

EvoWorld: A World-Model-Centric Framework for Continuous Self-Evolution of Modular Embodied Skills

Anonymous CVPR Workshop FMEA submission

Paper ID 42

Abstract

001 *Current progress in physical intelligence largely relies on*
002 *scaling monolithic Vision-Language-Action (VLA) models,*
003 *yet real-world policy data remain fragmented across scenes*
004 *and tasks. This mismatch limits transfer, exacerbates catas-*
005 *trophic forgetting, and impedes continual improvement. A*
006 *modular design that shares dynamics while specializing*
007 *skills is therefore a promising paradigm. We introduce*
008 **EvoWorld (EvoW)**, *a world-model-centric framework for*
009 *skill orchestration and iterative self-evolution. In EvoW,*
010 *VLA form an expandable library of pluggable experts. A*
011 *high-level router selects experts conditioned on scene and*
012 *task, while an action-conditioned video world model serves*
013 *as the cognitive core for rollout-based planning. The world*
014 *model provides counterfactual rollouts to score candidate*
015 *experts, while selected experts execute in the grounded*
016 *scene to generate trajectories for verification. A vision-*
017 *language evaluator delivers semantic scoring and diag-*
018 *nostic tags, enabling targeted updates to the world model,*
019 *router memory, or specific experts rather than global re-*
020 *training. This closes an automated loop that jointly im-*
021 *proves grounding, routing, and skill refinement without*
022 *manual task engineering. Experiments show that EvoW en-*
023 *ables automated task-to-policy synthesis with competitive*
024 *success rates and iterative improvement trends over itera-*
025 *tions, while producing valid and diverse trajectories that*
026 *support evaluation and skill refinement.*

027 1. Introduction

028 Physical intelligence is increasingly required to follow
029 open-ended user intent rather than a fixed list of benchmark
030 tasks [1, 4, 10]. To meet this demand, the field has largely
031 trended toward scaling monolithic Vision-Language-Action
032 (VLA) models and diffusion-style policies in order to cap-
033 ture a broad spectrum of robotic behaviors [5, 7, 25]. Large-
034 scale datasets such as DROID and Open X-Embodiment
035 have further fueled this progress by providing unprece-

036 dented task and visual coverage [8, 24]. Despite this, a po-
037 tential misalignment between universal physical dynamics
038 and highly task-specific execution strategies has been par-
039 tially overlooked in the pursuit of monolithic scaling, leav-
040 ing the challenge of pervasive real-world data silos some-
041 what under-addressed.

042 Current generalist control approaches rely predomi-
043 nantly on training large VLA models or flow-based policies
044 on heterogeneous data to achieve broad manipulation cov-
045 erage [3–5, 7, 25, 28], treating diverse robotic experiences
046 as tokens in a shared backbone and assuming that a single,
047 sufficiently large network generalizes across embodiments
048 via a shared latent space. In practice, however, monolithic
049 methods face limitations under fragmented scene distribu-
050 tions and long-horizon compounding errors [2, 26]: robotic
051 data remain distributed across labs and protocols [8, 24],
052 making continual updates costly and vulnerable to forget-
053 ting; and without explicit feasibility verification, these sys-
054 tems drift semantically in open-world settings [30, 44].

055 We identify the root of these challenges in the lack of
056 decoupling between universal physical dynamics and task-
057 specific execution strategies [15, 18, 38]. The laws of
058 physics are largely invariant across scenarios, while the
059 “skills” required for manipulation are diverse and highly
060 context-dependent [3, 25, 28]. Current paradigms lack an
061 explicit mapping that routes task semantics to a suitable spe-
062 cialist, forcing a single network to resolve competing gradi-
063 ents from disparate domains [40, 41]. We argue that a more
064 rational architecture treats the world model as a universal
065 cognitive foundation while maintaining policies as modu-
066 lar, pluggable skills updated independently [18, 38].

067 **EvoWorld (EvoW)** is proposed as a world-model-
068 centric framework for task-to-policy automation that en-
069 ables continuous skill evolution. Unlike monolithic ap-
070 proaches [3, 20, 25, 28], EvoW separates shared physical
071 grounding from task-level execution and treats VLAs as
072 modular components in an expandable library. The frame-
073 work operates through a principled pipeline: (i) anchoring
074 the task in a high-fidelity physical scene via semantic re-
075 trieval, (ii) performing counterfactual simulation using an

