

EVE: A Generator-Verifier System for Generative Policies

Anonymous Submission

Abstract—Visuomotor policies based on generative architectures such as diffusion and flow-based matching have shown strong performance but degrade under distribution shifts, demonstrating limited recovery capabilities without costly finetuning. In the language modeling domain, test-time compute scaling has revolutionized modern LLMs by leveraging foundation models as zero-shot verification modules to refine candidate solutions. We hypothesize that generative policies can similarly benefit from zero-shot VLM-based verifiers at inference time, a direction that remains relatively underexplored. To this end, we introduce **EVE** — a modular, generator-verifier interaction framework — that boosts the performance of pretrained generative policies at test time, with *no additional training*. **EVE** wraps a frozen base policy with multiple zero-shot, VLM-based verifier agents. Each verifier proposes action refinements, while an action incorporator fuses the aggregated verifier output into base policy denoising to produce the final action. Across a diverse suite of tasks and embodiments, **EVE** consistently improves task success rates without policy or verifier finetuning. Our ablations isolate contribution of verifier capabilities and action incorporator strategies.

I. INTRODUCTION

Foundation models for embodied tasks have demonstrated strong performance across complex tasks. These Vision-Language-Action (VLA) models are typically trained on a large set of manually collected robot demonstrations [1]–[5], using diffusion [6] or flow-matching based [2] generative architectures that capture the inherent multimodality of embodied tasks [7]. Although these models exhibit generalist robot manipulation capabilities, they struggle with slight deviations to operating conditions and do not exhibit strong recovery capabilities when encountering out-of-distribution states during deployment [8]–[11]. Improving the performance or robustness of such policies is typically done by finetuning or retraining with additional in-domain data (or recovery sequences), which is expensive to collect [10], [12], [13] and significantly affects the “generalist” performance of the policies.

To advance the performance and robustness of generative policies without additional *training or finetuning*, we propose to leverage state-of-the-art vision-language models (VLMs) within a unified generator-verifier architecture: **EVE** (Embodied Verifier Ensembles). Scaling test-time compute by leveraging learned reward models (or *verifiers*) has fundamentally redefined the capabilities of foundation LLMs without any additional finetuning or retraining [14]–[18], typically by sampling multiple candidate solutions from the base LLM which are then verified by additional LLMs for correctness. We argue that a similar shift is underway in the embodied domain, wherein frozen generative policies can be improved using *zero-shot verifiers* at test-time deployment — a paradigm

well-suited to robotics, where collecting high-quality real-world data is expensive and laborious [19]–[21]. While recent work has begun leveraging such verifiers for downstream embodied tasks, these approaches require training the verifier module or latent dynamics models *tabula rasa* [22], [23]. In contrast, **EVE** orchestrates multiple zero-shot, VLM-based verifier agents that are focused on distinct capabilities that boost the performance of the frozen base generator policy. Through extensive experimentation, we systematically study various design choices that affect the interaction between the generator policy and verifier modules to improve task performance on a diverse set of embodied manipulation tasks.

Building such systems for embodied policies involves several challenges: the generator and verifier often operate in mismatched modalities (language vs. low-level motor actions); aggregating diverse trajectory-level feedback across multiple verifiers is non-trivial; and it is unclear how to combine aggregated verifier output with base policy action predictions, as naive averaging or verifier overrides are suboptimal. To overcome these challenges, we propose leveraging multiple verifiers with distinct action-feedback strategies and combining the output action corrections with the base policy through an action incorporator.

The key contributions of our work are:

- 1) **EVE**, a **generator-verifier** system tailored for embodied policies, in which verifiers operate with different input modalities, capabilities, and action spaces.
- 2) An **action incorporator** module that employs guided diffusion to fuse aggregated verifier outputs with base policy action predictions.
- 3) We show that zero-shot **EVE** ensembles **outperform state-of-the-art embodied verifier baselines** trained with substantial in-domain data on SimplerEnv, and extend to long-horizon Maniskill-HAB and dual-arm RoboTwin tasks.
- 4) Analysis of how **EVE**’s ensemble **recovers failed roll-outs** and how **verifier ensembles** outperform individual counterparts.

II. RELATED WORK

Test-Time Scaling Through Verification. Recent work has found that spending additional compute during test-time deployment of LLMs can lead to large gains in performance on complex reasoning tasks [17], [24]. Many of these works leverage the generation-verification gap, wherein additional LLMs *verify* the base generator output through learned outcome [14] or process-based reward models [15]. A few recent works have begun leveraging this framework for embodied task

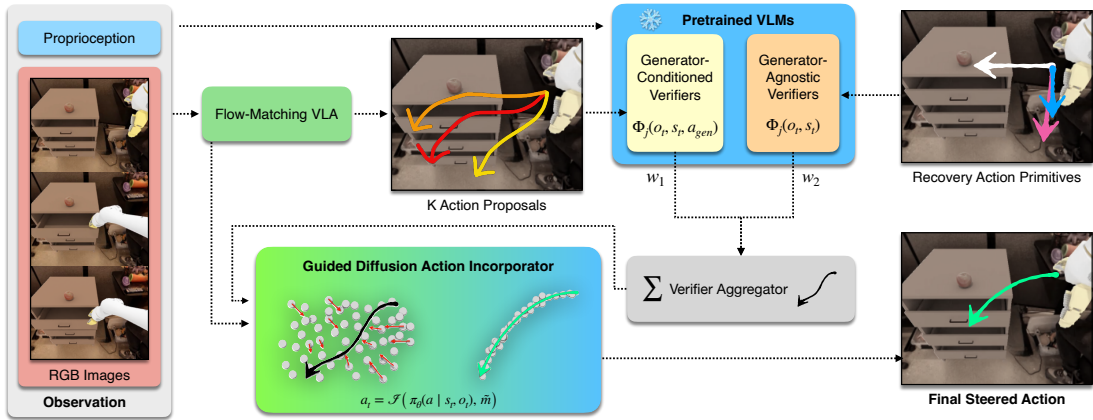


Fig. 1: EVE: A Generator-Verifier Interaction Framework for Generative Embodied Policies. Action feedback from ensemble of verifier agents is incorporated with base policy denoising through diffusion guidance.

performance: RoboMonkey [22] trains a reward model from scratch using synthetically-mined action preferences from a large-scale robotics dataset to score VLA action predictions, and HAVE [25] trains a history-conditioned verifier to score outputs from a diffusion-based generator policy. MAV [26] proposes a system of heterogeneous verifiers but constrains its study to text-only LLMs on mathematical and factual questions. In contrast, EVE constructs an ensemble of *zero-shot* verifiers with distinct capabilities and additionally introduces a systematic orchestration of these verifiers with the policy through an action incorporator that interpolates between base policy action generations and verifier feedback.

Reasoning/Steering to Improve Embodied Policies. Recent work trains policies with reasoning capabilities to improve generalization across diverse task settings [27]–[29], but these typically require fully finetuning the VLA policy on reasoning traces and do not leverage additional inference-time compute. Another line of work *steers* embodied policies towards desired objectives during task execution, typically by learning a latent dynamics model that simulates future states and computes alignment with the task goal; this error signal is then used in classifier-guidance style steering [30]–[32] or post-hoc action-proposal ranking [23], [33]. SAILOR [34] instead discovers recovery sequences within a learned world model and distills them back into the base policy via imitation. Unlike these, which require policy training from scratch or learned world models, we improve policy performance through a *generator-verification* framework comprising multiple *zero-shot verifiers*.

III. METHODOLOGY

We propose EVE a unified framework that augments pre-trained generative policies with modular verifier agents to improve action quality through multi-agent output aggregation and action incorporation.

A. Base Policy Candidate Generation

At each timestep t , given an instruction x , observation o_t , proprioceptive state s_t , and a frozen base policy π_θ , we generate K candidate actions $a_{gen} = \{a_t^{(k)}\}_{k=1}^K$ with $a_t^{(k)} \sim \pi_\theta(o_t, s_t)$. For a frozen diffusion policy, these correspond to

action sequences denoised from K independent noise samples, representing diverse plausible trajectories conditioned on the current observation and state.

B. Verifier Agents

We define a collection of **Verifier Modules** $\mathcal{V} = \{V_j\}_{j=1}^J$, each following a contract $V_j : \Phi_j(o_t, s_t, a_{gen}) \rightarrow m_j \in \mathcal{M}_j$, where Φ_j is a verifier-specific encoding of input context and candidate actions, and m_j is a message in a structured output space \mathcal{M}_j such as trajectory selections or text-based action corrections. We categorize verifiers based on the information available from the generator policy: **I) Generator-Agnostic:** These verifiers operate only on robot observations, receiving $a_{gen} = \emptyset$; conditioned on the task instruction, they select action sequences that maximally ensure task completion. **II) Generator-Conditioned:** These verifiers take a representation of the candidate action sequences as input alongside raw observations. In this work we use trajectory-based representations [38], though alternate representations are equally applicable [39], [40]. This categorization implicitly defines a *generator-verifier interface* in which each V_j interacts with the policy through its own encoding Φ_j .

Verifier Output Aggregation. Each m_j is first brought into an action trajectory representation: *generator-agnostic* verifiers select from predefined action primitives (each with a corresponding trajectory sequence), while *generator-conditioned* verifiers select directly from the available base policy actions (see App. Section C). An aggregation operator \mathcal{A} then projects outputs into a common semantic space, yielding $\tilde{m} = \mathcal{A}(\{m_j\}_{j=1}^J)$, a fused trajectory from weighted interpolation across verifier outputs. In App. Section G, we ablate weighting strategies in \mathcal{A} on downstream success rates.

C. Action Incorporator

We define an *action incorporator* \mathcal{I} that fuses the base policy output $\pi_\theta(a | s_t, o_t)$ with the aggregated verifier trajectory \tilde{m} to produce the executed action a_t . Instead of directly overriding or averaging with the base policy action, the incorporator progressively refines the action through a

Task	π_0	V-GPS	RoboMonkey	EVE		
	[2]	[33]	[22]	Pivot	Primitive	Ensemble
Verifier Training Budget	-	175K demos	20M synthetic	0 (zero-shot)		
WidowX						
Carrot on Plate	56.2	13.2	54.9	<u>56.9</u>	56.2	60.4
Eggplant Basket	91.7	20.1	88.2	91.0	86.1	94.4
Spoon on Towel	84.7	2.1	83.3	<u>87.5</u>	88.9	88.9
Stack Blocks	56.2	4.2	58.3	63.2	<u>61.1</u>	<u>61.1</u>
Average	72.2	9.9	71.2	74.7	73.1	76.2
Google Robot						
Move Near	77.1	77.8	<u>80.6</u>	<u>80.6</u>	79.9	84.7
Put Apple in Drawer	29.9	4.9	35.4	<u>36.1</u>	33.3	40.3
Close Drawer	73.6	68.8	75.7	<u>73.6</u>	70.8	75.7
Average	60.2	50.5	63.9	63.4	61.3	66.9
Total Average	67.1	27.3	68.1	69.8	68.0	72.2

TABLE I: EVE with zero-shot verifiers outperforms finetuned verifiers baselines on SimplerEnv [35]. Verifier training budgets are shown beneath method names. **Bold** is best, underline second-best.

guided denoising process that steers the policy toward verifier-consistent behavior while preserving the original policy prior. In the Guided Diffusion framework, action synthesis is directed by an objective $\xi(\tau, z)$ encoding the alignment between the generated trajectory τ and a verifier-derived feedback signal z . At each diffusion timestep k , the reverse diffusion step is $a_t^{k-1} = \alpha_k(a_t^k - \gamma_k(\epsilon_\theta(o_t, a_t^k, k) + \beta_k \nabla_{a_t^k} \xi(a_t^k, z))) + \sigma_k \eta$, where ϵ_θ is the denoising network, $\eta \sim \mathcal{N}(0, I)$ is Gaussian noise, and $\alpha_k, \gamma_k, \sigma_k$ come from the DDPM noise scheduler. The guidance coefficient β_k controls the influence of the alignment gradient. We use L2-norm discrepancy $\xi(a_t^k, z) = \frac{1}{2} \|a_t^k - z\|_2^2$, whose gradient $a_t^k - z$ biases the reverse diffusion toward verifier-consistent actions while maintaining stability within learned distribution $p(a_t | o_t, s_t)$ of pretrained policy.

