gap”
032 in current robot control systems, the disconnect between
033 high-level semantic understanding and low-level physical
034 execution [6, 13, 16, 25]. On one hand, large-scale models
035 represented by Vision-Language Models (VLMs) demon-
036 strate exceptional capabilities in high-level task planning
037 and semantic reasoning, understanding “what to do”; on
038 the other hand, low-level control methods based on imita-
039 tion learning or reinforcement learning can generate precise
040 joint-level actions, addressing “how to do it.” However, the
041 absence of a unified and scalable framework to bridge these
042 disparate capabilities often results in semantic plans that
043 are decoupled from physical constraints, conversely, low-
044 level control policies often lack the versatility to generalize
045 across complex, composed high-level tasks.

046 Current mainstream approaches exhibit distinct limita-
047 tions in addressing this problem. Although end-to-end re-
048 inforcement learning is theoretically capable of discover-
049 ing optimal policies, it suffers from extreme **sample in-**
050 **efficiency** due to the expansive search space of humanoid
051 whole-body coordination [24]. Imitation learning, while ef-
052 fective at acquiring natural motion patterns, typically yields
053 **poor generalization** when extrapolating limited demon-
054 strations to unseen, long-horizon scenarios [21]. Further-
055 more, direct application of VLMs to robot control faces se-
056 vere **physical inconsistency**; without kinematic or dynamic
057 constraints, the generated plans often fail to ground symbols
058 into executable physical actions. While existing literature
059 explores the integration of these methods, most approaches
060 remain restricted to simplistic task concatenation or a loose
061 coupling of independent modules, failing to establish a truly
062 unified, hierarchical interaction architecture [18, 20, 23].

063 To systematically address these challenges, this paper
064 presents **MetaWorld**, a hierarchical world model-based
065 robot control framework that synergistically combines the
066 semantic reasoning of VLMs, motion priors from imitation
067 learning, and online adaptation mechanisms from model-
068 based reinforcement learning. **The main contributions of**
069 **this paper include:**

- **Addressing Sample Efficiency via Motion Prior Fu-** 070
071 **sion:** We propose a hierarchical world model that lever- 072
073 ages a pre-trained multi-expert policy library. By synthe- 074
075 sizing actions within a compact **latent dynamics model**, 076
077 our framework achieves efficient reuse of skills and sig- 078
079 nificantly accelerates online adaptation. 080
• **Ensuring Physical Consistency via Semantic Inter-** 081
082 **faces:** We establish a reliable interaction paradigm using 083
084 VLMs as **semantic interfaces**. By mapping high-level 085
086 instructions directly to pre-validated physical skills, we 087
087 bypass the **symbol grounding** problem and ensure the
dynamic feasibility of generated motions.
- **Enhancing Generalization via Hierarchical Decou-** 082
083 **pling:** We introduce a modular architecture that de- 084
085 couples long-horizon semantic logic from local physi- 086
086 cal execution. This two-stage framework enables robust
cross-scene migration and achieves superior generaliza-
087 tion across diverse, unstructured environments.

088 2. Related Works

089 2.1. World Models in Robotics

090 World models facilitate policy learning by predicting en- 091
092 vironment state transitions, emerging as a critical tool 093
093 for enhancing sample efficiency and generalization in 094
094 robotics [17]. Mainstream approaches, such as the Dreamer 095
095 series [7–10], perform policy optimization through latent 096
096 space prediction, while model-based reinforcement learn- 097
097 ing methods like TD-MPC2 [11] further achieve high- 098
098 performance control in complex dynamic environments. 099
099 However, most existing world models lack hierarchical 100
100 depth, hindering the decomposition of high-level seman- 101
101 tic tasks into executable physical action sequences, thereby 102
102 limiting their direct application to long-horizon, multi-step 103
103 embodied tasks [5]. The hierarchical world model archi- 104
104 tecture proposed in this paper introduces explicit separation 105
105 between the semantic planning layer and physical execution 106
106 layer, enabling the world model to simultaneously handle 107
107 semantic parsing of language instructions and dynamic pre- 108
108 diction of physical environments, effectively bridging the
gap between task planning and action generation.

109 2.2. Skill Transfer Learning

110 Transfer learning facilitates robotic generalization through 111
111 the systematic reuse of pre-acquired skills. Current robot 112
112 control approaches include: domain adaptation (e.g., do- 113
113 main randomization) to bridge the sim-to-real gap [14]; pol- 114
114 icy fine-tuning (e.g., progressive networks) for local adjust- 115
115 ments [3]; and meta-learning (e.g., MAML) for fast few- 116
116 shot adaptation [4, 15]. However, these methods face two 117
117 major bottlenecks: a heavy reliance on target-domain data 118
118 or multi-stage training, both of which limit their practical 119
119 deployment; and a failure to achieve real-time, millisecond-