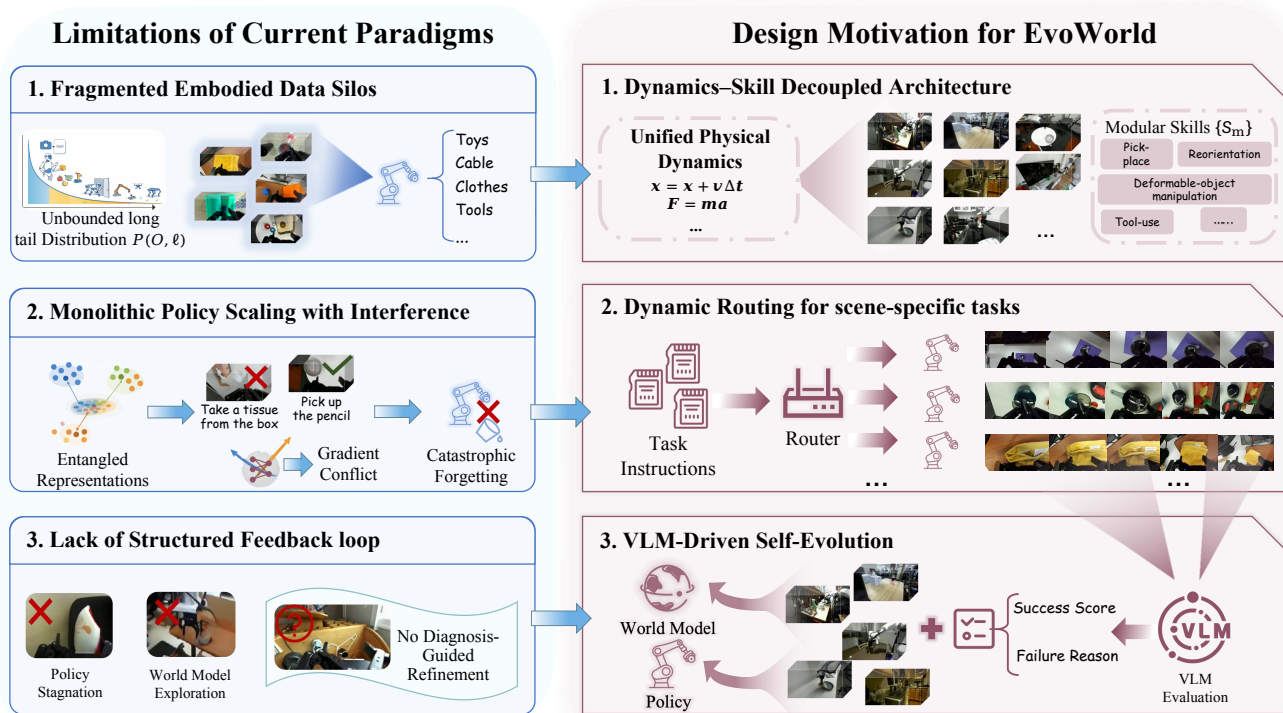


Figure 1. **Motivation and system overview of EvoW.** **Left:** current embodied-learning pipelines face three structural bottlenecks: (1) fragmented embodied data silos with long-tail task/scene distributions, (2) interference and gradient conflict in monolithic policy scaling that induce catastrophic forgetting, and (3) the lack of a diagnosis-guided feedback loop for iterative refinement. **Right:** the proposed framework resolves these bottlenecks via three coupled components: (1) a dynamics–skill decoupled architecture that separates reusable physical dynamics from modular task-specialized skills, (2) scene-aware dynamic routing that maps task instructions and grounded observations to the most suitable expert policy, and (3) a VLM-driven self-evolution loop that scores outcomes, identifies failure causes, and feeds diagnostic signals back to world-model adaptation and skill updates.

076 action-conditioned video world model [30, 44], (iii) routing
 077 the task to the optimal VLA expert through a scene-aware
 078 scheduler, and (iv) driving a closed-loop self-evolution cy-
 079 cle via Vision-Language Model (VLM) diagnosis. This de-
 080 sign allows the system to autonomously refine the world
 081 model and retrain specific skills without re-initializing the
 082 entire intelligence substrate.

083 Experiments demonstrate that EvoW outperforms mono-
 084 lithic baselines [3, 25, 28] in success rates, particularly
 085 across fragmented data silos, and achieves compounding
 086 improvements over iterations through VLM-based failure
 087 diagnosis and expert routing. We summarize our key con-
 088 tributions as follows:

- 089 • **Paradigm:** we propose a principled world-model-centric
 090 formulation that disentangles universal physical dynam-
 091 ics from task-specialized skills, shifting the focus from
 092 policy scaling to world-model-driven orchestration.
- 093 • **Method:** we implement the EvoW framework, which op-
 094 erationalizes this decoupling through retrieval-grounded
 095 scene construction and a VLM-driven diagnostic loop that
 096 enables autonomous, closed-loop self-evolution.

- 097 • **Validation:** extensive experiments across data silos ver-
 098 ify the necessity of expert routing and self-evolution,
 099 demonstrating that EvoW achieves compounding perform-
 100 ance gains without the overhead of full-system retrain-
 101 ing.

2. Related Work 102

103 **Scaling policies and data.** Generalist VLA policies such as
 104 RT-1, RT-2, and OpenVLA demonstrate strong instruction-
 105 following behavior through scale [4, 5, 25], and industrial
 106 foundation efforts (e.g., GR00T N1) push this direction with
 107 larger model and data pipelines [28]. In parallel, large real-
 108 robot datasets including Open X-Embodiment, DROID,
 109 and RH20T broaden scene and task coverage [8, 11, 24].
 110 However, most pipelines remain monolithic, where updates
 111 are largely global and specialist selection at deployment is
 112 weak. Under fragmented data silos, this often leads to in-
 113 terference and higher continual-update cost.

114 **Manipulation policies and structured primitives.** Be-
 115 yond model/data scaling, manipulation research focuses on

control structure, including end-to-end visuomotor learning, reinforcement learning, and imitation learning [23, 27, 35]. Structured formulations such as Transporter Networks and CLIPort improve compositionality and sample efficiency in multi-step tasks via stronger inductive bias for spatial reasoning and action decomposition [34, 42]. Recent dexterity-oriented works further extend policy capacity under richer embodiment and contact settings [43, 45]. Still, these methods are typically optimized within fixed policy-training pipelines and lack a diagnosis-guided, module-specific refinement loop coordinating grounding, routing, and world-model updates.

Grounding and memory for embodied reuse. Grounding and memory-based approaches improve generalization by reusing prior trajectories, grounded plans, and contextual cues [33]. Their core strength is efficient reuse under recurring task patterns, yet in many settings grounding is not tightly coupled with verifiable rollout generation and failure attribution. As a result, they are less suited to autonomous task-to-policy iteration when scenes and task semantics shift together in open-ended real-robot settings.