D. Intervention Detection

To avoid invoking expensive VLM verifier calls at every step, EVE uses an off-the-shelf failure detector for generative policies [41]: the verifier ensemble is triggered only when the Maximum Mean Discrepancy (MMD) of the base policy’s sampled action distribution exceeds a threshold, flagging erratic rollout segments for intervention (see App. Section H).

IV. EXPERIMENTS

We provide details on implementation and pseudocode in App. Section A and Algorithm 1. We focus on the following research questions: 1) Do ensembles of zero-shot verifiers outperform state-of-the-art embodied verifiers trained with in-domain data? (Section V-A), 2) Can VLM-based zero-shot verifiers in EVE generalize to new embodiments and tasks? (Section V-B), 3) How does verifier-based steering with EVE qualitatively recover failure trajectories? (see Section V-C), and 4) Can EVE generalize to real-world tasks? (Section V-D). We conduct detailed ablations of each component in EVE

Task	DP	DP + EVE	Δ
SetTable-OpenFridge	64.09	66.27	+2.18
SetTable-Place	56.85	58.63	+1.78
SetTable-Pick	25.40	26.19	+0.79
PrepareGroceries-Place	35.02	35.22	+0.20
PrepareGroceries-Pick	11.70	11.28	-0.42
TidyHouse-Pick	16.07	16.17	+0.10

TABLE II: Maniskill-HAB [36]: EVE steers a diffusion policy (DP) towards higher success.

Task	$\pi_{0.5}$	$\pi_{0.5} + \text{EVE}$	Δ
Beat Block Hammer	41.67%	45.49%	+3.82%
Place Can Basket	18.40%	21.88%	+3.48%
Place Container Plate	80.56%	83.68%	+3.12%
Place Object Stand	48.61%	49.31%	+0.70%
Move Can Pot	20.83%	20.14%	-0.69%

TABLE III: RoboTwin [37]: EVE improves $\pi_{0.5}$ performance across dual arm manipulation tasks.

in App. Section G. We present analysis on latency in App. Section M and failures in App. Section K.

V. RESULTS

A. Performance on SimplerEnv

We present the main quantitative results in Table I, evaluating EVE on SimplerEnv [35] across a diverse set of manipulation tasks on two different embodiments and analyze performance gains relative to prior approaches that leverage *trained* verifiers for policy steering.

EVE achieves the highest performance across diverse tasks and embodiments. Across both robot embodiments, EVE improves over the base policy π_0 and prior verifier-based baselines. In particular, EVE-Ensemble attains the best total average success rate of 72.2, outperforming the base π_0 (67.1) and improving upon steering gains through trained verifiers such as RoboMonkey (68.1) and V-GPS (27.3). EVE-Ensemble achieves **76.2** average success on WidowX and **66.9** on Google Robot, indicating robust performance across different platforms and tasks.

EVE outperforms trained verifier baselines. From Table I we see that EVE outperforms recent competitive verifier baselines that require substantial in-domain training budgets. For example, V-GPS [33] is trained with 175K demonstrations, while RoboMonkey [22] relies on 20M synthetic action preference samples for verifier training. More importantly, these verifiers are trained on task and embodiment-specific data which leads to limited generalization capabilities. In contrast, EVE uses zero training data, operating entirely with zero-shot VLM-based verifiers. Despite this dramatic reduction in data requirements, EVE variants achieve higher average performance across the task suite. We hypothesize that these benefits partially stem from EVE’s capability to provide verification feedback at the *action chunk level*. In contrast, the trained verifiers score individual actions on a per-step basis which can lead to poor temporal coherence.

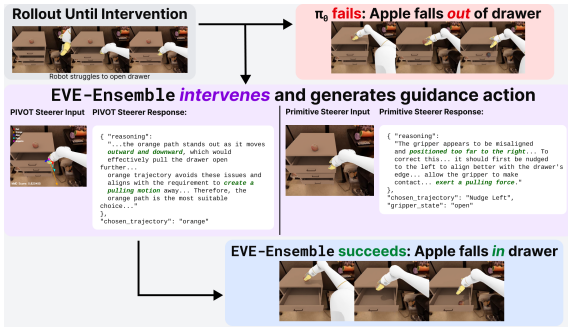


Fig. 2: **Complementary strengths.** The verifiers provide complementary guidance that helps recover a failed trajectory on Place Apple task in SimplerEnv.

Ensembling Improves Individual Verifiers. Even without ensembling, EVE-Pivot achieves 69.8% and EVE-Primitive achieves 68.0% on the total average, indicating that both verifier configurations effectively improve the base policy through our proposed generator-verifier framework. However, ensembling these zero-shot verifiers yields the strongest overall performance, suggesting complementary feedback from each verifier that better captures task progress and impending failures.

Large improvements lower success tasks. Analyzing task-wise performance of EVE-Ensemble, the largest jumps occur on lower-performing tasks, such as *Put Apple* in Drawer ($29.9 \rightarrow 40.3$) and *Stack Blocks* ($56.2 \rightarrow 63.2$). We also observe strong gains on *Move Near* ($77.1 \rightarrow 84.7$) and *Carrot on Plate* ($56.2 \rightarrow 60.4$). However, improvements are modest on near-saturated tasks like *Eggplant Basket* ($91.7 \rightarrow 94.4$).

B. Performance on Maniskill-HAB and Robotwin

For results in this section, we note that state-of-the-art verifiers, such as RoboMonkey [22] and V-GPS [33], are not applicable since they are **not trained** to provide verification feedback for novel tasks or embodiments. We report $\Delta SR = SR_{steered} - SR_{unsteered}$ and employ the EVE-Ensemble configuration.

Maniskill-HAB Results. In this section, we apply EVE to long-horizon mobile manipulation tasks from the Maniskill-HAB benchmark [36]. From Table II, we observe that EVE delivers consistent improvements in performance above the unsteered base policy rollouts by leveraging additional zero-shot verifiers. Specifically, we observe the largest benefits in SetTable-OpenFridge and SetTable-Place with improvements of 2.18% and 1.78% respectively. This empirically proves that EVE is able to boost the performance of the diffusion policy on complex tasks that require recovering from subtle execution degradation. We also present a detailed ablation study on MSHAB in App. Section G.

Robotwin-2.0 Results. To establish further generalization capability of EVE, we provide extended results on a bimanual arm embodiment task suite (using Aloha AgileX). Table III shows that EVE steering leads to performance gains over the $\pi_{0.5}$ VLA on RoboTwin tasks. We apply diffusion guidance to $\pi_{0.5}$ following the technique described in App. Section F2.

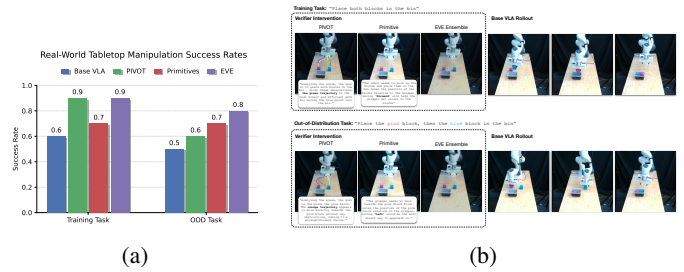


Fig. 3: (a) Each bar represents a success rate out of 10 trials. All EVE variants improve success rates over the base $\pi_{0.5}$. (b) Sample rollouts from real-world experiments demonstrate how EVE can be used to steer VLAs.

C. A Closer Look: How EVE recovers failed trajectories?

In this section, we perform a qualitative analysis of a scenario in which EVE steers a failed episode towards a successful completion. We analyse a failed rollout of Place Apple in Drawer, where the base policy π_0 only opens the drawer partially, leading to misaligned apple placement (Fig. 2). EVE detects a large spike in the MMD value during drawer opening and intervenes: the *Primitive* steerer proposes a “Nudge Left” action to guide the gripper toward the handle, while the *Pivot* steerer selects trajectories that *maximize outward pulling* of the drawer and discard misaligned ones. Averaged feedback from both verifiers guides the final action denoising, enabling the drawer to open fully and the apple to be placed successfully, qualitatively showing how complementary guidance from two different verifiers helps recover a failed trajectory. More qualitative examples in App. Section O.

D. Real Robot Experiments

We validate EVE with real-world robot tabletop manipulation experiments. We collect 50 demonstrations of a Franka Emika Panda arm placing two blocks into a bin. We finetune the $\pi_{0.5}$ VLA [42] and deploy it on the same setup. Our action chunk size is 100 and we query the verifier for the first two action chunks. Fig. 3a depicts the resulting success rates. The out-of-distribution task consists of the same layout of objects with an unseen language prompt. Instead of the training prompt of Place both blocks in the bin, we specify the VLA to Place the pink block, then the blue block in the bin or vice-versa. Without any verifier, the base VLA behavior is random in terms of which block is picked first; the language prompt is ignored. The *Pivot* and *Primitive* verifiers both steer the model to follow the OOD language instructions, with EVE obtaining the highest success rate.

VI. CONCLUSION

In summary, our results show that VLM-based verifier steering enhances policy performance across diverse embodiments and tasks. These findings demonstrate that large VLMs can provide semantically grounded feedback to improve control in open-ended environments when combined with a novel guided diffusion-based action incorporator.