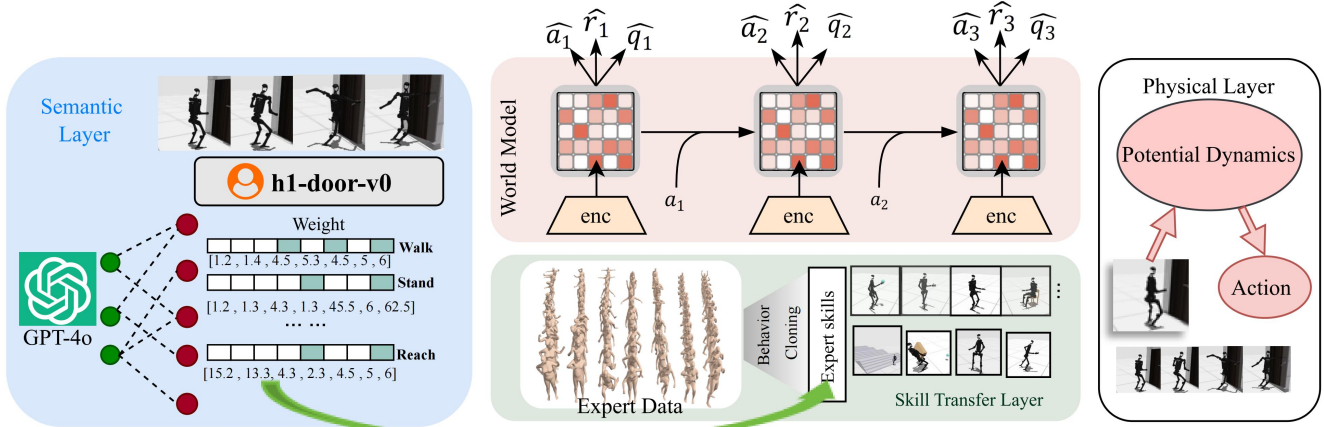


Figure 2. This illustrates the three-tier architecture of the **MetaWorld** framework: the semantic layer parses observation into executable skill sequences via a vision-language model; the skill transfer layer integrates expert policy priors through a hierarchical world model and enables dynamic adaptation; the physical layer performs precise control in a compact state space using a latent dynamics model.

120 level adaptation in the face of dynamic disturbances. To
 121 address these issues, we propose a dynamic expert selection
 122 and motion prior fusion mechanism. By constructing
 123 a multi-expert policy library, our approach dynamically
 124 selects the most relevant strategies within a Model Predictive
 125 Control (MPC) framework to enable real-time policy adjust-
 126 ment. This approach maintains the efficiency of expert
 127 policies while significantly enhancing adaptability and gen-
 128 eralization in unstructured environments.

129 2.3. VLM Applications in Real-time Task Parsing

130 Vision-Language Models (VLMs) have demonstrated re-
 131 markable proficiency in open-vocabulary perception, scene
 132 understanding, and high-level reasoning, and have been
 133 widely applied to high-level task planning and semantic
 134 guidance in robotics. Representative works such as VIMA
 135 and RT-2 attempt to directly use VLMs to generate robot
 136 action sequences or serve as front-ends for symbolic plan-
 137 ners [2, 13]. However, these methods commonly face the
 138 “symbol grounding” problem: VLM-generated plans of-
 139 ten ignore the robot’s kinematic and dynamic constraints
 140 as well as environmental physics, resulting in infeasible plans.
 141 In this work, we mitigate the physical limitations of VLMs
 142 by restricting their role to high-level semantic parsing and
 143 mapping their outputs onto a set of pre-validated, physi-
 144 cally feasible expert policies. Concurrently, our framework lever-
 145 ages the semantic understanding and task decomposition
 146 strengths of VLMs to process open-ended environmental
 147 semantics and achieve dynamic replanning through closed-
 148 loop feedback during training, guiding complex tasks to re-
 149 ference fundamental expert policies.

2.4. Transfer Learning and Policy Composition 150

151 Transfer learning aims to leverage knowledge from source
 152 tasks to accelerate learning in target domains, a paradigm
 153 particularly vital for high-DOF systems like humanoids [1,
 154 12]. Previous works have explored policy composition by
 155 linearly blending pre-trained primitives or using gating net-
 156 works to select specialized experts. However, these meth-
 157 ods often suffer from performance degradation when fac-
 158 ing long-horizon loco-manipulation tasks that require seam-
 159 less transitions between locomotion and manipulation [26].
 160 MetaWorld advances this field by integrating policy trans-
 161 fer directly into a hierarchical world model. Unlike con-
 162 ventional distillation-based methods that compress multiple
 163 experts into a single student policy, our approach employs
 164 a **dynamic expert selection** mechanism within a latent dy-
 165 namics space. By fusing motion priors from a multi-expert
 166 library, MetaWorld can compose novel motor skills on-the-
 167 fly, ensuring that the transferred knowledge is not only se-
 168 mantically relevant but also physically consistent with the
 169 robot’s hardware constraints.

3. Method 170

3.1. MetaWorld Hierarchical Architecture 171

172 The central idea of **MetaWorld** is to decompose the robotic
 173 control problem into two distinct layers: a semantic plan-
 174 ning layer responsible for interpreting task intent, and a
 175 physical execution layer responsible for generating physi-
 176 cally feasible actions. This hierarchical design can be for-
 177 malized as:

$$\pi(a_t | s_t, \mathcal{T}) = \pi_{\text{phys}}(a_t | s_t, \pi_{\text{sem}}(\mathcal{T})) \quad (1) \quad 178$$

179 where π_{sem} maps the task description \mathcal{T} to a semantic plan,
 180 and π_{phys} generates specific actions a_t based on the current

181 state s_t and the semantic plan. This architecture enables
182 independent optimization of semantic understanding and
183 physical control components while maintaining overall opti-
184 mality. The optimization objective is to maximize the ex-
185 pected cumulative reward $J(\pi) = \mathbb{E}[\sum \gamma^t r(s_t, a_t)]$ where
186 hierarchical optimization effectively addresses the distinct
187 challenges at the semantic and physical levels.

188 3.2. Semantic Planning and Symbol Grounding

189 Definition 1: *Symbol Grounding Error* refers to the state-
190 action divergence between high-level semantic intentions
191 and low-level physical feasibility. In our framework, we
192 bound this error by projecting VLM-generated semantic
193 sub-goals into a pre-trained *expert-action manifold*. Unlike
194 end-to-end methods that generate actions in unconstrained
195 spaces, MetaWorld ensures that even under semantic ambi-
196 guity, the resulting movement remains within a kinemati-
197 cally stable prior, thereby maintaining physical consistency.