World models for long-horizon control. World models support imagination-based rollouts for planning and policy improvement [15–17]. Recent latent-action and video-prediction models improve rollout horizon and perceptual realism [12, 13, 39], and other studies investigate tighter coupling between world models and control policies [6, 21, 22]. A core challenge remains long-horizon robustness: compounding prediction errors and semantic drift accumulate over time [2, 26].

Automated task generation and our position. RoboGen shows that closed-loop task generation can scale quickly in simulation [37]. For real-robot deployment, additional constraints are necessary: real-scene grounding, expert routing at deployment, and diagnosis-guided correction. EvoW addresses this gap by combining grounded initialization, task/scene-aware routing, and targeted diagnosis-guided updates in a world-model-centric loop, with explicit decoupling between shared dynamics and task-specialized skills.

3. Method

In this section we delineate the EvoW framework across three strategic dimensions: a formal problem formulation, the modular architectural design, and a diagnosis-driven closed-loop optimization. Our approach emphasizes how Contextual Scene Anchoring, Strategic Routing, and Action-conditioned World-model Rollouts are synergistically coupled to facilitate continuous, autonomous task-to-policy evolution.

3.1. Notations and Problem Setup

We formalize open-ended task-to-policy synthesis as a language-conditioned decision process. Let ℓ denote a natural language task instruction (e.g., “pick up the green object and put it in the bowl”). At each timestep t , o_t represents the visual observation and a_t the executed action. An interaction episode of horizon H is a trajectory $\tau = (o_0, a_0, \dots, o_H)$. To address the limitations of monolithic scaling in fragmented environments, EvoW decomposes this process into universal dynamics and specialized control:

A universal video world model (p_θ): acting as the system’s cognitive engine, this model captures transferable environmental dynamics. The action-conditioned latent transition is

$$z_{t+1} \sim p_\theta(z_{t+1} | z_t, a_t), \quad o_{t+1} = D(z_{t+1}), \quad (1)$$

where E encodes observations into a latent space $z_t = E(o_t)$, and D decodes predicted states back into the observation manifold.

A modular skill library (Π) aggregates task-specific execution strategies as a collection of M specialized VLA experts:

$$\Pi = \{\pi_m\}_{m=1}^M, \quad a_{t:t+H} \sim \pi_m(\cdot | o_t, \ell), \quad (2)$$

where a high-level Strategic Router q maps the grounded scene and instruction to the optimal expert π_{m^*} based on the experience memory \mathcal{M} :

$$m^* = \arg \max_m q(m | o_0, \ell, \mathcal{M}). \quad (3)$$

This world-model-centric decomposition allows physical dynamics to be shared as a universal prior, while policy capacity is allocated through modular specialization and MLLM-driven routing. As illustrated in Figure 2, this architectural decoupling lets EvoW bypass real-world data silos and achieve continuous self-improvement through an introspective diagnostic loop.

3.2. Contextual Scene Anchoring

To bridge the gap between abstract linguistic intent and executable physical behavior, EvoW adopts a Contextual Scene Anchoring paradigm whose objective is to initialize a semantically aligned real-world observation, thereby constraining the subsequent simulation to a physically plausible manifold rather than generating outcomes in a vacuum.

We define an offline experience manifold $\mathcal{D} = \{(o_i, \xi_i)\}_{i=1}^N$, where each entry represents a high-fidelity snapshot of real-robot interaction with observation o_i and auxiliary context ξ_i . Given an open-vocabulary task specification ℓ , we project the language instruction and visual candidates into a shared latent embedding space:

