

SwitchVLA: Execution-Aware Task Switching for Vision-Language-Action Models

Meng Li^{1,*}, Zhen Zhao^{1,*}, Zhengping Che^{1,*}, Fei Liao¹,
Kun Wu¹, Zhiyuan Xu¹, Pei Ren¹, Zhao Jin¹, Ning Liu¹, Jian Tang^{1,†}
¹Beijing Innovation Center of Humanoid Robotics
<https://switchvla.github.io>

Abstract: Robots deployed in dynamic environments must be able to not only follow diverse language instructions but flexibly adapt when user intent changes mid-execution. While recent Vision-Language-Action (VLA) models have advanced multi-task learning and instruction following, they typically assume static task intent, failing to respond when new instructions arrive during ongoing execution. This limitation hinders natural and robust interaction in dynamic settings, such as retail or household environments, where real-time intent changes are common. We propose **SwitchVLA**, a unified, execution-aware framework that enables smooth and reactive task switching without external planners or additional switch-specific data. We model task switching as a behavior modulation problem conditioned on execution state and instruction context. Expert demonstrations are segmented into temporally grounded contact phases, allowing the policy to infer task progress and adjust its behavior accordingly. A multi-behavior conditional policy is then trained to generate flexible action chunks under varying behavior modes through conditioned trajectory modeling. Experiments in both simulation and real-world robotic manipulation demonstrate that SwitchVLA enables robust instruction adherence, fluid task switching, and strong generalization—outperforming prior VLA baselines in both task success rate and interaction naturalness.

Keywords: Imitation Learning, Vision-Language-Action, Task Switching

1 Introduction

Vision-Language-Action (VLA) models have gained increasing attention for their ability to map natural language instructions and visual observations into executable robotic actions. Recent progress has significantly advanced instruction following and multi-task learning capabilities [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]. However, most of these models operate under the assumption that task intent remains fixed throughout execution. As a result, they struggle to respond when new instructions are issued mid-task—an assumption that falls short in real-world, interactive environments. For example, service robots in retail settings must flexibly adapt to customer hesitation, last-minute exchanges, or item returns that arise during ongoing interactions.

Despite its practical importance, the ability to dynamically respond to changing task intent has mainly been overlooked in prior work. Supporting such behavior requires VLA models to go beyond static instruction-to-action mapping. They must continuously track execution progress and reason about when to forward, rollback, or advance ongoing behaviors in response to changing user input. Crucially, this must be achieved while maintaining multi-task generalization, enabling robots to fluidly adapt across task boundaries without task-specific heuristics or handcrafted transitions.

*Equal contribution. †Corresponding author.

{gary.li, alex.zhao, z.che, leofly.liao, gongda.wu, eric.xu, chris.ren, mustafa.jin, neil.liu, jian.tang}@x-humanoid.com

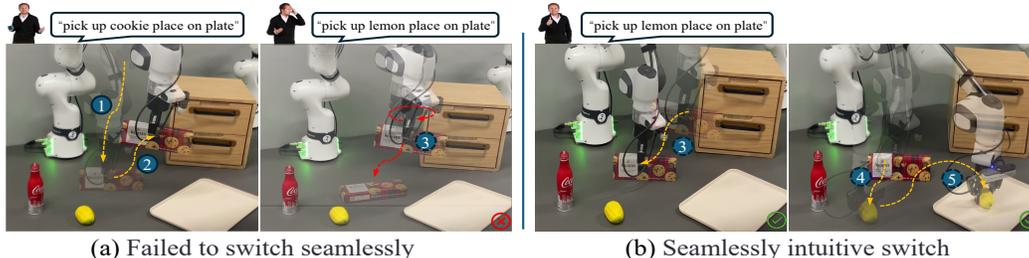


Figure 1: **(a)** Processes 1 and 2 show normal task execution. When the user changes their mind (e.g., “pick up lemon place on plate”), conventional VLA models cannot adjust its plan, leading to erratic behavior like oscillation or dropping items, as seen in Process 3. **(b)** A more natural response involves returning the previously held item (e.g., placing down the cookie box in Process 3 and then picking up the lemon and placing it on the plate in Processes 4 and 5).