198 The semantic planning layer employs a Vision-Language
199 Model (VLM) to map natural language task descriptions to
200 expert policy weights. Unlike traditional methods, we con-
201 strain the VLM output to an expert weight vector \mathbf{w} rather
202 than direct actions:

$$203 \quad \mathbf{w} = f_{\text{VLM}}(\mathcal{T}, \mathcal{E}) \quad (2)$$

204 The key innovation lies in treating the expert library as
205 a **functional basis**. By transforming the symbol ground-
206 ing problem into a **continuous manifold blending** of ex-
207 pert policies, we enable the generation of hybrid behav-
208 iors. The VLM generates expert weights through care-
209 fully engineered prompts, with the response R processed
210 by a parsing function to obtain normalized weights $w_i =$
211 $\exp(\text{extract}_i(R)) / \sum_j \exp(\text{extract}_j(R))$. Since each expert
212 policy π_{exp}^i is physically feasible, the generated semantic
213 plan $\pi_{\text{sem}}(\mathcal{T}) = \sum_i w_i \pi_{\text{exp}}^i$ naturally satisfies physical con-
214 straints. The symbol grounding error is bounded within the
215 range of expert policy differences, significantly outperform-
216 ing direct action generation methods.

217 3.3. Dynamic Adaptation Mechanism

218 To address dynamic environmental changes, we introduce
219 a **Zero-shot Behavioral Compositor** that performs state-
220 aware expert selection. Based on the current state s_t , we
221 construct a selection probability distribution:

$$222 \quad p(i | s_t) = \frac{\exp(\phi(s_t)^\top \psi(\pi_{\text{exp}}^i))}{\sum_{j=1}^K \exp(\phi(s_t)^\top \psi(\pi_{\text{exp}}^j))} \quad (3)$$

223 where ϕ is a state encoding function and ψ is an expert
224 feature extraction function. The VLM-generated seman-
225 tic weights w_i are fused with the dynamic selection prob-
226 abilities $p(i | s_t)$ to obtain the final weights $\tilde{w}_i(s_t, \mathcal{T}) =$
227 $\alpha w_i + (1 - \alpha)p(i | s_t)$. This fusion mechanism diverges from

classical discrete option-selection by enabling **continuous**
interpolation between experts. The parameter $\alpha \in [0, 1]$
controls the relative importance of semantic planning ver-
sus state awareness, enabling a balance between task consis-
tency and environmental adaptability. The reference expert
action $a_{\text{ref}} = \sum_i \tilde{w}_i(s_t, \mathcal{T}) \pi_{\text{exp}}^i(s_t)$ provides a high-quality
initial solution for the physical execution layer.

235 3.4. Physical Execution and Online Optimization

236 The physical execution layer employs the TD-MPC2 algo-
237 rithm, constructing a latent dynamics model for Model Pre-
238 dictive Control (MPC). Observations o_t are encoded into
239 latent states $z_t = f_{\text{enc}}(o_t)$ and state evolution is predicted
240 through the dynamics model $z_{t+1} = f_{\text{dyn}}(z_t, a_t)$. The
241 MPC optimization problem solves for the optimal action se-
242 quence over a future horizon H :

$$243 \quad \mathbf{a}_{t:t+H-1}^* = \arg \max_{\mathbf{a}_{t:t+H-1}} \mathbb{E} \left[\sum_{k=0}^{H-1} \gamma^k r(z_{t+k}, a_{t+k}) + \gamma^H V(z_{t+H}) \right] \quad (4)$$

244 the expert-guided action a_{ref} is incorporated into the
245 optimization objective $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{TD}} + \lambda \|a_t - a_{\text{ref}}\|^2$
246 where the Temporal-Difference (TD) learning loss $\mathcal{L}_{\text{TD}} =$
247 $\mathbb{E}[\|Q(z_t, a_t) - (r_t + \gamma Q(z_{t+1}, \pi(z_{t+1})))\|^2]$ ensures accu-
248 rate value function estimation. This design maintains online
249 adaptation capabilities while leveraging expert knowledge
250 to accelerate the learning process.

251 3.5. Theoretical Analysis and Implementation

252 Based on contraction mapping theory, we prove algorithm
253 convergence under appropriate parameter selection. The
254 value iteration operator \mathcal{T} satisfies the contraction property
255 $\|\mathcal{T}Q_1 - \mathcal{T}Q_2\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty$ guaranteeing conver-
256 gence of the value function to the optimal solution. Com-
257 pared to traditional methods, the sample complexity is re-
258 duced from $\mathcal{O}(|\mathcal{S}||\mathcal{A}|/[(1 - \gamma)^2 \epsilon^2])$ to $\mathcal{O}(1/[(1 - \gamma)^2 \epsilon^2] +$
259 $K)$ demonstrating the efficiency advantage of knowledge
260 reuse.