$$e_\ell = \Phi_L(\ell), \quad e_{o_i} = \Phi_O(o_i), \quad (4)$$

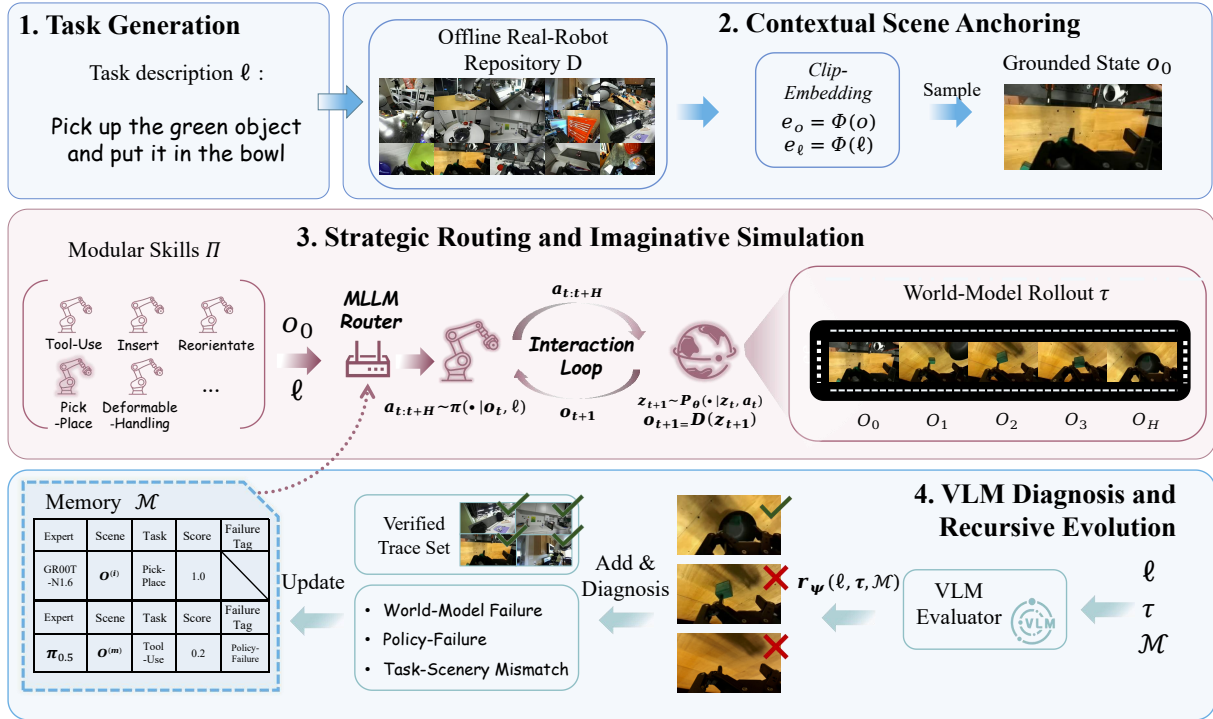


Figure 2. **EvoWorld pipeline overview.** (1) **Task Generation:** construct candidate task specifications for each round. (2) **Contextual Scene Anchoring:** select a semantically aligned real scene from offline experience. (3) **Strategic Routing and Imaginative Simulation:** score candidates via world-model imagination, route to a suitable VLA expert, and execute the selected expert. (4) **VLM Diagnosis and Recursive Evolution:** use success score and failure type to trigger targeted module updates rather than global retraining.

with Φ_L, Φ_O pre-trained dual encoders optimized for cross-modal alignment. To ensure diverse yet relevant initialization, we ground ℓ by sampling an initial state o_0 from a semantic-similarity-based softmax:

$$o_0 \sim p(i | \ell) = \frac{\exp(\langle e_\ell, e_{o_i} \rangle / \tau_g)}{\sum_{j=1}^N \exp(\langle e_\ell, e_{o_j} \rangle / \tau_g)}, \quad (5)$$

where τ_g controls the exploration–exploitation trade-off during grounding.

Once o_0 is anchored, the action-conditioned video world model serves as the primary cognitive engine for counterfactual simulation. Encoding o_0 into $z_0 = E(o_0)$, the system generates a temporal rollout via recursive latent transitions $z_{t+1} \sim p_\theta(z_{t+1} | z_t, a_t)$ and decoding $o_{t+1} = D(z_{t+1})$. This factorization explicitly disentangles scene grounding from dynamics prediction, mitigating the compounding errors and semantic drift that typically plague long-horizon generative models.

3.3. Strategic Routing and Imaginative Simulation

Given the fragmented nature of real-world robotic data, a single monolithic policy often suffers from catastrophic forgetting and performance degradation across diverse embodiments. To circumvent this, EvoW treats each specialized

VLA as a modular expert π_m within a library $\{\pi_m\}_{m=1}^M$, where each expert is optimized for specific task families or physical constraints. The core challenge then shifts to Strategic Routing: mapping an open-vocabulary instruction ℓ and a grounded observation o_0 to the most compatible expert.

We implement the strategic router via a Multi-modal Large Language Model (MLLM) that leverages its semantic knowledge to evaluate the compatibility between task requirements and the functional descriptions of VLA experts. The router is conditioned on a long-term Experience Memory:

$$\pi_{m^*} = q(\ell, o_0, \mathcal{M}), \quad \mathcal{M} = \{(\ell^i, o_0^i, \pi_m^i, y^i, \kappa^i)\}_{i=1}^N, \quad (6)$$

where each entry captures task ℓ , initial state o_0 , selected expert π_m , performance score y , and diagnostic failure type κ .

The routing process operates in a dual-mode logic of exploration and exploitation. **Exploration:** for unfamiliar task–scene pairs with low semantic density in \mathcal{M} , the router encourages diversity by sampling from a broader categorical distribution to discover new “scene–expert” correspondences. **Exploitation:** for recurrent contexts, the MLLM synthesizes historical success/failure cues from \mathcal{M}

Algorithm 1 EvoW closed-loop routing and self-evolution

Require: Task set \mathcal{L} , manifold \mathcal{D} , expert library Π , world model p_θ , router q , evaluator Ψ , memory \mathcal{M} , threshold δ .

Ensure: Verified trajectory set $\widehat{\mathcal{T}}$, updated Π , refined \mathcal{M} .

```

1:  $\widehat{\mathcal{T}} \leftarrow \emptyset$ 
2: for each round do
3:   (S1) sample task  $\ell$ 
4:   (S2) ground  $o_0$  from  $\mathcal{D}$  conditioned on  $\ell$ 
5:   (S3) select  $\pi_{m^*} \leftarrow q(\cdot | \ell, o_0, \mathcal{M})$ ; rollout-score candidates; execute  $\pi_{m^*}$  to obtain  $\tau$ 
6:   (S4)  $(r_\psi, \kappa) \leftarrow \Psi(\tau, \ell, \mathcal{M})$ ; update  $\mathcal{M} \leftarrow \mathcal{M} \cup \{(\ell, o_0, m^*, r_\psi, \kappa)\}$ 
7:   if  $r_\psi \geq \delta$  then
8:      $\widehat{\mathcal{T}} \leftarrow \widehat{\mathcal{T}} \cup \{(\tau, \ell, o_0, m^*, r_\psi)\}$ 
9:   else if  $\kappa = \text{TSM}$  then
10:    update grounding module
11:  else if  $\kappa = \text{PF}$  then
12:    update routing and/or selected expert
13:  else
14:    update world model
15:  end if
16: end for
17: return  $\widehat{\mathcal{T}}, \mathcal{M}$ 

```

into memory-weighted logits, biasing selection toward experts with proven reliability in similar manifolds. To further solidify the decision, for each candidate expert π_m identified by the MLLM, the world model generates an imagined rollout τ_{π_m} . The router computes an imaginative-rollout score that reflects predicted task-consistency, and the final expert π_{m^*} maximizes the joint of semantic compatibility and simulated success, ensuring execution is grounded in both high-level intent and low-level physical feasibility.