REFERENCES

- [1] J. Bjorck, F. Castañeda, N. Cherniadev, X. Da, R. Ding, L. Fan, Y. Fang, D. Fox, F. Hu, S. Huang *et al.*, “Gr00t n1: An open foundation model for generalist humanoid robots,” *arXiv preprint arXiv:2503.14734*, 2025.
- [2] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, S. Jakobczak, T. Jones, L. Ke, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, L. X. Shi, J. Tanner, Q. Vuong, A. Walling, H. Wang, and U. Zhilinsky, “ π_0 : A vision-language-action flow model for general robot control,” 2024. [Online]. Available: <https://arxiv.org/abs/2410.24164>
- [3] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi *et al.*, “Openvla: An open-source vision-language-action model,” *arXiv preprint arXiv:2406.09246*, 2024.
- [4] J. Barreiros, A. Beaulieu, A. Bhat, R. Cory, E. Cousineau, H. Dai, C.-H. Fang, K. Hashimoto, M. Z. Irshad, M. Itkina *et al.*, “A careful examination of large behavior models for multitask dexterous manipulation,” *arXiv preprint arXiv:2507.05331*, 2025.
- [5] G. R. Team, S. Abeyruwan, J. Ainslie, J.-B. Alayrac, M. G. Arenas, T. Armstrong, A. Balakrishna, R. Baruch, M. Bauza, M. Blokzijl *et al.*, “Gemini robotics: Bringing ai into the physical world,” *arXiv preprint arXiv:2503.20020*, 2025.
- [6] S. Liu, L. Wu, B. Li, H. Tan, H. Chen, Z. Wang, K. Xu, H. Su, and J. Zhu, “Rdt-1b: a diffusion foundation model for bimanual manipulation,” *arXiv preprint arXiv:2410.07864*, 2024.
- [7] J. Urain, A. Mandlekar, Y. Du, M. Shafiqullah, D. Xu, K. Fragkiadaki, G. Chalvatzaki, and J. Peters, “Deep generative models in robotics: A survey on learning from multimodal demonstrations,” *arXiv preprint arXiv:2408.04380*, 2024.
- [8] Q. Gu, Y. Ju, S. Sun, I. Gilitschenski, H. Nishimura, M. Itkina, and F. Shkurti, “Safe: Multitask failure detection for vision-language-action models,” *arXiv preprint arXiv:2506.09937*, 2025.
- [9] S. Zhai, Q. Zhang, T. Zhang, F. Huang, H. Zhang, M. Zhou, S. Zhang, L. Liu, S. Lin, and J. Pang, “A vision-language-action-critic model for robotic real-world reinforcement learning,” *arXiv preprint arXiv:2509.15937*, 2025.
- [10] Y. Yang, Z. Duan, T. Xie, F. Cao, P. Shen, P. Song, P. Jin, G. Sun, S. Xu, Y. You *et al.*, “Fpc-vla: A vision-language-action framework with a supervisor for failure prediction and correction,” *arXiv preprint arXiv:2509.04018*, 2025.
- [11] K. Black, M. Y. Galliker, and S. Levine, “Real-time execution of action chunking flow policies,” *arXiv preprint arXiv:2506.07339*, 2025.
- [12] Y. Dai, J. Lee, N. Fazeli, and J. Chai, “Racer: Rich language-guided failure recovery policies for imitation learning,” in *2025 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2025, pp. 15 657–15 664.
- [13] Z. Lin, J. Duan, H. Fang, D. Fox, R. Krishna, C. Tan, and B. Wen, “Failsafe: Reasoning and recovery from failures in vision-language-action models,” *arXiv preprint arXiv:2510.01642*, 2025.
- [14] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano *et al.*, “Training verifiers to solve math word problems,” *arXiv preprint arXiv:2110.14168*, 2021.
- [15] H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe, “Let’s verify step by step,” in *The Twelfth International Conference on Learning Representations*, 2023.
- [16] J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins, “Solving math word problems with process-and outcome-based feedback,” *arXiv preprint arXiv:2211.14275*, 2022.
- [17] C. Snell, J. Lee, K. Xu, and A. Kumar, “Scaling llm test-time compute optimally can be more effective than scaling model parameters,” *arXiv preprint arXiv:2408.03314*, 2024.
- [18] L. Zhang, A. Hosseini, H. Bansal, M. Kazemi, A. Kumar, and R. Agarwal, “Generative verifiers: Reward modeling as next-token prediction,” *arXiv preprint arXiv:2408.15240*, 2024.
- [19] A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain *et al.*, “Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0,” in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 6892–6903.
- [20] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, M. K. Srirama, L. Y. Chen, K. Ellis *et al.*, “Droid: A large-scale in-the-wild robot manipulation dataset,” *arXiv preprint arXiv:2403.12945*, 2024.
- [21] Q. Bu, J. Cai, L. Chen, X. Cui, Y. Ding, S. Feng, S. Gao, X. He, X. Hu, X. Huang *et al.*, “Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems,” *arXiv preprint arXiv:2503.06669*, 2025.
- [22] J. Kwok, C. Agia, R. Sinha, M. Foutter, S. Li, I. Stoica, A. Mirhoseini, and M. Pavone, “Robomonkey: Scaling test-time sampling and verification for vision-language-action models,” *arXiv preprint arXiv:2506.17811*, 2025.
- [23] Y. Wu, R. Tian, G. Swamy, and A. Bajcsy, “From foresight to forethought: Vlm-in-the-loop policy steering via latent alignment,” *arXiv preprint arXiv:2502.01828*, 2025.
- [24] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, “Chain-of-thought prompting elicits reasoning in large language models,” *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, 2022.
- [25] Y. Li, X. Mao, Y. Yuan, K. Sim, B. Eisner, and D. Held, “Learn from what we have: History-aware verifier that reasons about past interactions online,” *arXiv preprint arXiv:2509.00271*, 2025.
- [26] S. Lifshitz, S. A. McIlraith, and Y. Du, “Multi-agent verification: Scaling test-time compute with multiple verifiers,” *arXiv preprint arXiv:2502.20379*, 2025.
- [27] M. Zawalski, W. Chen, K. Pertsch, O. Mees, C. Finn, and S. Levine, “Robotic control via embodied chain-of-thought reasoning,” *arXiv preprint arXiv:2407.08693*, 2024.
- [28] W. Chen, S. Belkhal, S. Mirchandani, O. Mees, D. Driess, K. Pertsch, and S. Levine, “Training strategies for efficient embodied reasoning,” *arXiv preprint arXiv:2505.08243*, 2025.
- [29] J. Clark, S. Mirchandani, D. Sadigh, and S. Belkhal, “Action-free reasoning for policy generalization,” *arXiv preprint arXiv:2502.03729*, 2025.
- [30] M. Du and S. Song, “Dynaguide: Steering diffusion policies with active dynamic guidance,” *arXiv preprint arXiv:2506.13922*, 2025.
- [31] Z. Sun and S. Song, “Latent policy barrier: Learning robust visuomotor policies by staying in-distribution,” *arXiv preprint arXiv:2508.05941*, 2025.
- [32] Y. Wang, L. Wang, Y. Du, B. Sundaralingam, X. Yang, Y.-W. Chao, C. Pérez-D’Arpino, D. Fox, and J. Shah, “Inference-time policy steering through human interactions,” in *2025 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2025, pp. 15 626–15 633.
- [33] M. Nakamoto, O. Mees, A. Kumar, and S. Levine, “Steering your generalists: Improving robotic foundation models via value guidance,” *arXiv preprint arXiv:2410.13816*, 2024.
- [34] A. K. Jain, V. Mohta, S. Kim, A. Bhardwaj, J. Ren, Y. Feng, S. Choudhury, and G. Swamy, “A smooth sea never made a skilled sailor: Robust imitation via learning to search,” *arXiv preprint arXiv:2506.05294*, 2025.
- [35] X. Li, K. Hsu, J. Gu, K. Pertsch, O. Mees, H. R. Walke, C. Fu, I. Lunawat, I. Sieh, S. Kirmani *et al.*, “Evaluating real-world robot manipulation policies in simulation,” *arXiv preprint arXiv:2405.05941*, 2024.
- [36] A. Shukla, S. Tao, and H. Su, “Maniskill-hab: A benchmark for low-level manipulation in home rearrangement tasks,” *arXiv preprint arXiv:2412.13211*, 2024.
- [37] T. Chen, Z. Chen, B. Chen, Z. Cai, Y. Liu, Z. Li, Q. Liang, X. Lin, Y. Ge, Z. Gu *et al.*, “Robotwin 2.0: A scalable data generator and benchmark with strong domain randomization for robust bimanual robotic manipulation,” *arXiv preprint arXiv:2506.18088*, 2025.
- [38] S. Nasiriany, F. Xia, W. Yu, T. Xiao, J. Liang, I. Dasgupta, A. Xie, D. Driess, A. Wahid, Z. Xu *et al.*, “Pivot: Iterative visual prompting elicits actionable knowledge for vlms,” *arXiv preprint arXiv:2402.07872*, 2024.
- [39] F. Liu, K. Fang, P. Abbeel, and S. Levine, “Moka: Open-world robotic manipulation through mark-based visual prompting,” *arXiv preprint arXiv:2403.03174*, 2024.
- [40] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei, “Voxposer: Composable 3d value maps for robotic manipulation with language models,” *arXiv preprint arXiv:2307.05973*, 2023.
- [41] C. Agia, R. Sinha, J. Yang, Z.-a. Cao, R. Antonova, M. Pavone, and J. Bohg, “Unpacking failure modes of generative policies: Runtime monitoring of consistency and progress,” *arXiv preprint arXiv:2410.04640*, 2024.

- [42] K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail, M. R. Equi, C. Finn, N. Fusai, M. Y. Galliker, D. Ghosh, L. Groom, K. Hausman, brian ichter, S. Jakubczak, T. Jones, L. Ke, D. LeBlanc, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, A. Z. Ren, L. X. Shi, L. Smith, J. T. Springenberg, K. Stachowicz, J. Tanner, Q. Vuong, H. Walke, A. Walling, H. Wang, L. Yu, and U. Zhilinsky, “ π_0 : a vision-language-action model with open-world generalization,” in *9th Annual Conference on Robot Learning*, 2025. [Online]. Available: <https://openreview.net/forum?id=vlhoswksBO>
- [43] W. Kwon, Z. Zhuang, Y. Shen, S. Liang, K. Li, S. Li, D. Wu, X. Lin, M. Stone, S. Moritz *et al.*, “Efficient memory management for large language model serving with pagedattention,” *arXiv preprint arXiv:2309.06180*, 2023.
- [44] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang *et al.*, “Qwen2.5-vl technical report,” 2025.
- [45] A. Ren, “open-pi-zero: Re-implementation of pi0 vision-language-action (vla) model from physical intelligence,” <https://github.com/allenzren/open-pi-zero>, 2025, commit main branch.
- [46] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song, “Diffusion policy: Visuomotor policy learning via action diffusion,” *The International Journal of Robotics Research*, vol. 44, no. 10-11, pp. 1684–1704, 2025.
- [47] Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le, “Flow matching for generative modeling,” *arXiv preprint arXiv:2210.02747*, 2022.
- [48] Y. Zhou, A. Xu, Y. Zhou, J. Singh, J. Gui, and S. Joty, “Variation in verification: Understanding verification dynamics in large language models,” *arXiv preprint arXiv:2509.17995*, 2025.
- [49] Y. J. Ma, J. Hejna, C. Fu, D. Shah, J. Liang, Z. Xu, S. Kirmani, P. Xu, D. Driess, T. Xiao, O. Bastani, D. Jayaraman, W. Yu, T. Zhang, D. Sadigh, and F. Xia, “Vision language models are in-context value learners,” in *The Thirteenth International Conference on Learning Representations*, 2025. [Online]. Available: <https://openreview.net/forum?id=friHA15ofG>

APPENDIX

A. Implementation Details

EVE Verifier Details. For the *Generator-Conditioned* verifier (*Pivot steerer*), we employ the PIVOT [38] prompting strategy with 40 samples from the base policy, drawing on the observation image the 5 most visually distinct trajectories by cosine similarity. For the *Generator-Agnostic* verifier (*Primitive steerer*), we mark the goal location on the RGB image at the intervention point and ask the VLM to suggest a recovery action from a set of action primitives, using only camera views accessible to the base policy (see App. Section D for details). We provide all prompts used in our experiments in App. Section B and details on verifiers in App. Section C. We use vLLM [43] as the inference engine (see App. Section I) and Qwen-2.5-VL-72B [44] as the backbone VLM for all verifiers.

Benchmark and Evaluation Setup. We present results on the SimplerEnv benchmark [35], a real-to-sim tabletop evaluation benchmark with high correlation to real-world success rates, across 7 tasks and two embodiments (4 WidowX, 3 Google Robot). We use π_0 [2] as the base policy (see App. Section F2) and report 48 rollouts across 3 seeds. We also evaluate on the *ManiSkill-HAB* long-horizon mobile manipulation benchmark [36] and *Robotwin-2.0* [37]) bimanual arm task suite. Further details on base policies and task setup are in App. Section D, baselines and evaluation protocols in App. Section E, and hyperparameter settings in App. Section P.