261 4. Experiment

262 We evaluate our method on HumanoidBench [22], a com-
263 prehensive benchmark for humanoid locomotion and ma-
264 nipulation. For locomotion, we select three tasks: walk,
265 stand, and run; for manipulation, we include door open-
266 ing. The locomotion tasks are learned via imitation learn-
267 ing from the AMASS [19] expert dataset, shaped with
268 trajectory-tracking reward signals. Other skills in the base
269 expert policy repository (e.g., reach) are inherited from TD-
270 MPC2 and are therefore excluded from evaluation. While
271 our architecture is designed with a comprehensive library of
272 eight specialized experts—namely *walk, carry, sit, crawl,*

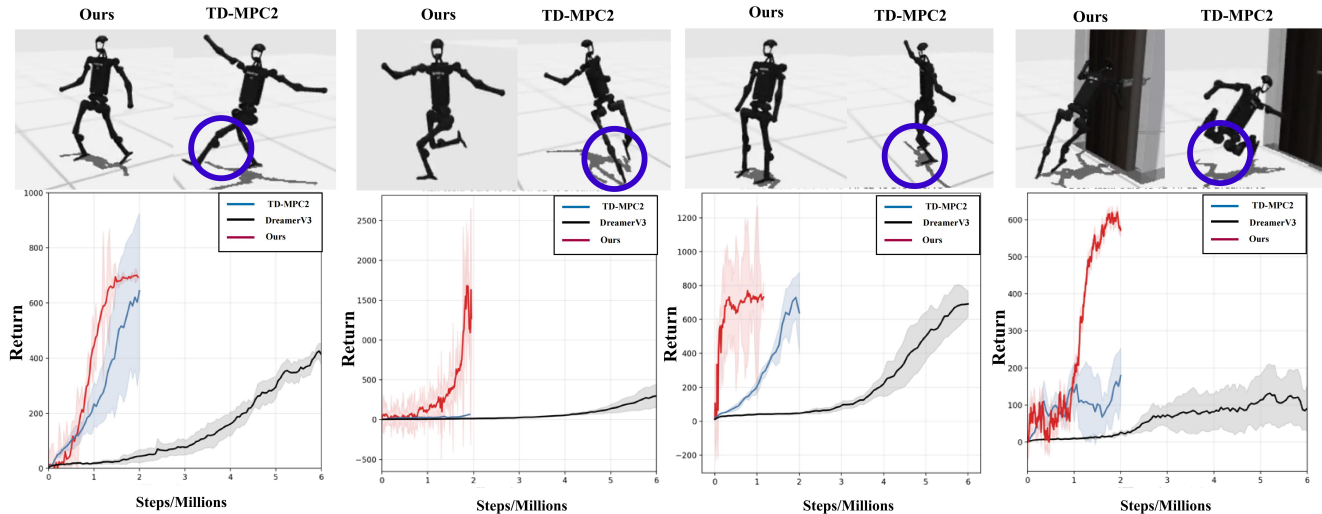


Figure 3. We evaluate our method on locomotion (run, walk, stand) and manipulation (door) tasks against baselines TD-MPC2 and DreamerV3.

Task	Metric	DreamerV3	TD-MPC2	Ours	Imp. (%)
Stand	Ret. \uparrow (\pm Std)	699.3 \pm 62.7	749.8 \pm 54.3	793.4 \pm 13.5	9.5
	Conv. \downarrow (M)	5.5	1.9	0.8	78.3
Walk	Ret. \uparrow (\pm Std)	428.2 \pm 14.5	644.2 \pm 162.3	701.2 \pm 7.6	30.77
	Conv. \downarrow (M)	6.0	1.8	1.4	64.1
Run	Ret. \uparrow (\pm Std)	298.5 \pm 84.5	66.1 \pm 4.7	1689.9 \pm 13.6	826.0
	Conv. \downarrow (M)	6.0	2.0	1.9	52.5
Door	Ret. \uparrow (\pm Std)	165.8 \pm 50.2	179.8 \pm 52.9	680.0 \pm 50.0	293.5
	Conv. \downarrow (M)	9.0	2.0	1.7	69.1
Avg.	Ret. \uparrow	398.0	410.0	966.1	139.1
	Conv. \downarrow (M)	6.6	1.9	1.5	64.7

Table 1. Performance comparison on locomotion and manipulation tasks. All values are reported to one decimal place. The **Ours** column is highlighted in light blue, and the rightmost column shows the percentage improvement.

273 *stairs, stand, run, and reach*—this evaluation focuses on a
 274 representative subset: *walk, stand, and run*. These tasks are
 275 selected to rigorously benchmark the model’s fundamental
 276 locomotion stability and its ability to manage high-dynamic
 277 transitions. Our policy library consists of $N = 8$ special-
 278 ized experts (e.g., walk, run, stand, etc.). Each expert is a
 279 neural policy trained on thousands of motion clips from the
 280 AMASS dataset, rather than a single episode mapping. This
 281 ensures that each primitive expert possesses sufficient gen-
 282 eralization within its specific modal domain (e.g., various

walking speeds and gaits).

4.1. Baseline

We evaluate our method against two state-of-the-art model-
 based reinforcement learning algorithms to demonstrate its
 superiority: TD-MPC2 [11], a scalable model-predictive
 control framework that utilizes local trajectory optimization
 with learned implicit world models, and DreamerV3 [9],
 a robust actor-critic algorithm that learns behaviors within
 a latent imagination space across diverse domains. These

292 baselines represent the current frontier in sample efficiency
 293 and asymptotic performance, providing a rigorous bench-
 294 mark for our proposed components.

295 *Distinction from Classical Hierarchical Control.* While
 296 hierarchical planning conditioned on task priors has been
 297 explored, our approach diverges from traditional *discrete*
 298 *option-selection* frameworks. Instead of switching between
 299 isolated controllers, MetaWorld treats the expert library as
 300 a **continuous functional basis**. By leveraging VLM-driven
 301 dynamic weights, we enable the **emergence of hybrid be-**
 302 **haviors** that are not present in any single expert’s training
 303 distribution. This formulation transforms the hierarchical
 304 interface from a simple router into a generative action com-
 305 positor, significantly enhancing the zero-shot adaptability in
 306 non-stationary environments.