3.4. VLM Diagnosis and Recursive Evolution

To achieve continuous self-improvement while avoiding the cost and forgetting of global retraining, EvoW implements an introspective diagnosis mechanism leveraging the cross-modal reasoning of VLMs. This stage performs fine-grained credit assignment, pinpointing the bottleneck within the decoupled architecture. Upon executing a routed expert π_m that yields trajectory τ , the VLM evaluator Ψ contrasts the actual outcome with the task specification ℓ and the initial physical grounding:

$$(r_\psi, \kappa) = \Psi(\tau, \ell, \mathcal{M}), \quad r_\psi \in [0, 1], \quad \kappa \in \{\text{TSM}, \text{PF}, \text{WMF}\}, \quad (7)$$

where r_ψ is a continuous success metric and κ a categorical failure attribution: Task–Scene Mismatch (TSM) flags semantic-grounding errors; Policy Failure (PF) flags suboptimal control from the VLA; World-Model Failure (WMF)

flags divergence between predicted dynamics and physical reality. EvoW then orchestrates targeted updates above a reliability threshold δ :

Verified experience accumulation ($r_\psi \geq \delta$): high-confidence trajectories are integrated into the verified trace set as high-fidelity demonstrations for self-imitation and router calibration.

Semantic grounding refinement ($\kappa = \text{TSM}$): the system updates router memory \mathcal{M} to refine the semantic-to-scene grounding manifold and prevent recurrence of analogous grounding errors.

Modular policy plasticity ($\kappa = \text{PF}$): autonomous retraining is triggered for the specific VLA expert π_{m^*} using the diagnostic feedback as a corrective signal.

Dynamics grounding refinement ($\kappa = \text{WMF}$): when the simulation fails to capture valid physics, world-model parameters are updated using the discrepancy between τ and the physical ground truth.

This targeted evolutionary strategy allows EvoW to bypass the one-size-fits-all training bottleneck and achieve continuous compounding improvements through recursive, module-specific plasticity. The full procedure is summarized in Algorithm 1.

4. Experiments

We evaluate EvoW along three dimensions: **(1) skill-level performance** (single-skill SR and iterative policy improvement); **(2) task diversity and scene validity**; and **(3) generative quality** of imagined trajectories.

4.1. Experimental Setup

To isolate the contributions of our method, we keep VLA backbones, world-model data sources, and evaluation splits fixed across compared variants. We sample tasks from skill families and use grounded initial scenes anchored to a shared offline repository. Across variants, we fix the rollout horizon, task split, scene split, and trajectory budget per task. We first compare base policies and routed variants under identical initialization, then run iterative refinement with matched update budgets per round.

The benchmark covers five representative manipulation families \mathcal{S}_1 – \mathcal{S}_5 : \mathcal{S}_1 (Pick-and-Place), \mathcal{S}_2 (Tool-Use), \mathcal{S}_3 (Reorientation), \mathcal{S}_4 (Deformable-Object Manipulation), and \mathcal{S}_5 (Articulation).

4.2. Skill-Level Performance

We first study whether routing alone enhances policy success under fixed model and data constraints, then analyze cumulative benefits from iterative refinement. Table 1 reports static SR across four VLA backbones and the five skill families. Routing with EvoW consistently improves success rates: averaged over backbones, SR rises from 0.259 to

Table 1. Policy success rates across five skill families and four VLA backbones. **EvoW*** denotes our routed variant; Gain($\uparrow\%$) is relative to the corresponding base average.

Backbone	Setting	\mathcal{S}_1	\mathcal{S}_2	\mathcal{S}_3	\mathcal{S}_4	\mathcal{S}_5	Avg	Gain($\uparrow\%$)
GR00T-N1.6	Base	0.38	0.24	0.18	0.28	0.20	0.256	–
	EvoW*	0.48	0.50	0.46	0.42	0.48	0.468	$\uparrow 82.8\%$
PI-0.5	Base	0.32	0.14	0.24	0.18	0.28	0.232	–
	EvoW*	0.48	0.28	0.38	0.44	0.48	0.412	$\uparrow 77.6\%$
PI-0	Base	0.40	0.22	0.26	0.32	0.22	0.284	–
	EvoW*	0.62	0.38	0.40	0.40	0.50	0.460	$\uparrow 62.0\%$
PI-0-fast	Base	0.38	0.18	0.26	0.24	0.26	0.264	–
	EvoW*	0.66	0.38	0.44	0.36	0.42	0.452	$\uparrow 71.2\%$

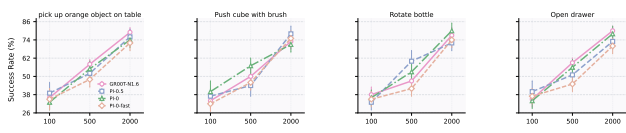


Figure 3. Task-level refinement trends. Iterative diagnosis-driven updates raise SR; variance bars indicate stable gains across runs.