Algorithm 1 EVE: Embodied Verifier Ensemble Inference Pseudocode

Require: Horizon H , Observations $\{o_t\}$, states $\{s_t\}$, Frozen base policy π_θ with N denoising steps, Verifiers $\mathcal{V} = \{V_j\}_{j=1}^J$, MMD threshold τ , MMD computation samples M

- 1: **for** $t = 1$ **to** H **do**
- 2: // Candidate Action Generation
- 3: Sample K base-policy candidates $a_{\text{gen}} = \{a_t^{(k)}\}_{k=1}^K$ using π_θ
- 4: Compute MMD score η_t from overlapping segments using K action samples
- 5: **if** $\eta_t < \tau$ **then**
- 6: Execute nominal action from π_θ and **continue**
- 7: **end if**
- 8: // Verifier Inference
- 9: **for** $j = 1$ **to** J **do**
- 10: Build verifier-specific encoding $\Phi_j(o_t, s_t, a_{\text{gen}})$
- 11: $m_j \leftarrow V_j(\Phi_j(o_t, s_t, a_{\text{gen}}))$ ▷ Selections/corrections over a_{gen}
- 12: **end for**
- 13: $\tilde{m} \leftarrow \mathcal{A}(\{m_j\}_{j=1}^J)$ ▷ Aggregate verifier outputs
- 14: // Guided Diffusion Action Incorporation
- 15: Initialize noisy action sample a_t^N (from DDPM prior)
- 16: **for** $k = N, \dots, 1$ **do**
- 17: Set $z \leftarrow \tilde{m}$ and define
- 18: $\xi(a_t^k, z) = \frac{1}{2} \|a_t^k - z\|_2^2$
- 19: Compute alignment gradient
- 20: $g_k \leftarrow a_t^k - z$
- 21: Denoise with guidance:
- 22: $a_t^{k-1} \leftarrow \alpha_k(a_t^k - \gamma_k(\epsilon_\theta(o_t, a_t^k, k) + \beta_k g_k)) + \sigma_k \eta$
- 23: **end for**
- 24: $a_t \leftarrow a_t^0$ ▷ Final executable control
- 25: Execute a_t
- 26: **end for**

B. Prompts for Verifier Steering

We provide detailed prompts that are used in the Simpler-Env task suites Section P2 and ManiSkill-HAB Section P3.

C. Verifier Ensemble Details

Recall from Section III-B that each verifier module V_j interacts with the generator policy through a verifier-specific encoding

$$V_j : \Phi_j(o_t, s_t, a_{\text{gen}}) \rightarrow m_j \in \mathcal{M}_j, \quad (1)$$

where Φ_j defines the policy-verifier interface, o_t and s_t denote the current observation and proprioceptive state, $a_{\text{gen}} = \{a_t^{(k)}\}_{k=1}^K$ is the set of candidate actions (or trajectories) from the base policy, and m_j is the structured message produced by the verifier (e.g., a trajectory selection or an action primitive).

Below, we instantiate Φ_j for the *Generator-Conditioned Pivot steerer* and the *Generator-Agnostic Primitive steerer* used in our experiments. We also include specific design decisions for both MSHAB and Simpler-Env environments (see App. Section D for task details).

Generator-Conditioned Interface: Pivot Steerer. The *Pivot steerer* is a *Generator-Conditioned* verifier that operates on candidate trajectories produced by the frozen base policy. In all experiments, we use $K = 40$ candidate trajectories sampled from the diffusion policy:

$$a_{\text{gen}} = \{a_{t:t+H}^{(k)}\}_{k=1}^{40}, \quad (2)$$

where each $a_{t:t+H}^{(k)}$ is a horizon- H action sequence produced by the base policy.

The Pivot steerer interface Φ_{pivot} converts $(o_t, s_t, a_{\text{gen}})$ into a compact, diverse set of trajectory visualizations in the robot’s image space, suitable for VLM prompting:

$$\Phi_{\text{pivot}}(o_t, s_t, a_{\text{gen}}) = (x, o_t, \{\hat{\tau}^{(i)}\}_{i=1}^{K_{\text{pivot}}}), \quad (3)$$

where x is the task instruction and $\{\hat{\tau}^{(i)}\}_{i=1}^{K_{\text{pivot}}}$ is a subset of $K_{\text{pivot}} = 5$ visually distinct trajectory overlays in the RGB image frame, derived from the original $K = 40$ samples. Concretely, we proceed as follows:

- 1) **Trajectory decoding in task space.** Each candidate sequence $a_{t:t+H}^{(k)}$ is represented either as (i) joint position deltas (for MSHAB tasks) or (ii) end-effector poses (for Simpler-Env tasks). If the sequence is in joint space, we apply forward kinematics to obtain the corresponding sequence of end-effector poses $\{T_{t+h}^{(k)}\}_{h=0}^H$. If the sequence is already given in end-effector space, we use it directly.
- 2) **Projection into the RGB image frame.** For each candidate trajectory, we project the end-effector poses into the camera frame using known intrinsics and extrinsics, obtaining a 2D path in pixel coordinates. We render this path as an overlay (e.g., a polyline and/or waypoints) on top of the current RGB observation o_t , producing a trajectory visualization $\hat{\tau}^{(k)}$ that indicates how the end effector would move in the image plane.
- 3) **Representative Trajectory Selection.** We greedily select $K_{\text{pivot}} = 5$ trajectories that are maximally diverse under cosine distance, yielding the final set $\{\hat{\tau}^{(i)}\}_{i=1}^{K_{\text{pivot}}}$.
- 4) **Prompt construction.** The interface consists of: (i) the textual task instruction x , (ii) the current RGB frame o_t , and (iii) the K_{pivot} selected trajectory overlays $\{\hat{\tau}^{(i)}\}$. These elements are serialized into a multi-modal prompt to the VLM, which is asked to select the trajectory that best completes the task.

Given this interface, the Pivot steerer verifier V_{pivot} produces a message

$$m_{\text{pivot}} \in \mathcal{M}_{\text{pivot}}, \quad (4)$$

which we instantiate as a discrete selection over the K_{pivot} candidates (e.g., an index of the preferred trajectory) together with a natural language rationale. The candidate action sequence corresponding to the selected trajectory is then used to steer the base policy (lines 17 – 22 of Algorithm 1).

Generator-Agnostic Interface: Primitive Steerer. The Primitive steerer is a Generator-Agnostic verifier and therefore does not consume generator proposals ($a_{\text{gen}} = \emptyset$ in Eq. (1)). Instead, it directly reasons over the current observation and task instruction to select a recovery primitive from a set of predefined ones. Its interface is given by

$$\Phi_{\text{prim}}(o_t, s_t, \emptyset) = (x, \hat{o}_t, \mathcal{A}_{\text{prim}}), \quad (5)$$

where $\mathcal{A}_{\text{prim}}$ denotes the discrete set of action primitives and \hat{o}_t is an augmented visual observation encoding the task goal.

Concretely, we construct Φ_{prim} as follows:

- 1) **Goal marking.** Using the task specification in MSHAB, we extract the goal location of the target object in image

coordinates and overlay this location on the current RGB observation o_t (e.g., by drawing a marker or highlight). The resulting goal-annotated image is denoted \hat{o}_t . For the Simpler-Env tasks, the task instruction is completely described only through language so we do not require goal marking.

- 2) **Primitive set specification.** We define a fixed vocabulary of action primitives $\mathcal{A}_{\text{prim}}$ (e.g., discrete end-effector motions and gripper commands). For the MSHAB tasks, we include primitives that allow for base movement and gripper action. For the Simpler-Env tasks, we include “nudge” primitives which move the end-effector by pre-defined amount in specific directions.
- 3) **Prompt construction.** The interface output consists of: (i) the textual task instruction x , (ii) the goal-marked image \hat{o}_t (or just current image), and (iii) a textual description of the available primitives $\mathcal{A}_{\text{prim}}$. These are serialized into a multi-modal prompt that asks the VLM to choose the most appropriate primitive that provides recovery.

The Primitive steerer verifier V_{prim} then returns a message

$$m_{\text{prim}} \in \mathcal{M}_{\text{prim}}, \quad (6)$$

which we instantiate as a single selected primitive from $\mathcal{A}_{\text{prim}}$ (and optionally a natural language explanation). This primitive is mapped to its low-level control command and used for recovery steering using the guided diffusion incorporator in EVE (see Section III-C).

D. Task Setup and Policy Details

In this section, we provide the task setup details for the ManiSkill-HAB task suite [36] and the SimplerEnv tasks [35].

1) *SimplerEnv Details:* We use subtasks from the SimplerEnv benchmark [35] which has shown strong correlation between simulator and real world evaluations. Specifically we conduct evaluations on the following tasks across two embodiments:

- **Close Drawer:** The robot is spawned in front of a cabinet which has multiple articulated drawers (top/middle/bottom). The robot is tasked to push and close a specific drawer in the cabinet. The robot can be spawned over 9 unique location around the cabinet and can be tasked to close one of the 3 drawers.
- **Move Object:** The robot is tasked to pick an object and place it near another specified target object. Each trial spawns 3 objects in a triangular arrangement placed on the cabinet tabletop.
- **Place Apple in Closed Top Drawer:** The Google robot is spawned in front of a cabinet with an apple on the countertop. The robot is tasked to open the initially closed top drawer, pick up the apple and place it inside the drawer.
- **Carrot on Plate:** The WidowX robot is spawned in a tabletop environment containing a carrot and a plate. The robot is tasked to grasp the carrot and place it onto the plate.

- **Put Eggplant in Basket:** The WidowX robot is spawned in a tabletop setting and is tasked to grasp an eggplant object and accurately place it into the target basket.
- **Spoon on Towel:** The WidowX robot is tasked to pick up a spoon from the workspace and place it onto a towel.
- **Stack Cube:** The WidowX robot is presented with two cubes on a tabletop. The robot is tasked to grasp a specific cube and stack it stably on top of another target cube.

We leverage the evaluation codebase provided in the open-pi-zero repository [45]. We use π_0 for all experiments (see App. Section F2 for details) and use the provided checkpoints in the codebase. For all experiments, we run 48 rollouts spread over 3 random seeds and ensure that the unsteered and steered runs use the exact same random seeding and episode configurations.

2) *ManiSkill-HAB Details:* We refer to the robot end-effector as ee , and its rest position as r . The end-effector resting position is $r = (0.5 \text{ m}, 0 \text{ m}, 1.25 \text{ m})$ relative to the *base*. Let q_{arm} be the arm joint positions, r_{arm} the arm resting joint positions, and \dot{q}_{arm} the arm joint velocities. Similarly, for the torso we define q_{tor} , r_{tor} , and \dot{q}_{tor} . Let v_{base} be the base linear velocity in m s^{-1} (with components $v_{\text{base},x}, v_{\text{base},y}$) and ω_{base} the base angular velocity in rad s^{-1} . We initialize the robot at $(r_{\text{pos}}, r_{\text{arm}}, r_{\text{tor}})$ with $\dot{q}_{\text{arm}} = 0$, $\dot{q}_{\text{tor}} = 0$, $v_{\text{base}} = 0$, and $\omega_{\text{base}} = 0$, and then add clipped Gaussian noise:

$$\begin{aligned} q_{\text{arm}} &\leftarrow q_{\text{arm}} + \text{clip}(\mathcal{N}(0, 0.1), -0.2, 0.2), \\ p_{\text{base}} &\leftarrow p_{\text{base}} + \text{clip}(\mathcal{N}(0, 0.1), -0.2, 0.2), \\ \theta_{\text{base}} &\leftarrow \theta_{\text{base}} + \text{clip}(\mathcal{N}(0, 0.25), -0.5, 0.5). \end{aligned}$$

The z -axis is ‘‘up’’ in ManiSkill3.

Subtask definitions. We use the following shorthand:

$$\begin{aligned} d_t^a &= \|a_{\text{pos}} - b_{\text{pos}}\|_2 \\ &\text{(distance between } ee \text{ and its rest position),} \\ j_k &= \max_{1 \leq i \leq |q_k|} |q_{k,i} - r_{k,i}| \\ &\text{(max deviation from rest for joint group } k\text{).} \end{aligned}$$

We also define $C_{[0:t]}$ to be the sum of cumulative collisions from time 0 to t in N. Below, we provide the ManiSkill task definitions, success and failure criteria:

Task A: Pick[a , optional] (x_{pose})

Description. Pick object x from articulation a (if provided).

Initialization. Spawn robot facing x , within 2 m of x , with noise, and without collisions.