307 4.2. Experimental setup

308 Our architecture employs the following hyperparameters:
 309 the VLM temperature parameter $\tau = 0.3$ controls the
 310 stochasticity in expert weight generation; the expert guid-
 311 ance weight $\lambda = 0.05$ balances TD-learning and expert pri-
 312 ors; the MPC planning horizon is set to $H = 3$ to ensure
 313 real-time performance; the discount factor $\gamma = 0.99$ trades
 314 off long-term rewards; the learning rate is fixed at 0.001
 315 to ensure stable convergence; the batch size of 256 balances
 316 training efficiency and memory usage; and the fusion coeffi-
 317 cient $\alpha = 0.7$ regulates the relative importance of semantic
 318 planning and state adaptation. For task success criteria, we
 319 follow the definitions provided in HumanoidBench.

320 4.3. Validation of Hierarchical Architecture Effec- 321 tiveness

322 **Purpose of the experiment.** In this section, we evaluate
 323 the effectiveness of our hierarchical architecture against tradi-
 324 tional model-based reinforcement learning methods. We
 325 first assess the performance of the base expert policy across
 326 three locomotion tasks: walk, run, and stand. Subsequently,
 327 on the walk task, we further test the model’s capability to
 328 perform complex manipulation tasks after dynamic adapta-
 329 tion. To quantitatively evaluate the performance, we record
 330 the cumulative *Return* to measure the task completion qual-
 331 ity and the number of *Convergence Steps* to reflect the sam-
 332 ple efficiency. Specifically, a higher Return indicates su-
 333 perior task execution, while fewer Steps demonstrate the
 334 model’s ability to achieve optimal performance with min-
 335 imal environment interactions, highlighting the efficiency
 336 of our hierarchical architecture.

337 **Result.** As illustrated in Table 1 and Fig. 3, the
 338 **MetaWorld** framework exhibits overwhelming superiority
 339 across four representative tasks, achieving an average return
 340 improvement of 139.1%. In basic locomotion tasks, particu-
 341 larly the remarkable **826.0% boost in the Run task**, the
 342 performance gains stem from the *manifold constraints* im-

Table 2. Success rate evaluation of base expert policies and complex manipulation tasks. We evaluate each policy using checkpoints at 2.0M steps across 10 independent trials. Our method significantly outperforms state-of-the-art MBRL baselines.

Method	Locomotion (Success / 10)			Manipulation
	Stand	Walk	Run	Door (Success / 10)
DreamerV3 [?]	2/10	2/10	1/10	0/10
TD-MPC2 [11]	3/10	3/10	2/10	1/10
Ours	9/10	9/10	9/10	8/10

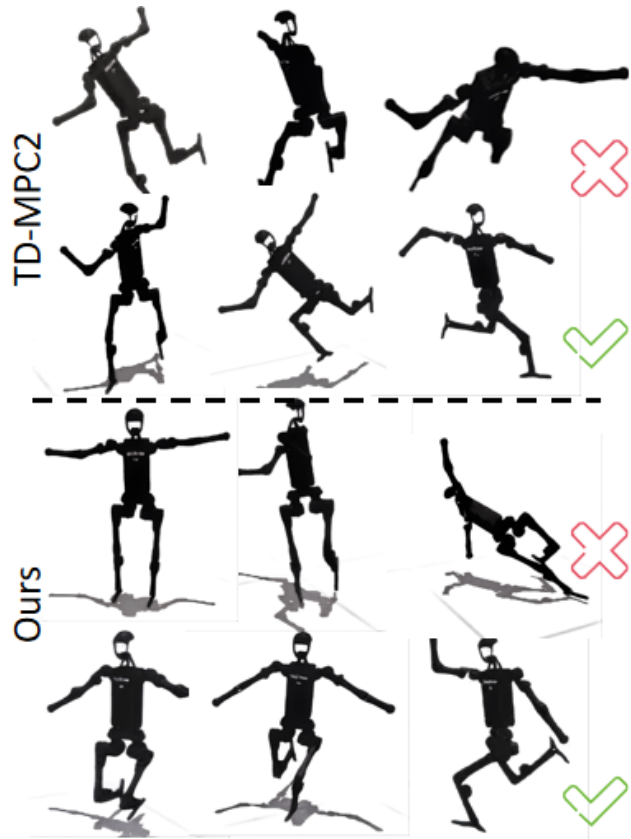


Figure 4. We visualized our method against baselines TD-MPC2.

343 posed by the imitation-learned expert library on the high-
 344 dimensional action space. By transforming continuous ac-
 345 tion search into a weighted composition of expert trajec-
 346 tories, our architecture provides the VLM planner with a
 347 physically feasible and low-dimensional action subspace,
 348 effectively mitigating the *curse of dimensionality* inherent
 349 in humanoid control.

350 For the complex Door task, the **293.5% improvement**
 351 validates our hierarchical innovation: the VLM acts as a se-
 352 mantic decoder to decompose ambiguous instructions into
 353 structured sub-goal sequences, while the dynamic selection
 354 mechanism orchestrates corresponding motion primitives