330 0.448 (absolute +18.9 points, relative +73.0%), with per-
 331 backbone relative gains spanning +62.0% to +82.8%, indicating the gain is stable rather than concentrated on a single
 332 backbone.
 333

334 Figure 3 reports SR over refinement iterations on repre-
 335 sentative tasks. SR improves across tasks, indicating that
 336 diagnosis-conditioned updates in Stage 4 reduce recurrent
 337 failure patterns over time. These trends are consistent with
 338 the formulation: expert selection contributes immediate
 339 gains under fixed backbones, while diagnosis-conditioned
 340 refinement yields additional improvements over successive
 341 rounds.

342 4.3. Task Diversity, Scene Validity, and Generative 343 Quality

344 Beyond execution success, EvoW must scale task genera-
 345 tion while preserving scene validity, and its imagined roll-
 346 outs must remain visually faithful and temporally consist-
 347 ent.

348 **Diversity vs. executability.** Compared with simulation
 349 generators including RoboGen [37], RL Bench, MetaWorld,
 350 ManiSkill2 and GenSim, EvoW achieves the lowest Self-
 351 BLEU [29] and CLIP image similarity [31], with competi-
 352 tive ViT similarity [9]. EvoW is weaker on Sentence-BERT
 353 similarity [32], which is expected because grounded exe-
 354 cutability constraints reduce free-form language variation.
 355 This trade-off is acceptable: tasks are intended for phys-
 356 ically realizable rollouts rather than unconstrained textual
 357 diversity.

358 **Generative quality.** Against world-model baselines

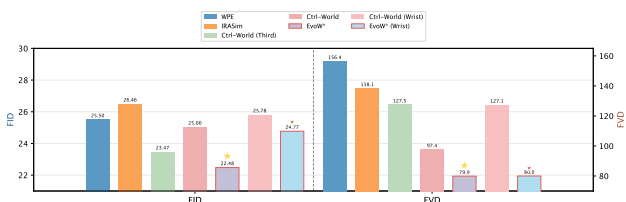


Figure 4. Generative-quality comparison under a view-conditioned setting. EvoW achieves lower FID and FVD than the compared baselines; the larger reduction on FVD is consistent with improved temporal coherence.

(WorldGym, IRASim, Ctrl-World) [14, 30, 44], we report
 FID [19] and FVD [36]. Figure 4 shows that EvoW achieves
 lower FID and FVD than the compared baselines; the larger
 relative gain on FVD indicates stronger temporal consist-
 ency over long-horizon rollouts. Under view-conditioned
 settings, EvoW remains better than Ctrl-World for both
 third-view and wrist-view comparisons, consistent with im-
 proved rollout stability and better preservation of object
 identity and contact progression across frames.

368 5. Conclusion

369 We presented **EvoW**, a framework that maps open-ended
 370 task semantics to verified training trajectories and exe-
 371 cutable manipulation policies. EvoW grounds each task
 372 in task-consistent real scenes, synthesizes scene-consistent
 373 episodes with a decoupled world model, and closes the loop
 374 via VLM-based routing and verification. This design de-
 375 couples shared dynamics from task-specialized skills and
 376 enables targeted updates through failure diagnosis. Scene
 377 grounding and automated VLM feedback are complemen-
 378 tary for scaling open-ended manipulation, yielding reli-
 379 able execution and continual improvement over iterations.
 380 Future work includes (i) expanding the anchor repository,
 381 (ii) uncertainty-aware long-horizon rollouts and filtering,
 382 and (iii) calibrated routing/verification with task-aware crit-
 383 ics.

384

References

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

[1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do as i can, not as i say: Grounding language in robotic affordances. In *CORL*, 2022.

[2] Kavosh Asadi, Dipendra Misra, Seungchan Kim, and Michel L. Littman. Combating the compounding-error problem with a multi-step model, 2019.

[3] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky. π_0 : A vision-language-action flow model for general robot control, 2026.

[4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspier Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *CORL*, 2023.

[5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspier Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale. In *RSS*, 2023.

[6] Jun Cen, Chaohui Yu, Hangjie Yuan, Yuming Jiang, Siteng Huang, Jiayan Guo, Xin Li, Yibing Song, Hao Luo, Fan Wang, Deli Zhao, and Hao Chen. Worldvla: Towards autoregressive action world model, 2025. 441
442
443
444

[7] 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. In *RSS*, 2023. 445
446
447
448

[8] Embodiment Collaboration, Abby O’Neill, Abdul Rehman, Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi ”Jim” Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minh Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Muhammad Zubair Irshad, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, 449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498