Success.

$$\begin{aligned} 1_{\text{grasped}(x)} \wedge d_{ee}^r \leq 0.05 \wedge j_{\text{arm}} \leq 0.6 \\ \wedge 1_{\text{is_static}} \wedge C_{[0:t]} \leq 5000 \end{aligned}$$

Failure.

$$C_{[0:t]} > 5000 \text{ N.}$$

Task B: Place[a , optional] ($x_{\text{pose}}, g_{\text{pos}}$)

Description. Place object x at goal g (in articulation a , if provided).

Initialization. Spawn with grasp pose sampled from Pick(x_{pose}) policy, robot facing g , within 2 m of g , with noise and without collisions.

Success.

$$\begin{aligned} \neg 1_{\text{grasped}(x)} \wedge d_x^g \leq 0.15 \wedge d_{ee}^r \leq 0.05 \wedge j_{\text{arm}} \leq 0.2 \\ \wedge j_{\text{tor}} \leq 0.01 \wedge 1_{\text{is_static}} \wedge C_{[0:t]} \leq 7500 \end{aligned}$$

Failure.

$$C_{[0:t]} > 7500 \text{ N.}$$

Task C: Open[a] (a_{pos})

Description. Open articulation a with handle at a_{pos} .

Initialization. Spawn the robot facing a . If a is a fridge, sample the base pose uniformly from the region $[0.933, -0.6] \times [1.833, 0.6]$ in front of a ; otherwise use $[0.3, -0.6] \times [1.5, 0.6]$. Add noise and ensure no collisions.

Success. Let a_q , $a_{q_{\text{max}}}$, and $a_{q_{\text{min}}}$ be the current, maximum, and minimum joint positions for the target articulation (drawer or fridge). Define the required opening fraction

$$a_{\text{ofrac}} = \begin{cases} 0.75, & \text{if } a \text{ is a fridge,} \\ 0.9, & \text{otherwise.} \end{cases}$$

We set

$$1_{\text{open}(a)} = 1\{a_q \geq a_{\text{ofrac}}(a_{q_{\text{max}}} - a_{q_{\text{min}}}) + a_{q_{\text{min}}}\},$$

and declare success if

$$\begin{aligned} 1_{\text{open}(a)} \wedge d_{ee}^r \leq 0.05 \\ \wedge j_{\text{arm}} \leq 0.2 \\ \wedge j_{\text{tor}} \leq 0.01 \\ \wedge 1_{\text{is_static}} \\ \wedge C_{[0:t]} \leq 10000 \end{aligned}$$

Failure.

$$C_{[0:t]} > 10000 \text{ N.}$$

Diffusion Policy Baseline. To serve as the base policy, we train diffusion policy (DP) baselines. We use the setup from the MS-HAB paper, with a UNet backbone, a DDPM scheduler, and a 4-layer CNN for visual encoders. For visual observations, the policy relies on two onboard camera views: an egocentric head camera and a wrist-mounted (gripper) camera. For consistency, we use the same architecture and hyperparameters for all subtasks.

Hyperparameter	Value
Learning Rate	0.0001
Batch Size	256
Observation Horizon	2
Action Horizon	1
Prediction Horizon	16
Diffusion Step Embedding Dim	256
UNet Dimensions	[256, 512, 1024]
Number of Groups	8
Number of Training Iterations	500,000

TABLE IV: Diffusion Policy Hyperparameters

Evaluation Setup. We conduct evaluations using pretrained diffusion policy [46] checkpoints. In our analysis, we consider a subset of 6 subtasks where the base diffusion policy has a non-trivial success rate. We report *Success-Once* rates, which computes the percentage of trajectories (out of 1000) that achieve success at least once in an episode with 200 maximum steps. All experiments are reported by running 24 environments in parallel, each with 42 episodes. All rollouts of a particular EVE configuration are done in pairs of back-to-back steered and unsteered runs by ensuring the exact same random seeding.

E. Baselines and Evaluation Protocol

We provided detailed description of baselines and evaluation procedure in this section.

1) *Baselines:* We provide brief description of the primary baselines we compare against on the SimplerEnv benchmark tasks (see Table I). The baselines are defined as follows:

- 1) **RoboMonkey** [22]: Construct a large-scale synthetic action preference dataset by applying augmentations to the BridgeV2 dataset dataset. This dataset is used to finetune a LLaVA-7B scale to assign higher scores to correct samples and lower scores to incorrect ones. During inference, multiple action samples are sampled from the base policy and gaussian perturbations are applied to construct the set of actions that will be scored by the verifier. The action with highest score assigned by verifier is then executed and this process is repeated at each time step.
- 2) **V-GPS** [33]: Trains a value function using offline RL using a mix of BridgeV2 and fractal datasets [19]. During inference, the value function is used to generate scores (or *Q-values*) for multiple actions sampled by the base policy. This is followed by an action re-ranking step where the final action is sampled from a “re-ranked” categorical distribution obtained by computing a temperature-controlled softmax over *Q-values*.

2) *Evaluation Details: SimplerEnv Evals.* For EVE, we finetune the ensemble weights (*Pivot:Primitive*) by performing a search for *Pivot* weights with values from the discrete set: [0.1, 0.2, 0.3, 0.5]. We use the best ensemble weight ratio for each task. For RoboMonkey, we sample 9 action candidates from the base policy and augment to 32 samples using gaussian perturbations. For VGPS, we sample 10 actions from the base policy and use a softmax temperature of 0.1 for action re-ranking. We additionally provide detailed hyperparameters

in Table VIII. These hyperparameters are kept constant across tasks.

MSHAB parameters. We outline hyperparameters used in Maniskill-HAB evaluations in Table IX.

Robotwin parameters. We outline hyperparameters used in Robotwin 2.0 evaluations (see Section F1) in Table VIII.

F. Extended Results and Diffusion Guidance

1) *RoboTwin 2.0 Task Suite Results: Evaluation Details.* For all experiments presented in Section V-B of main text, we use the EVE-Ensemble variant. All evaluated tasks use the exact same hyper-parameters selected from a discrete search over following values: Guidance Scale: [0.5, 1, 1.5, 2, 2.5, 3, 5, 10, 30], MMD Threshold: [0.8, 0.9, 1.0, 1.05, 1.1, 1.2], Ensemble Ratio (*Pivot:Primitive*): [3 : 7, 1 : 1, 7 : 3, 9 : 1]. The final evaluations use the hyperparameters listed in Table VIII. Each task was run over 6 seeds with 48 rollouts per seed.

2) *EVE Steering with Flow-Based VLAs:* In the main paper, we design an action incorporator module that leverages guided diffusion (see Section III-C) to steer diffusion policies, but we can also apply EVE to base policies trained with conditional flow matching [47]. Specifically we present an adaptation of EVE to a large flow-based VLA policy π_0 [2] in the following section and test performance on the SimplerEnv benchmark (see Table I). Additionally we also apply the same technique to the more recent $\pi_{0.5}$ policy and present results on Robotwin 2.0 in Section F1.

Flow-matching policies. Given observation o_t , a flow policy generates an action chunk by first sampling Gaussian noise $A_t^0 \sim \mathcal{N}(0, I)$ and then integrating the learned velocity field v_π over a “flow time” variable $\tau \in [0, 1]$:

$$A_t^{\tau+\Delta} = A_t^\tau + \Delta v_\pi(A_t^\tau, o_t, \tau), \quad \Delta = \frac{1}{n}, \quad (7)$$

where n is the number of denoising steps. The final action is A_t^1 after n Euler steps of Eq. (7).

Guided Inference for Flow Policies. To steer the flow generation towards a specified action sequence A_{ref} (e.g., an action suggested by the EVE verifier system), we use the guided inference scheme proposed in [11]. Let $v_\pi(A_\tau, o, \tau)$ denote the velocity predicted by the policy at flow step τ . We first compute the estimated terminal (clean) action \hat{A}_1 :

$$\hat{A}_1 = A_\tau + (1 - \tau)v_\pi(A_\tau, o, \tau). \quad (8)$$

We define the guidance objective as minimizing the squared error between this estimate and the reference action A_{ref} . We can compute the gradient of the loss $\mathcal{L} = \frac{1}{2} \|\hat{A}_1 - A_{\text{ref}}\|^2$ (similar to line 18 in Algorithm 1):

$$\nabla_v \mathcal{L} = (1 - \tau)(\hat{A}_1 - A_{\text{ref}}). \quad (9)$$

We then adjust the velocity using a guidance scale γ and perform the standard Euler integration step:

$$\hat{v}_\tau = v_\pi(A_\tau, o, \tau) - \gamma \nabla_v \mathcal{L}, \quad (10)$$

$$A_{\tau+\Delta\tau} = A_\tau + \Delta\tau \hat{v}_\tau. \quad (11)$$

Integrating EVE with Flow Policies. When the base policy is a diffusion model, we use Algorithm 1 which performs guided diffusion using the verifier-ensemble action output in lines 15 – 23 via a DDPM-style reverse process. For a flow-based base policy, the outer structure of Algorithm 1 is unchanged: we still draw K candidate action chunks from the frozen base policy (lines 1 – 4), run intervention detection using MMD-based detector (lines 5 – 7), and run the verifier ensemble to obtain an aggregated correction signal \tilde{m} (lines 8 – 13). The only modification is that lines 15 – 23 are replaced by the guided flow integration of Eq. (11), where \hat{v}_τ is computed by setting A_{ref} to \tilde{m} in Eq. (9). Thus EVE can be applied to both diffusion and flow-matching policies with a minimal change to the inference procedure. For all experiments, we use the exact same *Primitive* and *Pivot* steerers as defined in Section A.

G. Detailed Ablations

In this section we present extensive ablations of each module in the EVE framework. For these experiments, we use the `SetTable-Place` task from MSHAB using the `EVE-Ensemble` configuration, unless otherwise stated. In all ablation results, we report the *Delta in Success Rate %* which measures the delta gain in steered and unsteered runs.

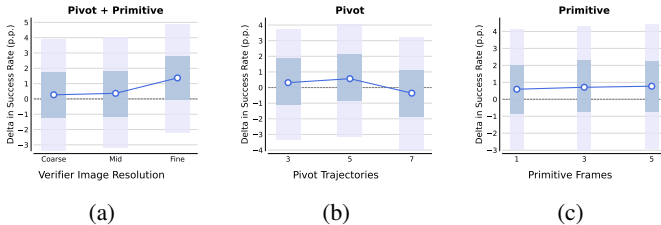


Fig. 4: EVE ablations on the MSHAB task suite (verifier configurations).

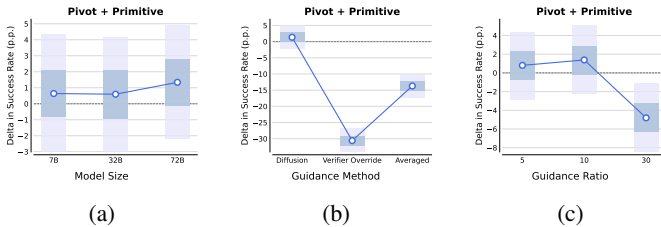


Fig. 5: EVE ablations on the MSHAB task suite (system-level components).

Verifier Model Scaling. We investigate how the size of the underlying VLM affects the verifier’s ability to steer the base policy effectively. We compare the performance impact using Qwen-2.5-VL at different parameter scales: 7B, 32B, and 72B. As illustrated in Fig. 5a, increasing the parameter scale of the verifier model size generally correlates with improved task success rates. Specifically, we find that the 72B model consistently outperforms the smaller variants. This finding is in-line with recent work from the language modelling literature [18], [48] which find increasing benefits with stronger verifier models.

Action Incorporator Design. In Fig. 5b, we ablate various strategies to incorporate verifier-aggregated action feedback into the base policy action predictions. From the results, it is clear that the *Verifier Override* leads to a drastic reduction in performance. This is potentially because the *Primitive* verifier only selects from a predefined list of recovery primitives which are used for guidance over the base action dimensions only (see App Section C for details). Additionally, we observe that direct averaging between verifier-aggregated and base policy actions performs poorly in comparison to the guided diffusion strategy employed in EVE. This is because direct averaging does not ensure that the output action is close to the marginal action distribution of the base policy. This suggests that the action incorporator in EVE integrates “just” the right amount of verifier feedback to prevent task failure but still ensure task completion.