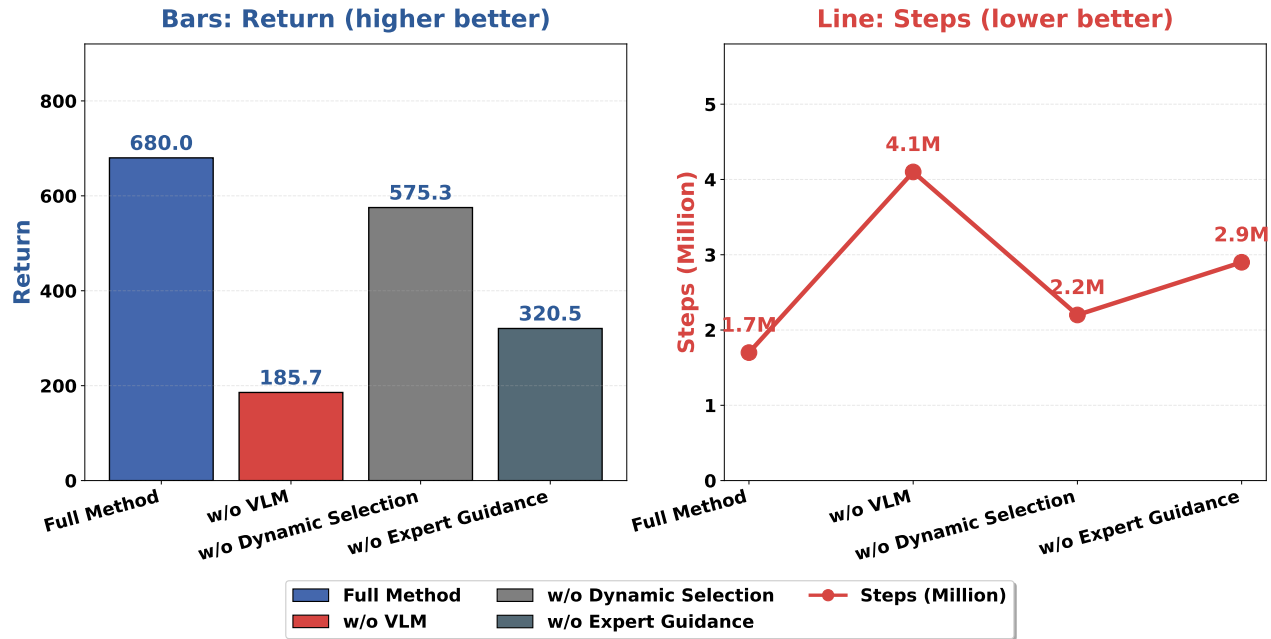


Figure 5. Ablation study results on the Door task (maximum reward and convergence steps).

355 for execution. The pivotal insight is that complex manipulation
 356 is achieved through the *decoupling and recombination*
 357 of high-level reasoning and low-level motor priors. This
 358 strategy not only resolves the long-standing *symbol ground-*
 359 *dilemma* but also bypasses the *sample efficiency bottle-*
 360 *neck* typical of end-to-end RL paradigms, offering a robust
 361 new paradigm for the seamless integration of logical infer-
 362 ence and reactive control in intelligent robotics.

363 While the use of motion priors inherently accelerates ex-
 364 ploration, the primary challenge addressed here is the *tem-*
 365 *poral composition* of these priors. The efficiency gain stems
 366 not just from the experts themselves, but from the latent dy-
 367 namics model’s ability to smoothly transition between dis-
 368 crete expertise domains, a task that remains prohibitively
 369 slow for flat RL architectures.

370 4.4. Evaluation of Task Success Rate and Stability.

371 To rigorously evaluate the stability and success rate of our
 372 base expert policies and complex task execution, we con-
 373 duct assessments using model checkpoints saved at 2.0M
 374 environment steps. For the base locomotion experts, a trial
 375 is considered successful if the humanoid maintains its bal-
 376 ance without falling for a duration of 30s. For complex ma-
 377 nipulation tasks, success is defined by the successful com-
 378 pletion of the door-opening sequence. We evaluate each
 379 policy across 10 independent trials with different random
 380 seeds and report the total number of successful completions
 381 to reflect the robustness of each method.

382 **Result.** The quantitative results in Table 2 and Fig. 4 re-
 383 veal the core advantages of our hierarchical architecture

in complex humanoid control. First, *the inferior perfor-*
mance of DreamerV3 and TD-MPC2 in base locomotion
tasks can be attributed to the challenges of exploration in
high-dimensional state spaces. Conventional Model-based
 RL algorithms attempt to learn world models and whole-
 body control policies from scratch, leading to a failure to
 maintain dynamic balance within the 2.0M sampling bud-
 get. In contrast, our method leverages **expert guidance**
 to constrain the policy search within a high-performance man-
 ifold, achieving superior sample efficiency.

Second, *the performance gap in the “Door” task (80% vs. 10%)*
highlights the necessity of coupling semantic reason-
ing with motor execution. While baselines often lose
 balance while attempting manipulation, our VLM-driven
 layer dynamically adapts expert weights based on visual
 feedback, maintaining center-of-mass stability during com-
 plex sequences. This demonstrates that our **semantic-**
guided dynamic selection mechanism is pivotal for solv-
 ing long-horizon tasks that require both precision and sta-
 bility.

4.5. Ablation study and stability analysis

Purpose of the experiment. To rigorously quantify the
 contribution of each module within the MetaWorld frame-
 work, we perform a systematic series of ablation experi-
 ments. First, we ablate semantic planning by replacing the
 VLM-based planner with a randomly initialized MLP-based
 weighting layer, where expert weights remain learnable via
 RL gradients but lack semantic guidance. Next, we ablate
 dynamic expert selection by fixing $\alpha = 1.0$ during train-

413 ing, thereby disabling state-aware expert selection. Finally,
414 we ablate skill-expert guidance by setting $\lambda = 0$ in the TD-
415 MPC2 optimization objective, removing expert action guid-
416 ance. We present experimental observations and analysis
417 for each ablation.