499	Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu,	gas, David Ha, Honglak Lee, and James Davidson. Learning	557
500	Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter	latent dynamics for planning from pixels. In <i>ICML</i> , 2019.	558
501	Abbeel, Priya Sundareshan, Qiuyu Chen, Quan Vuong, Rafael	[16] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Moham-	559
502	Rafailov, Ran Tian, Ria Doshi, Roberto Martín-Martín, Ro-	mad Norouzi. Dream to control: Learning behaviors by la-	560
503	han Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Run-	latent imagination. In <i>ICLR</i> , 2020.	561
504	jia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah,	[17] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and	562
505	Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kir-	Jimmy Ba. Mastering atari with discrete world models. In	563
506	mani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl,	<i>ICLR</i> , 2021.	564
507	Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shu-	[18] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy	565
508	ran Song, Sichun Xu, Siddhant Haldar, Siddharth Karam-	Lillicrap. Mastering diverse domains through world models,	566
509	cheti, Simeon Adebola, Simon Guist, Soroush Nasiriany,	2024.	567
510	Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian	[19] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner,	568
511	Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae	Bernhard Nessler, and Sepp Hochreiter. Gans trained by a	569
512	Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tan-	two time-scale update rule converge to a local nash equilib-	570
513	may Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao,	rium. In <i>NeurIPS</i> , 2017.	571
514	Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev,	[20] Physical Intelligence, Kevin Black, Noah Brown, James	572
515	Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity	Darpanian, Karan Dhabalia, Danny Driess, Adnan Esmail,	573
516	Chung, Vidhi Jain, Vikash Kumar, Vincent Vanhoucke, Vi-	Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Gal-	574
517	tor Guizilini, Wei Zhan, Wenxuan Zhou, Wolfram Burgard,	liker, Dibya Ghosh, Lachy Groom, Karol Hausman, Brian	575
518	Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu,	Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin	576
519	Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yan-	LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri,	577
520	song Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen	Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi,	578
521	Chebatar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yix-	Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz,	579
522	uan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho,	James Tanner, Quan Vuong, Homer Walke, Anna Walling,	580
523	Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin	Haohuan Wang, Lili Yu, and Ury Zhilinsky. $\pi_{0.5}$: a	581
524	Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang	vision-language-action model with open-world generaliza-	582
525	Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma,	tion, 2025.	583
526	Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and	[21] Joel Jang, Seonghyeon Ye, Zongyu Lin, Jiannan Xiang, Jo-	584
527	Zipeng Lin. Open x-embodiment: Robotic learning datasets	han Bjorck, Yu Fang, Fengyuan Hu, Spencer Huang, Kaushil	585
528	and rt-x models, 2025.	Kundalia, Yen-Chen Lin, Loic Magne, Ajay Mandlkar,	586
529	[9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,	Avnish Narayan, You Liang Tan, Guanzhi Wang, Jing Wang,	587
530	Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,	Qi Wang, Yinzen Xu, Xiaohui Zeng, Kaiyuan Zheng, Rui-	588
531	Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-	jie Zheng, Ming-Yu Liu, Luke Zettlemoyer, Dieter Fox, Jan	589
532	vain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is	Kautz, Scott Reed, Yuke Zhu, and Linxi Fan. Dreamgen: Un-	590
533	worth 16x16 words: Transformers for image recognition at	locking generalization in robot learning through video world	591
534	scale. In <i>ICLR</i> , 2021.	models, 2025.	592
535	[10] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey	[22] Yuxin Jiang, Yuchao Gu, Ivor W. Tsang, and Mike Zheng	593
536	Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid,	Shou. Olaf-world: Orienting latent actions for video world	594
537	Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong	modeling, 2026.	595
538	Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duck-	[23] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz,	596
539	worth, Sergey Levine, Vincent Vanhoucke, Karol Hausman,	Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan	597
540	Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch,	Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, and Sergey	598
541	and Pete Florence. Palm-e: An embodied multimodal lan-	Levine. Qt-opt: Scalable deep reinforcement learning for	599
542	guage model. In <i>ICML</i> , 2023.	vision-based robotic manipulation, 2018.	600
543	[11] Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu,	[24] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Bal-	601
544	Chenxi Wang, Junbo Wang, Haoyi Zhu, and Cewu Lu.	akrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush	602
545	Rh20t: A comprehensive robotic dataset for learning diverse	Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang	603
546	skills in one-shot, 2023.	Chen, Kirsty Ellis, Peter David Fagan, Joey Hejna, Masha	604
547	[12] Shenyan Gao, Siyuan Zhou, Yilun Du, Jun Zhang, and	Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree	605
548	Chuang Gan. Adaworld: Learning adaptable world models	Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha,	606
549	with latent actions. In <i>ICML</i> , 2025.	Arhan Jain, Abraham Lee, Youngwoon Lee, Marius Mem-	607
550	[13] Quentin Garrido, Tushar Nagarajan, Basile Terver, Nicolas	mell, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Al-	608
551	Ballas, Yann LeCun, and Michael Rabbat. Learning latent	bert Zhan, Kevin Black, Cheng Chi, Kyle Beltran Hatch,	609
552	action world models in the wild, 2026.	Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pan-	610
553	[14] Yanjiang Guo, Lucy Xiaoyang Shi, Jianyu Chen, and	nag R Sanketi, Archit Sharma, Cody Simpson, Quan Vuong,	611
554	Chelsea Finn. Ctrl-world: A controllable generative world	Homer Rich Walke, Blake Wulfe, Ted Xiao, Jonathan Hee-	612
555	model for robot manipulation. In <i>ICLR</i> , 2026.	won Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia,	613
556	[15] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Ville-		