Guidance Ratio Ablation. We conduct ablations with diffusion guidance coefficient β_k (see line 22 in Algorithm 1), which controls amount of verifier feedback incorporated into base policy denoising. In Fig. 5c, we observe that guidance coefficient significantly affects performance of the task with an optimal value of 10 with sharp drops in neighboring values. This result suggests that a very large value of guidance pushes the denoising too far away from the base policy action distribution, causing large temporal inconsistencies leading to reduced task performance.

Verifier Image Resolution. In this experiment, we re-render images at a higher resolution of 256p and 512p from the simulator and pass them to the verifier ensemble. From Fig. 4a, we see that using images with higher visual resolution enhances performance significantly. We note that higher resolution images are important for contact-rich mobile manipulation tasks in MSHAB and enable the verifier to provide fine-grained action feedback.

Verifier Information Ablations. In this experiment, we analyze the density of information that is passed to the individual verifiers through their respective policy-verifier interfaces (see Φ_j defined in Section III). For the *Pivot* verifier, we ablate the number of trajectories that the verifier can select from. From Fig. 4b we observe an increase as the number of drawn trajectories is increased from 3 to 5. But increasing the number of trajectories to 7 leads to degradation in the performance, potentially because the VLM can no longer effectively discern between the trajectories. For the *Primitive* verifier, we ablate the number of history frames that are passed to the VLM. In Fig. 4c, we ablate the number of history frames that are passed to the *Primitive* verifier for recovery primitive selection. From the results, we observe that increasing the frame history doesn’t affect performance significantly. We hypothesize this is because the robot doesn’t have very large movements in adjacent frames once it approaches the receptacle from which it needs to either pick or place an object.

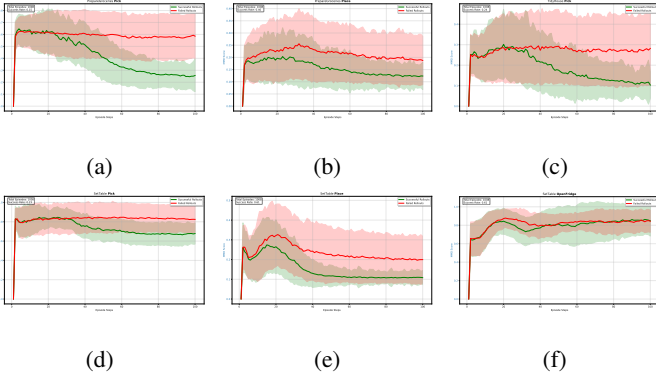


Fig. 6: Base policy MMD scores for successful and failed rollouts on 6 mobile manipulation tasks from the MS-HAB [36] benchmark. Across tasks, we consistently find that MMD scores discriminate between failed and successful rollouts.

H. Intervention Detection using Mean Maximum Discrepancy Scores

To detect distribution shifts and potential erratic behaviors during test-time deployment, we employ the Maximum Mean Discrepancy (MMD) metric. This formulation is based on the STAC metric proposed by Agia et al. [41], which quantifies the distance between the distribution of action trajectories generated at the current timestep t and those generated at the future timestep $t+k$. In Fig. 6 we qualitatively show that MMD scores consistently discern between success and failure rollouts. This section details the mathematical formulation of the temporal consistency check used in Section III-D.

Marginal Action Distributions: We consider a receding horizon control setting where the policy π_θ predicts a sequence of actions (an action chunk) of length H . Let k denote the execution horizon (the number of steps the robot executes before replanning). At timestep t , the policy generates a distribution of trajectories. Following the formulation in [41], we isolate the segment of this trajectory that overlaps with the next planning step $t+k$. The overlapping temporal window has a length of $H-k$.

We define the two marginal distributions over this overlapping window as follows:

- 1) $\bar{\pi}_t$: The distribution of action sequences generated at time t , restricted to the window $[t+k, t+H-1]$. Formally, $\bar{\pi}_t := p(a_{t+k:t+H-1} \mid o_t, s_t)$.
- 2) $\tilde{\pi}_{t+k}$: The distribution of action sequences generated at time $t+k$, restricted to the same window $[t+k, t+H-1]$. Formally, $\tilde{\pi}_{t+k} := p(a_{t+k:t+H-1} \mid o_{t+k}, s_{t+k})$.

Under nominal conditions, the policy’s plan at time t for the future window should remain consistent with the updated plan generated at $t+k$. A high divergence between $\bar{\pi}_t$ and $\tilde{\pi}_{t+k}$ indicates erratic behavior or distribution shift.

Maximum Mean Discrepancy (MMD): We measure the distance between these two distributions using the squared MMD in a Reproducing Kernel Hilbert Space (RKHS) \mathcal{H} associated

with a kernel function $k(\cdot, \cdot)$. The squared population MMD is defined as [41]:

$$D^2(\bar{\pi}_t, \tilde{\pi}_{t+k}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim \bar{\pi}_t} [k(\mathbf{x}, \mathbf{x}')] + \mathbb{E}_{\mathbf{y}, \mathbf{y}' \sim \tilde{\pi}_{t+k}} [k(\mathbf{y}, \mathbf{y}')] - 2\mathbb{E}_{\mathbf{x} \sim \bar{\pi}_t, \mathbf{y} \sim \tilde{\pi}_{t+k}} [k(\mathbf{x}, \mathbf{y})] \quad (12)$$

where \mathbf{x}, \mathbf{y} represent the flattened vectors of the action sequences $a_{t+k:t+H-1}$ sampled from their respective distributions. We employ the Radial Basis Function (RBF) as the kernel function.

Empirical Estimation: Since the analytical densities of the diffusion policy are intractable, we approximate Equation 12 using a finite batch of samples, consistent with the STAC implementation [41]. We draw B samples from the policy at timestep t (denoted as $X = \{\mathbf{x}_i\}_{i=1}^B$) and B samples at timestep $t+k$ (denoted as $Y = \{\mathbf{y}_j\}_{j=1}^B$). The empirical MMD estimate \hat{D} is computed as:

$$\hat{D}(\bar{\pi}_t, \tilde{\pi}_{t+k}) = \frac{1}{B^2} \sum_{i=1}^B \sum_{j=1}^B k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{B^2} \sum_{i=1}^B \sum_{j=1}^B k(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{B^2} \sum_{i=1}^B \sum_{j=1}^B k(\mathbf{x}_i, \mathbf{y}_j) \quad (13)$$

This empirical estimate provides a differentiable and computationally efficient signal of temporal inconsistency. During inference, if the MMD score computed across adjacent timesteps exceeds the threshold τ , we trigger the intervention mechanism described in Section III-D.

I. vLLM Inference Details

We run all VLM-based verifier inferencing using the vLLM library [43]. vLLM provides highly optimized KV cache management using paged attention which enables high throughput handling of concurrent API requests from multiple clients. Specifically, we run each Qwen-2.5-VL-72B instance on a single node of 8 NVIDIA 48GB A40 GPUs. Each vLLM single node server is run with a tensor-parallel rank of 8. We limit gpu memory utilization on each server to 0.85 to prevent crashes due to large burst in number of verification API requests from parallel simulator environments.

J. Task Success and Failure Categories

We use the trajectory categorization system proposed in MSHAB [36]. The authors define various event-based heuristics which enables automated categorization of policy trajectory rollouts into specific success and failure modes based on the chronological sequence of events. In Tables X and XI we provide the qualitative descriptions for the modes associated with the Place and Open subtasks (reproduced from the original MSHAB paper).

K. Failure Analysis and Limitations

Fig. 7 showcases the exact distribution of transitions of failure types to successes through EVE steering on MSHAB. EVE-Ensemble consistently redirects catastrophic failures such as Place-in-goal failure and Excessive Collision, toward stable success modes. This suggests that primary gains through steering arise not just from correcting rare anomalies but from restructuring the policy’s dominant error pathways. However, not all tasks benefit from verifier-based steering through EVE. For example, in PrepareGroceries-Pick, we observe minor drops in performance with the EVE-Ensemble configuration (see Table II). This is potentially due to the task requiring picking objects placed inside refrigerators where effective feedback from verifiers is limited. One potential avenue to improve this is to acquire additional feedback from verifiers, such as progress estimation signals [49], so that the verifier can help disambiguate if the policy is stuck. We provide details of success and failure categories in App. Section J and failure sankey plots for all tasks in App. Section L.

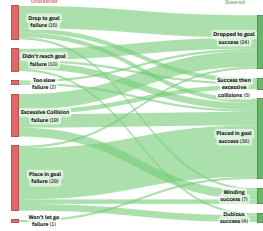


Fig. 7: Sankey plot showing failure episodes switching to successful cases on SetTable-Place.



(a) Failure episodes switching to successful cases on SetTable-OpenFridge. (b) Success to failure episode switches on SetTable-Place. See Fig. 7 for inverted analysis.

Fig. 8: Sankey visualizations of verifier steering effects on SetTable tasks using EVE-Ensemble verifier steering.

Task	Execution Time (s) ↓					Success Rate (%) ↑			Intervention Rate (%) ↓		
	V-GPS	RoboMonkey	EVE			V-GPS	RoboMonkey	EVE (Ens.)	V-GPS	RoboMonkey	EVE
			Prim	Pivot	Ensemble				Prim	Pivot	Ensemble
Move Near	108.8	259.1	195.9	198.0	223.1	77.8	80.6	84.7	100.00	100.00	3.36 2.89 3.07
Put Apple in Drawer	272.0	644.4	<u>558.0</u>	<u>555.3</u>	648.5	4.9	35.4	40.3	100.00	100.00	2.61 2.61 2.95
Close Drawer	155.7	368.5	298.6	<u>295.6</u>	301.0	68.8	75.7	75.7	100.00	100.00	0.23 0.22
Carrot on Plate	48.9	105.6	141.8	136.4	169.1	13.2	54.9	60.4	100.00	100.00	9.48 6.29 9.54
Eggplant Basket	52.6	<u>203.3</u>	257.4	259.1	277.9	20.1	88.2	94.4	100.00	100.00	2.84 2.83 2.99
Spoon on Towel	26.3	<u>108.2</u>	148.9	145.0	182.8	2.1	83.3	88.9	100.00	100.00	11.63 7.87 9.99
Stack Blocks	26.6	<u>107.2</u>	135.5	131.1	147.7	4.2	58.3	61.1	100.00	100.00	4.93 3.05 4.81
Average	98.7	256.7	248.0	<u>245.8</u>	278.6	27.3	68.1	72.2	100.00	100.00	5.01 3.68 4.80

TABLE V: Execution time, success rate, and intervention rate across tasks shared with Table I. **Bold** is fastest and underline is second fastest (execution time). Red highlights show the worst average performance and intervention rates.

L. Failure Sankey Plots

For all success and failure definitions please refer to Section J for details. In this section, we provide additional sankey plots for additional analysis on EVE verifier steering.

Failure-to-Success Transitions. In Fig. 8a we observe that for articulation-focused tasks such as SetTable-OpenFridge, steering nearly eliminates “too slow” and “cant reach” articulation failures, with almost all trajectories ending as open successes. Overall, steering behaves as a robust correction mechanism that enables recovery from catastrophic failure scenarios.

Success-to-Failure Transitions. Fig. 8b analyzes how verifier-based steering perturbs successful SetTable-Place rollouts by re-running unsteered success episodes with steering enabled and categorizing the resulting failures. We observe that high-quality “placed in goal” successes are reclassified as “place in goal failure” or “excessive collision failures” which indicates that slight errors in execution can lead to slowly decaying failures (see Table X), even if the object is correctly placed at the target location initially. This suggests that the verifier systems need to be improved to provide precise feedback which can help prevent such delayed failures in the policy rollout.