418 **Result.** As shown in Table 3 To systematically decou-
419 ple the contributions within the **MetaWorld** framework, we
420 conduct a series of ablation experiments on the challeng-
421 ing “Door” task. As illustrated in Fig. 5, the catastrophic
422 72.7% performance collapse when replacing the VLM with
423 a naive learnable weighting layer reveals the framework’s
424 absolute reliance on *high-level task decomposition*; even
425 with learnable parameters, without semantic grounding, the
426 agent fails to discover the complex manifold of expert com-
427 binations required to translate abstract symbols into exe-
428 cutable sub-goal sequences. This confirms that VLM-based
429 planning provides a critical *semantic warm-start* for solv-
430 ing the symbol grounding problem. Furthermore, the ex-
431 pert guidance module serves as a critical **behavioral an-
432 chor**, where its **52.9% contribution** demonstrates that in
433 high-DOF humanoid systems, naive exploration often fails
434 due to complex contact dynamics. By providing a *behav-
435 ioral manifold constraint*, expert policies restrict the search
436 space to kinematically feasible regions, facilitating a vi-
437 tal synergy between imitation-based priors and online RL.
438 Lastly, while dynamic expert selection contributes a rela-
439 tively smaller **15.4%** to the success rate, it provides es-
440 sential **closed-loop resilience**, enabling reactive recovery
441 against external perturbations or distribution shifts. Ulti-
442 mately, the framework’s superiority stems from the seam-
443 less integration of semantic reasoning, behavioral priors,
444 and dynamic adaptation.

Table 3. Ablation Study: Contribution of each component to the performance.

Method	Return (\uparrow)	Steps (M, \downarrow)
Full Method	680.0	1.7
w/o VLM	185.7	4.1
w/o Dynamic Selection	575.3	2.2
w/o Expert Guidance	320.5	2.9

445 5. Limitations

446 Despite its superior performance, MetaWorld exhibits sev-
447 eral limitations that warrant further investigation. First, the
448 current **expert policy generation** relies on a basic imitation
449 learning module with trajectory-matching reward mecha-
450 nisms, which lacks the ability to dynamically weigh fine-
451 grained motion errors at the joint level. Second, the **ex-
452 pert policy selection** mechanism employs a weighted fu-
453 sion approach rather than a truly intelligent routing archi-

454 tecture, such as a Mixture-of-Experts (MoE) system. This
455 limits the precision of skill composition and the fluidity of
456 semantically-conditioned switching between disparate be-
457 haviors. Third, the framework’s **generalization and scala-
458 bility** remain constrained; it lacks explicit few-shot adapta-
459 tion capabilities for novel, complex task combinations and
460 does not yet address potential gradient interference during
461 multi-skill coordination. These factors may hinder its per-
462 formance when scaling to even more diverse and unstruc-
463 tured real-world environments.

464 6. Conclusion

465 This study presents **MetaWorld**, a hierarchical world
466 model framework designed to bridge the semantic-physical
467 gap in humanoid robot control. By integrating VLM-
468 based semantic planning, dynamic expert policy transfer,
469 and latent dynamics models, MetaWorld establishes a cohe-
470 sive pipeline from high-level reasoning to low-level physi-
471 cal execution. Experimental results on **Humanoid-Bench**
472 demonstrate that our framework achieves a 139.1% im-
473 provement in average reward, significantly enhancing both
474 task completion efficiency and motion coherence. Mov-
475 ing forward, we plan to iteratively upgrade the frame-
476 work by incorporating dynamic reward shaping for imita-
477 tion learning, an MoE-based semantic-aware routing mech-
478 anism, and enhanced few-shot transfer mechanisms. These
479 advancements will move MetaWorld toward a more ro-
480 bust, adaptive, and scalable solution for the complex loco-
481 manipulation demands of next-generation humanoid robots.

482 References

- 483 [1] Marina Y. Aoyama, Sethu Vijayakumar, and Tetsuya
484 Narita. Few-shot transfer of tool-use skills using hu-
485 man demonstrations with proximity and tactile sens-
486 ing. 2025. 3
- 487 [2] Anthony Brohan, Noah Brown, Justice Carbajal, Yev-
488 gen Chebotar, Xi Chen, Krzysztof Choromanski,
489 Tianli Ding, Danny Driess, Avinava Dubey, Chelsea
490 Finn, et al. Rt-2: Vision-language-action models
491 transfer web knowledge to robotic control. *Science
492 Robotics*, 2023. 3
- 493 [3] Lukas Fehring, Theresa Eimer, and Marius Lindauer.
494 Growing with experience: Growing neural networks
495 in deep reinforcement learning. In *2025 Multi-
496 disciplinary Conference on Reinforcement Learning
497 and Decision Making (RLDM)*, 2025. 2
- 498 [4] Chelsea Finn, Pieter Abbeel, and Sergey Levine.
499 Model-agnostic meta-learning for fast adaptation of
500 deep networks. In *Proceedings of the 34th Inter-
501 national Conference on Machine Learning (ICML)*,
502 pages 1126–1135, 2017. 2
- 503 [5] Kentaro Fujii and Shingo Murata. Real-world robot