- 614 Rohan Bajjal, Mateo Guaman Castro, Daphne Chen, Qi-
615 yu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul Fos-
616 ter, Jensen Gao, Vitor Guizilini, David Antonio Herrera,
617 Minh Heo, Kyle Hsu, Jiaheng Hu, Muhammad Zubair Ir-
618 shad, Donovan Jackson, Charlotte Le, Yunshuang Li, Kevin
619 Lin, Roy Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mir-
620 chandani, Daniel Morton, Tony Nguyen, Abigail O'Neill,
621 Rosario Scalise, Derick Seale, Victor Son, Stephen Tian,
622 Emi Tran, Andrew E. Wang, Yilin Wu, Annie Xie, Jingyun
623 Yang, Patrick Yin, Yunchu Zhang, Osbert Bastani, Glen
624 Berseth, Jeannette Bohg, Ken Goldberg, Abhinav Gupta,
625 Abhishek Gupta, Dinesh Jayaraman, Joseph J Lim, Jitendra
626 Malik, Roberto Martín-Martín, Subramanian Ramamoorthy,
627 Dorsa Sadigh, Shuran Song, Jiajun Wu, Michael C. Yip,
628 Yuke Zhu, Thomas Kollar, Sergey Levine, and Chelsea Finn.
629 Droid: A large-scale in-the-wild robot manipulation dataset.
630 In *RSS*, 2024.
- 631 [25] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao,
632 Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Fos-
633 ter, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kol-
634 lar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey
635 Levine, Percy Liang, and Chelsea Finn. Openvla: An open-
636 source vision-language-action model. In *CORL*, 2024.
- 637 [26] Nathan Lambert, Kristofer Pister, and Roberto Calandra. In-
638 vestigating compounding prediction errors in learned dy-
639 namics models, 2022.
- 640 [27] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter
641 Abbeel. End-to-end training of deep visuomotor policies.
642 *JMLR*, 17(39):1–40, 2016.
- 643 [28] NVIDIA, Johan Bjorck, Fernando Castañeda, Nikita Cher-
644 niadev, Xingye Da, Runyu Ding, Linxi "Jim" Fan, Yu
645 Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, Joel Jang,
646 Zhenyu Jiang, Jan Kautz, Kaushil Kundalia, Lawrence Lao,
647 Zhiqi Li, Zongyu Lin, Kevin Lin, Guilin Liu, Edith Llon-
648 top, Loic Magne, Ajay Mandlekar, Avnish Narayan, Soroush
649 Nasiriany, Scott Reed, You Liang Tan, Guanzhi Wang, Zu
650 Wang, Jing Wang, Qi Wang, Jiannan Xiang, Yuqi Xie,
651 Yinzhen Xu, Zhenjia Xu, Seonghyeon Ye, Zhiding Yu, Ao
652 Zhang, Hao Zhang, Yizhou Zhao, Ruijie Zheng, and Yuke
653 Zhu. Gr00t n1: An open foundation model for generalist
654 humanoid robots, 2025.
- 655 [29] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing
656 Zhu. BLEU: a method for automatic evaluation of machine
657 translation. In *Proceedings of ACL*, 2002.
- 658 [30] Julian Quevedo, Ansh Kumar Sharma, Yixiang Sun, Varad
659 Suryavanshi, Percy Liang, and Sherry Yang. Worldgym:
660 World model as an environment for policy evaluation, 2025.
- 661 [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
662 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
663 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen
664 Krueger, and Ilya Sutskever. Learning transferable visual
665 models from natural language supervision. In *ICML*, 2021.
- 666 [32] Nils Reimers and Iryna Gurevych. Sentence-BERT: Sen-
667 tence embeddings using siamese BERT-networks. In *Conf.*
668 *Empirical Methods Natural Lang. Process.*, 2019.
- 669 [33] Gabriel Sarch, Yue Wu, Michael J. Tarr, and Katerina
670 Fragkiadaki. Open-ended instructable embodied agents with
memory-augmented large language models. In *EMNLP*,
2023. 671
- [34] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport:
What and where pathways for robotic manipulation. In
CORL, 2021. 672
- [35] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral
cloning from observation, 2018. 673
- [36] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach,
Raphael Marinier, Marcin Michalski, and Sylvain Gelly. To-
wards accurate generative models of video: A new metric &
challenges, 2018. 674
- [37] Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang,
Yian Wang, Katerina Fragkiadaki, Zackory Erickson, David
Held, and Chuang Gan. Robogen: Towards unleashing infi-
nite data for automated robot learning via generative simula-
tion. In *ICML*, 2024. 675
- [38] Yucen Wang, Fengming Zhang, De-Chuan Zhan, Li Zhao,
Kaixin Wang, and Jiang Bian. Co-evolving latent action
world models, 2025. 676
- [39] Zizhao Wang, Chang Shi, Jiaheng Hu, Kevin Rohling,
Roberto Martín-Martín, Amy Zhang, and Peter Stone. Fac-
tored latent action world models, 2026. 677
- [40] Shihan Wu, Xu Luo, Ji Zhang, Junlin Xie, Jingkuan Song,
Heng Tao Shen, and Lianli Gao. Policy contrastive decoding
for robotic foundation models, 2025. 678
- [41] Jiazhi Yang, Kunyang Lin, Jinwei Li, Wencong Zhang, Tian-
wei Lin, Longyan Wu, Zhizhong Su, Hao Zhao, Ya-Qin
Zhang, Li Chen, Ping Luo, Xiangyu Yue, and Hongyang Li.
Rise: Self-improving robot policy with compositional world
model, 2026. 679
- [42] Andy Zeng, Pete Florence, Jonathan Tompson, Stefan
Welker, Jonathan Chien, Maria Attarian, Travis Armstrong,
Ivan Krasin, Dan Duong, Ayzaan Wahid, Vikas Sindhwani,
and Johnny Lee. Transporter networks: Rearranging the vi-
sual world for robotic manipulation. In *CORL*, 2020. 680
- [43] Tony Z. Zhao, Jonathan Tompson, Danny Driess, Pete Flo-
rence, Kamyar Ghasemipour, Chelsea Finn, and Ayzaan
Wahid. Aloha unleashed: A simple recipe for robot dexterity.
In *CORL*, 2024. 681
- [44] Fangqi Zhu, Hongtao Wu, Song Guo, Yuxiao Liu, Chilam
Cheang, and Tao Kong. Irasim: A fine-grained world model
for robot manipulation, 2024. 682
- [45] Minjie Zhu, Yichen Zhu, Jinming Li, Junjie Wen, Zhiyuan
Xu, Ning Liu, Ran Cheng, Chaomin Shen, Yaxin Peng,
Feifei Feng, and Jian Tang. Scaling diffusion policy in trans-
former to 1 billion parameters for robotic manipulation. In
ICRA, 2025. 683
- 684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717