M. Latency Analysis

Recent work such as [22], [33] train verifiers using large-scale datasets to score candidate actions at every timestep. Consequently, during inference, these methods require intervening at every step, incurring 100% verification overhead regardless of policy confidence of the current action prediction.

In contrast, EVE leverages a Maximum Mean Discrepancy (MMD) trigger to selectively invoke the VLM-based verifier only when the distribution of sampled actions from the base policy deviates significantly. From Table V, we observe that across all evaluated tasks, our method intervenes on average 4.8% of the steps. While V-GPS achieves the lowest latency among the compared methods, this efficiency comes at the cost of substantially worse task success (see Table V). This suggests that reducing verification overhead alone does not yield reliable performance. We provide a detailed latency and throughput analysis on MSHAB task suite in Section N.

N. Verifier Inference Overhead

Environments	Primitive Steering	Ensemble Throughput
1	0.11	0.04
10	0.26	0.10
20	0.34	0.10
40	0.41	0.12

TABLE VI: Average Throughput (server responses/second) comparison across environments for Primitive Steering and Ensemble methods on a single vllm server.

In this section, we analyze verifier inference overheads in terms of inference latency and system throughput across different configurations and scale in the MSHAB setting. The goal is to demonstrate how the EVE system latency scales with respect to server resources to demonstrate that compute can effectively be shared between parallel policy rollouts simultaneously running EVE steered inference.

vLLM Request Batching Table VI effectively demonstrates that as the number of environments increases, the

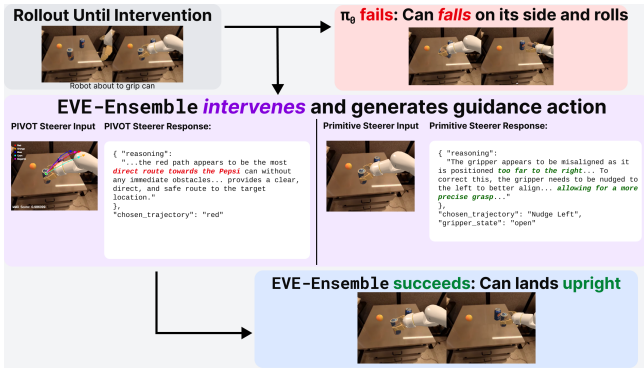


Fig. 9: **Robustness via ensembling.** Aggregating verifier feedback remains corrective even when one steerer is misleading at an intervention point during failed rollout of Move Near task in SimplerEnv.

average latency (average time a robot spends waiting for a vllm server response) for an environment increases, yet, as multiple requests are running simultaneously, throughput increases (0.04 to 0.12 responses per second) by $3\times$. We observe that increasing the batch size (number of parallel environments) effectively improves system throughput. For the Primitive steerer, increasing the environment count from 1 to 40 results in a near $4\times$ increase in throughput (0.11 to 0.41 responses/s). This demonstrates that the underlying vLLM serving infrastructure effectively leverages batching to amortize the cost of large model inference.

System Scalability To mitigate the inference bottleneck during large-scale evaluation, we experiment with using multiple inference servers. For this experiment we fix the total number of episodes to 1008 across 24 envs (42 episodes each) as in the main evaluations. We can now distribute the 24 envs to multiple servers. Table VII illustrates the reduction in total evaluation time as we scale the compute resources from 1 to 6 servers. We observe a strong linear scaling trend. By increasing the server count from 1 to 6, the evaluation time for a fixed set of rollouts drops from ~ 91 minutes to ~ 30 minutes.

Servers	Eval. Time (min)
1	91.68
2	62.88
4	37.20
6	30.35

TABLE VII: Experiment time across different server counts for 1008 rollouts in MSHAB.

O. Qualitative Examples

Robustness in verifier intervention via Ensembling. We analyze another failed rollout on the Move Near task, which requires the robot to grasp the Red Bull can and place it upright near the Coke can. As shown in Fig. 9, the base policy π_0 attempts to grasp the can but executes a faulty grip, causing the can to tip over and roll off the table. In contrast, EVE detects a spike in the MMD value during the *approach phase* and triggers an intervention. At the intervention point,

the *Primitive* steerer proposes a corrective *Nudge Left* action to better align the gripper with the can, improving the grasp configuration. Interestingly, the *Pivot* steerer incorrectly selects a trajectory that does not correspond to meaningful task progress. However, because EVE aggregates feedback across verifiers, averaging the ensemble feedback yields a small leftward correction that leads to achieving a stable grasp. This enables the robot to grasp the Red Bull can correctly and place it upright near the Coke can. Overall, this example illustrates EVE’s robustness to imperfect verifier feedback: even if one steerer is misleading at an intervention point, the aggregated ensemble signal remains corrective and recovers the trajectory.

Maniskill-HAB steering example. Fig. 10 shows an example of steering with EVE on an instance of the SetTable-Place task in the MSHAB benchmark.

P. Hyperparameter Information and Prompts

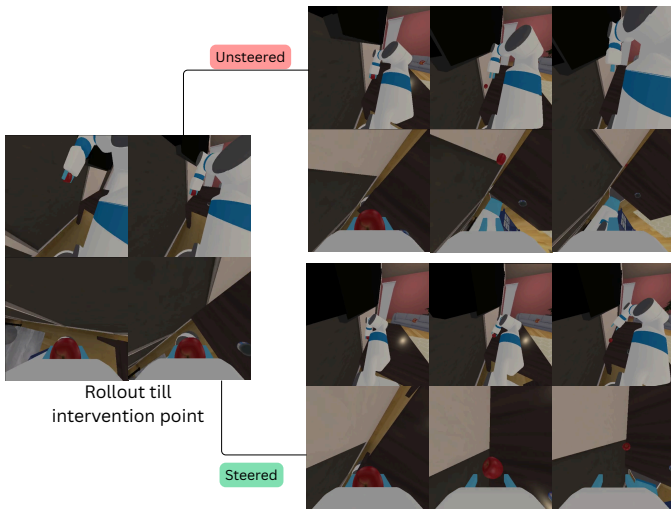
Embodiment	MMD Thresh.	Guidance Ratio	Guidance Steps	Sensor Res	Img Res	Num. Frames	Traj. Perturb	Num. Traj. Drawn
WidowX	0.9	40	4	640x512	1	0.01	5	
Google Robot	0.8	40	2	640x512	1	0.01	5	
Agile-X Dual Arm	1.05	2.5	15	1485x1182	1	0.01	5	

TABLE VIII: Hyper-parameter settings for Simpler-Env evaluations (see Table I) and Robotwin evaluations (see App. Section F1). Traj Perturb refers to the standard deviation of gaussian noise applied to *Primitive* trajectories. Num frames refers to number of frames used as history in the *Primitive* steerer. Num Traj Drawn refers to number of trajectories we overlay on image for the *Primitive* steerer.

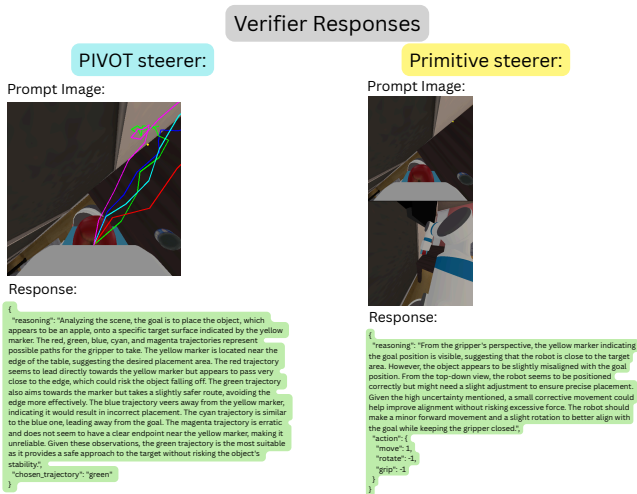
Method	Task	Subtask	MMD Thresh.	Guidance Ratio	Guidance Steps	Sensor Res	Img Res	Num. Frames	Traj. Perturb	Num. Traj. Drawn	Ensemble Ratio
Pivot	Prepare Groceries	Pick	0.7	10	8	512x512	1	0.25	5	N/A	
		Place	0.48								
	Tidy House	Pick	0.7								
		Place	0.48								
		OpenFridge	0.48								
Primitive	Prepare Groceries	Pick	0.7								
		Place	0.48								
	Tidy House	Pick	0.7								
		Place	0.48								
		OpenFridge	0.48								
Pivot + Primitive	Prepare Groceries	Pick	0.7								
		Place	0.48								
	Tidy House	Pick	0.7								
		Place	0.48								
		OpenFridge	0.48								

TABLE IX: Hyper-parameter settings for the MSHAB evaluations (see Section V-B) Note that Guidance Ratio, Steps, Resolution, and Frames are consistent across all methods. Traj Perturb refers to standard deviation of gaussian noise added to the PIVOT trajectories before overlaying onto RGB image.

- 1) *Hyperparameter Tables:*
- 2) *SimplerEnv Prompts:*



(a) Visual comparison of rollouts. The left side shows the state up to the intervention point, leading to Unsteered (top) and Steered (bottom) outcomes.



(b) Verifier responses showing the PIVOT steerer trajectories and the Primitive steerer reasoning output.

Fig. 10: Comparison of the steered vs unsteered rollouts and the corresponding verifier guidance. (a) Comparison of steered vs. unsteered trajectories. (b) Visualization of trajectory choices and VLM reasoning output.

Primitive Steering Prompt (Simpler-Env)

You are an expert AI controller for a mobile manipulator robot in a home environment.

SITUATION

- The primary camera view is from the robot's overhead or wrist camera.
- The robot has paused execution because the base policy is highly uncertain, likely due to misalignment, a potential collision, or being stuck.
- The name of each primitive is defined below in the Primitive List.

The task is to: <TASK_DESCRIPTION/>

SITUATION/

PRIMITIVE_LIST

- “Nudge Left”: translate the gripper to the left
- “Nudge Right”: translate the gripper to the right
- “Nudge Up”: translate the gripper vertically upwards

- “Nudge Down”: translate the gripper vertically downwards
- “Nudge Forward”: move the gripper forward into the scene
- “Retreat”: move the gripper backward outward from the scene
- “Gripper Open”: open the gripper
- “Gripper Close”: close the gripper

PRIMITIVE_LIST/
Analyze the scene to diagnose the error state:

- Misalignment:** Is the gripper too far left, right, up, or down relative to the target object?
- Collision:** Is the gripper pressing against a surface it shouldn't be? (Needs Retreat” or moving backward)
- Air Pushing:** Is the gripper moving in free space without touching the object? (Needs Nudge Forward”)

Select the primitive that best corrects this error to allow the robot to resume the task.

OUTPUT_FORMAT
You must conclude your response with a single, well-formed JSON object and nothing else. Do not use markdown formatting (like ``json) or add any text before or after the JSON block. The JSON object must contain three keys: "reasoning", "chosen_trajectory", and "gripper_state".

- “reasoning”: A string containing your detailed analysis of the error state. Explicitly mention if the gripper is misaligned (and in which direction) or if it is stuck. Explain why the chosen primitive corrects this specific error.
- “chosen_trajectory”: A string containing the name of the best primitive from the Primitive List. If no primitive helps, choose "none".
- “gripper_state”: A string containing the current state of the gripper, either "open" or "close".

Example of a perfect response format:

```
{
  "reasoning": str.
  "chosen_trajectory": str.
  "gripper_state": str.
}
```

OUTPUT_FORMAT/

Pivot Steering (Simpler-Env)

You are an expert AI controller for a mobile manipulator robot in a home environment.

SITUATION

- The primary camera view is from the robot's gripper. This view is overlaid with visualizations of potential future actions.
- The image provided shows five potential future trajectories for the gripper, colored red, orange, blue, cyan, and magenta. These trajectories represent different options the robot's base policy is considering.
- A marker at the end of each trajectory indicates the final predicted position and orientation of the gripper for that path. This marker may not always be visible.