- 504 control by deep active inference with a temporally hi- 556
505 erarchical world model. *IEEE Robotics and Automa- 557*
506 *tion Letters (RA-L)*, 2025. 2 558
- 507 [6] Haoran Geng, Songlin Wei, Congyue Deng, Bokui 559
508 Shen, He Wang, and Leonidas Guibas. SAGE: Bridg- 560
509 ing Semantic and Actionable Parts for GGeneralizable 561
510 Articulated-Object Manipulation under Language In- 562
511 structions. In *Proceedings of Robotics: Science and 563*
512 *Systems (RSS)*, 2024. 2 564
- 513 [7] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and 565
514 Mohammad Norouzi. Dream to control: Learning be- 566
515 haviors by latent imagination. In *International Con- 567*
516 *ference on Learning Representations (ICLR)*, 2020. 2 568
- 517 [8] Danijar Hafner, Timothy Lillicrap, Mohammad 569
518 Norouzi, and Jimmy Ba. Mastering atari with discrete 570
519 world models. In *International Conference on Learn- 571*
520 *ing Representations (ICLR)*, 2021. 572
- 521 [9] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Tim- 573
522 othy Lillicrap. Mastering diverse control tasks through 574
523 world models. *Nature*, 2025. 5 575
- 524 [10] Danijar Hafner, Wilson Yan, and Timothy Lillicrap. 576
525 Training agents inside of scalable world models. *arXiv 577*
526 *preprint arXiv:2509.24527*, 2025. 2 578
- 527 [11] Nicklas Hansen, Hao Su, and Xiaolong Wang. Td- 579
528 mpc2: Scalable, robust world models for continuous 580
529 control. In *International Conference on Learning Rep- 581*
530 *resentations (ICLR)*, 2024. 2, 5, 6 582
- 531 [12] M. Hou, K. Hindriks, A. E. Eiben, and K. Baraka. 583
532 Robot policy transfer with online demonstrations: An 584
533 active reinforcement learning approach, 2025. 3 585
- 534 [13] Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan 586
535 Zhang, and Fei-Fei Li. Rekep: Spatio-temporal reason- 587
536 ing of relational keypoint constraints for robotic 588
537 manipulation. In *Conference on Robot Learning 589*
538 *(CoRL)*, 2024. 2, 3 590
- 539 [14] Yuankun Jiang, Chenglin Li, Wenrui Dai, Junni Zou, 591
540 and Hongkai Xiong. Variance reduced domain ran- 592
541 domization for reinforcement learning with policy 593
542 gradient. *IEEE Transactions on Pattern Analysis and 594*
543 *Machine Intelligence (TPAMI)*, 2024. 2 595
- 544 [15] M. Kayaalp, S. Vlaski, and A. Sayed. Dif-maml: 596
545 Decentralized multi-agent meta-learning. *IEEE Open 597*
546 *Journal of Signal Processing*, 3:71–93, 2022. 2 598
- 547 [16] Lars Kunze, Tobias Roehm, and Michael Beetz. To- 599
548 wards semantic robot description languages. In *2011 600*
549 *IEEE International Conference on Robotics and Au- 601*
550 *tomation (ICRA)*, pages 5589–5595, 2011. 2 602
- 551 [17] Chenhao Li, Andreas Krause, and Marco Hutter. 603
552 Robotic world model: A neural network simulator for 604
553 robust policy optimization in robotics. In *Proceedings 605*
554 *of The 8th Conference on Robot Learning (CoRL)*, 606
555 2025. 2 607
- [18] Yi Li, Yuquan Deng, Jesse Zhang, Joel Jang, Marius 608
Mommel, Raymond Yu, Caelan Reed Garrett, Fabio 609
Ramos, Dieter Fox, Anqi Li, Abhishek Gupta, and 610
Ankit Goyal. Hamster: Hierarchical action models for 611
open-world robot manipulation. In *International Con- 612*
ference on Learning Representations (ICLR), 2025. 2 613
- [19] Naureen Mahmood, Nima Ghorbani, Nikolaus F. 614
Troje, Gerard Pons-Moll, and Michael J. Black. 615
Amass: Archive of motion capture as surface shapes. 616
In *IEEE/CVF International Conference on Computer 617*
Vision (ICCV), 2019. 4 618
- [20] Mingjie Pan, Jiyao Zhang, Tianshu Wu, Yinghao 619
Zhao, Wenlong Gao, and Hao Dong. Omnimanip: To- 620
wards general robotic manipulation via object-centric 621
interaction primitives as spatial constraints. In *Pro- 622*
ceedings of the IEEE/CVF Conference on Computer 623
Vision and Pattern Recognition (CVPR), 2025. 2 624
- [21] Silvia Saporá, Gokul Swamy, Chris Lu, Yee Whye 625
Teh, and Jakob Nicolaus Foerster. Evil: Evolution 626
strategies for generalisable imitation learning. In *In- 627*
ternational Conference on Machine Learning (ICML), 628
2024. 2 629
- [22] Carmelo Sferrazza, Dun-Ming Huang, Xingyu Lin, 630
Youngwoon Lee, and Pieter Abbeel. Humanoid- 631
bench: Simulated humanoid benchmark for whole- 632
body locomotion and manipulation. *arXiv preprint 633*
arXiv:2403.10506, 2024. 4 634
- [23] Yutong Shen, Hangxu Liu, Lei Zhang, Penghui Liu, 635
Ruizhe Xia, Tianyi Yao, and Tongtong Feng. Detach: 636
Cross-domain learning for long-horizon tasks via mix- 637
ture of disentangled experts, 2025. 2 638
- [24] Kai Tai Song and Hsiang Hsi Chen. Hagrasp: Hy- 639
brid action grasp control in cluttered scenes using 640
deep reinforcement learning. In *2024 IEEE In- 641*
ternational Conference on Robotics and Automation 642
(ICRA), 2024. 2 643
- [25] Zichen Zhang, Yunshuang Li, Osbert Bastani, Ab- 644
hishek Gupta, Dinesh Jayaraman, Yecheng Jason Ma, 645
and Luca Weihs. Universal visual decomposer: Long- 646
horizon manipulation made easy. In *2024 IEEE In- 647*
ternational Conference on Robotics and Automation 648
(ICRA), pages 6973–6980. IEEE, 2024. 2 649
- [26] Xinghao Zhu, Yuxin Chen, Lingfeng Sun, Farzad 650
Niroui, Simon Le Cleac’h, Jiuguang Wang, and Kuan 651
Fang. Relic: Versatile loco-manipulation through flex- 652
ible interlimb coordination. In *Conference on Robot 653*
Learning (CoRL), 2025. 3 654