The task is to: <TASK_DESCRIPTION/>

- The robot will fail the task if it encounters too high a cumulative force.
- The robot has detected high uncertainty in its next action and requires your expert guidance to select the best path forward.
- If all proposed trajectories appear unsafe or incorrect, it is best to reject all of them.

SITUATION/

TRAJECTORY_CHOICES

- “red”: Choose this to command the robot to follow the red path.
- “orange”: Choose this to command the robot to follow the orange path.
- “blue”: Choose this to command the robot to follow the blue path.
- “cyan”: Choose this to command the robot to follow the cyan path.
- “magenta”: Choose this to command the robot to follow the magenta path.
- “none”: Choose this to reject all proposed trajectories. This is the safest option if all paths lead to failure (e.g., collision, incorrect placement).

TRAJECTORY_CHOICES/
Analyze the scene and the proposed trajectories to determine which, if any, is the best path for the robot to follow to achieve its current task.

OUTPUT_FORMAT
You must conclude your response with a single, well-formed JSON object and nothing else. Do not use markdown formatting (like ``json) or add any text before or after the JSON block. The JSON object must contain two keys: "reasoning" and "chosen_trajectory".

- “reasoning”: A string containing your detailed analysis of the scene and each trajectory, incorporating fine-grained visual details to justify your

decision. Explain why the chosen trajectory is superior and why the others are suboptimal.

- **"chosen_trajectory"**: A string containing one of the valid choices: "red", "orange", "blue", "cyan", "magenta", or "none".

Example of a perfect response format:

```
{
  "reasoning": "str",
  "chosen_trajectory": "str"
}
```

OUTPUT_FORMAT/

3) Maniskill-HAB Prompts:

Primitive Steering Prompt (ManiSkill-HAB)

You are an expert AI controller for a mobile manipulator robot in a home environment.

SITUATION

- You have been given a set of images.
- One camera view is from the robot's gripper.
- The other camera view is from on top of the robot's head.
- If there are more than one images given for each perspective, then the images are in a sequence leading up to the current moment.

TASK DESCRIPTION/

- The robot will fail the task if it encounters too high a cumulative force over the duration of the task.
- The robot has detected high uncertainty in its next action and may require a corrective maneuver to ensure success.
- If it is difficult to ascertain the correct action, it is best to not influence the policy at all (all null action).

SITUATION/

AVAILABLE_ACTIONS

Available Base Movement Actions:

- 1: Move the robot backwards relative to where its arm is pointing.
- 0: Keep the robot in place; do not move.
- 1: Move the robot forwards relative to where its arm is pointing.
- null: Do not influence the current action, allow the robot to continue with its current trajectory.

Available Base Rotation Actions:

- 1: Rotate the robot clockwise relative to a top down perspective (i.e., turn right).
- 0: Keep the robot in place; do not rotate.
- 1: Rotate the robot counter-clockwise relative to a top down perspective (i.e., turn left).
- null: Do not influence the current action, allow the robot to continue with its current trajectory.

Available Gripper Actions:

- 1: Continue to hold the object.
- 1: Drop the object.
- null: Do not influence the current action, allow the robot to continue with its current trajectory.

AVAILABLE_ACTIONS/

Analyze the scene and determine the best action for the robot to pursue to achieve its current task.

OUTPUT_FORMAT

You must conclude your response with a single, well-formed JSON object and nothing else. Do not use markdown formatting (like ``json) or add any text before or after the JSON block.

The JSON object must contain two keys: "reasoning" and "action".

- **"reasoning"**: A string containing your detailed analysis incorporating fine-grained visual details to justify your decisions.
- **"action"**: An object containing the keys "move", "rotate", and "grip", with their corresponding numerical action values.

Example of a perfect response format:

```
{
  "reasoning": "The gripper is too far to the right of the goal. The robot needs to move its base forward and rotate slightly counter-clockwise to align properly before attempting to place the object. The gripper should remain closed.",
```

```
"action": {
  "move": 0,
  "rotate": null,
  "grip": -1
}
```

OUTPUT_FORMAT/

Pivot Steering Prompt (ManiSkill-HAB)

You are an expert AI controller for a mobile manipulator robot in a home environment.

SITUATION

- The primary camera view is from the robot's gripper. This view is overlaid with visualizations of potential future actions.
- The image provided shows three potential future trajectories for the gripper, colored red, green, and blue. These trajectories represent different options the robot's base policy is considering.
- A marker at the end of each trajectory indicates the final predicted position and orientation of the gripper for that path. This marker may not always be visible.

TASK DESCRIPTION/

- The robot will fail the task if it encounters too high a cumulative force.
- The robot has detected high uncertainty in its next action and requires your expert guidance to select the best path forward.
- If all proposed trajectories appear unsafe or incorrect, it is best to reject all of them.

SITUATION/

TRAJECTORY_CHOICES

"red": Choose this to command the robot to follow the red path.

"green": Choose this to command the robot to follow the green path.

"blue": Choose this to command the robot to follow the blue path.

"none": Choose this to reject all proposed trajectories. This is the safest option if all paths lead to failure (e.g., collision, incorrect placement).

Analyze the scene and the proposed trajectories to determine which, if any, is the best path for the robot to follow to achieve its current task.

TRAJECTORY_CHOICES/

OUTPUT_FORMAT

You must conclude your response with a single, well-formed JSON object and nothing else. Do not use markdown formatting (like ``json) or add any text before or after the JSON block.

The JSON object must contain two keys: "reasoning" and "chosen_trajectory".

- **"reasoning"**: A string containing your detailed analysis of the scene and each trajectory, incorporating fine-grained visual details to justify your decision. Explain why the chosen trajectory is superior and why the others are suboptimal.
- **"chosen_trajectory"**: A string containing one of the four valid choices: "red", "green", "blue", or "none".

Example of a perfect response format:

```
{
  "reasoning": "str",
  "chosen_trajectory": "str"
}
```

OUTPUT_FORMAT/

TABLE X: **Place Task Definitions.** d_x^g denotes distance to goal. Collision threshold: 7500.

Event Definitions		
grasped : $\neg \mathcal{K}_{g,t-1} \wedge \mathcal{K}_{g,t}$ at goal : $d_{x,t-1}^g > 0.15 \wedge d_{x,t}^g \leq 0.15$ left goal : $d_{x,t-1}^g \leq 0.15 \wedge d_{x,t}^g > 0.15$ rel. at goal : $d_x^g \leq 0.15 \wedge$ $\mathcal{K}_{g,t-1} \wedge \neg \mathcal{K}_{g,t}$ rel. out goal : $d_x^g > 0.15 \wedge \mathcal{K}_{g,t-1} \wedge \neg \mathcal{K}_{g,t}$		
Mode	Description	Condition
<i>Success Modes</i> (if $success \in E_{place}$)		
i. Place in goal	Releases in goal region; returns to rest.	$ E_{pl} \leq 4 \wedge (\text{rel. at goal} \in E_{pl} \vee d_{x,0}^g \leq 0.15) \wedge \text{idx}_{\text{left goal}} \leq \text{idx}_{\text{at goal}} \wedge \text{excessive} \notin E_{pl}$
ii. Dropped to goal	Releases outside, rolls/falls in; returns to rest.	$ E_{pl} \leq 4 \wedge (\text{rel. out goal} \in E_{pl} \vee d_{x,0}^g > 0.15) \wedge \text{idx}_{\text{left goal}} \leq \text{idx}_{\text{at goal}} \wedge \text{excessive} \notin E_{pl}$
iii. Dubious	In goal region and rest, but leaves before timeout.	$\text{idx}_{\text{at goal}} < \text{idx}_{\text{left goal}} \wedge \text{excessive} \notin E_{pl}$
iv. Winding	Leaves goal once, but eventually placed/dropped in.	$ E_{pl} > 4 \wedge \text{idx}_{\text{at goal}} > \text{idx}_{\text{left goal}} \wedge \text{excessive} \notin E_{pl}$
v. Success then coll.	Success followed by excessive collisions.	$\text{excessive} \in E_{pl}$
<i>Failure Modes</i> (if $success \notin E_{place}$)		
vi. Excessive coll.	Collision threshold exceeded.	$\text{excessive} \in E_{pl}$
vii. Didn't grasp	Fails to grasp at initialization.	$E_{pl} = () \wedge \text{excessive} \notin E_{pl}$
viii. Didn't reach	Grasps but never reaches goal region.	$ E_{pl} > 0 \wedge \text{at goal} \notin E_{pl} \wedge \text{excessive} \notin E_{pl}$
ix. Place in goal fail	Placed in goal, but rolls/falls out.	$\text{at goal} \in E_{pl} \wedge \mathcal{K}_{\text{placed latest}} \wedge \text{idx}_{\text{left goal}} > \text{idx}_{\text{at goal}} \wedge \text{excessive} \notin E_{pl}$
x. Dropped to goal fail	Dropped outside, rolls in, then rolls out.	$\text{at goal} \in E_{pl} \wedge \mathcal{K}_{\text{dropped latest}} \wedge \text{idx}_{\text{left goal}} > \text{idx}_{\text{at goal}} \wedge \text{excessive} \notin E_{pl}$
xi. Won't let go	In goal region, but never released.	$\text{at goal} \in E_{pl} \wedge \text{idx}_{\text{grasped}} > \text{idx}_{\text{rel. at goal}} \wedge \text{idx}_{\text{grasped}} > \text{idx}_{\text{rel. out goal}} \wedge \text{excessive} \notin E_{pl}$
xii. Too slow	Released in goal, but times out before rest.	Implies $\text{idx}_{\text{at goal}} > \text{idx}_{\text{left goal}} \wedge \text{excessive} \notin E_{pl}$

TABLE XI: **Open Task Definitions.** a_q is articulation position. Collision threshold: 10000.

Event Definitions		
opened : $\neg \mathcal{K}_{\text{open},t-1} \wedge \mathcal{K}_{\text{open},t}$ closed : $\mathcal{K}_{\text{open},t-1} \wedge \neg \mathcal{K}_{\text{open},t}$ slightly : $\neg \mathcal{K}_{\text{slight},t-1} \wedge \mathcal{K}_{\text{slight},t}$		
Mode	Description	Condition
<i>Success Modes</i> (if $success \in E_{open}$)		
i. Open success	Opens and returns to rest.	$\text{excessive} \notin E_{open} \wedge \text{idx}_{\text{opened}} > \text{idx}_{\text{closed}}$
ii. Dubious	Opens, returns to rest, then closes.	$\text{excessive} \notin E_{open} \wedge \text{idx}_{\text{opened}} < \text{idx}_{\text{closed}}$
iii. Success then coll.	Opens, then excessive collisions.	$\text{excessive} \in E_{open}$
<i>Failure Modes</i> (if $success \notin E_{open}$)		
iv. Excessive coll.	Collision threshold exceeded.	$\text{excessive} \in E_{open}$
v. Can't reach	Cannot reach articulation handle.	$\text{contact} \notin E_{open} \wedge \text{excessive} \notin E_{open}$
vi. Closed after open	Opens, but closes before rest.	$\text{closed} \in E_{open} \wedge \text{idx}_{\text{closed}} > \text{idx}_{\text{opened}} \wedge \text{idx}_{\text{closed}} > \text{idx}_{\text{slightly}} \wedge \text{excessive} \notin E_{open}$
vii. Slightly opened	Slightly opens, but not fully.	$\text{idx}_{\text{slightly}} > \text{idx}_{\text{opened}} \wedge \text{idx}_{\text{slightly}} > \text{idx}_{\text{closed}} \wedge \text{excessive} \notin E_{open}$
viii. Too slow	Opens, but times out before rest.	$\text{opened} \in E_{open}$
ix. Can't open	Reaches but cannot open.	$\text{contact} \in E_{open} \wedge \text{opened} \notin E_{open}$