Some recent approaches leverage large language or vision-language models for high-level task planning and interactive intent grounding [12, 13, 14, 15, 16], enabling re-planning when a new instruction is issued. However, these methods typically lack low-level execution flexibility: even if a task switching is inferred, the robot must wait for the current task to complete before re-planning, limiting real-time adaptability. Moreover, inference latency in large models makes high-frequency reactive reasoning infeasible. Other works attempt to enhance failure recovery and task switching by collecting additional demonstrations [17, 18], yet this strategy scales poorly—switching behaviors grow combinatorially with task diversity, making exhaustive data collection impractical. Therefore, a key challenge remains: *how to endow VLA models with the ability to smoothly and reactively switch tasks during execution using only existing trajectory data, without high-level planning or extra demonstrations*. As shown in Fig. 1, conventional VLA policies often struggle to naturally handle dynamic human instructions, such as discarding or mishandling objects when the target changes mid-task. In contrast, the expected model behavior enables intuitive responses by automatically executing recovery actions before proceeding with new commands—a behavior pattern closely mirroring natural human operation.

To address the limitations of prior work, we propose **SwitchVLA**, a unified execution-aware framework that enables dynamic behavior switching within a VLA policy. Instead of relying on external planners or manual scripts, SwitchVLA gives the model fine-grained execution awareness, enabling it to reason about task progress and adaptively forward, rollback, or advance actions in response to changing instructions. We formulate task switching as a state-conditioned behavioral modulation problem, where local execution dynamics drive continuous decision-making. To support this, we segment expert trajectories into time-based contact phases, enabling the policy to infer the current stage of the task and adjust its behavior accordingly. By employing execution-aware conditional policies, SwitchVLA can select behavior modes and predict future actions without requiring additional switch-specific data. Empirical results show that SwitchVLA enables coherent transitions, improves robustness under dynamic task changes, and generalizes effectively across multi-stage manipulation tasks.

In summary, our contributions are as follows:

- We propose SwitchVLA, a unified execution-aware framework that employs a novel training paradigm and innovative architecture to support dynamic task switching without relying on additional switch-specific data.
- We design a multi-behavior conditional policy capable of smoothly forwarding, rolling back, or advancing actions within a single policy backbone.
- We empirically validate SwitchVLA in both simulated and real-world robotic manipulation tasks, demonstrating substantial improvements in task transition smoothness, recovery effectiveness, and instruction adherence over existing VLA baselines.

2 Related Works

Vision-Language-Action Models for Robotic Control. VLA models have transformed robotic control by mapping multimodal inputs to actions [1, 19, 3, 20, 21, 22, 23, 8, 9, 24, 25, 10, 11, 26, 27]. Many works integrate language with robotic actions via imitation learning, where language serves as either goal specification [1, 20, 8, 9, 11, 27] or process-level planning [19, 28]. While these methods achieve strong generalization and instruction-following via architectural and data-scale improvements, they often assume a static execution flow—treating instruction following as a single-shot action prediction task, without modeling execution dynamics or supporting real-time adaptation. For instance, OpenVLA [8, 11] and π_0 [9] predict the next action from current observations and language but lack mechanisms for adjusting behaviors mid-execution. SayCan [19] and ViLa [28] depend on external modules to replan when states diverge.

Task Switching and Interactive Execution. Recent work explores interactive execution with increasing focus on task switching, categorized into modular systems and learning-based methods. Modular systems employ planners, hierarchical policies, or feedback loops [12, 13, 28, 14, 15, 29], but often rely on manually designed rollback, advance, or recovery strategies, limiting scalability and integration with end-to-end VLA models. They usually assume task completion before new instructions are issued. For example, ReKep [14] focuses on high-level interaction but suffers from latency issues; YAY Robot [30] supports context-aware correction but confines rollback to scripted demonstrations; Hi Robot [16] enables multi-behavior learning but requires complex, behavior-specific training. Learning-based methods [17, 31, 18] support model-driven interaction but depend on additional specialized data. RACER [17] introduces failure-aware recovery at the high level, relying on curated failure datasets, while ADC [18] augments learning with extra demonstrations.

While most existing methods focus on high-level planning, they often neglect the challenges of dynamic execution and real-time task adaptation. In contrast, our work targets the low-level execution of VLA models, proposing a unified policy that enables smooth, reactive, and adaptive task switching during execution—without requiring additional demonstrations.

3 Method

We introduce SwitchVLA, a unified learning framework for dynamic task-switchable action generation within Vision-Language-Action (VLA) settings. Our method treats task-switching not as a hard re-planning problem, but as a conditional behavior prediction challenge, where action dynamics adapt in response to evolving task intents and execution feedback.

3.1 Problem Formulation

Standard VLA Execution. Given the robot’s expert trajectory, $\tau = \{(l|o_t, q_t)\}_{t=0}^T$, where l is the task language instruction for the trajectory, o_t and q_t represent the visual observation and robot state (e.g., joint angles) at time t , respectively. The goal is to learn a policy that maps $(l|o_t, q_t) \mapsto a_{t+1}$ in a behaviorally coherent manner, where a_{t+1} denotes robot’s action at time $t + 1$.

Task Switching. In practical deployment, the robot may receive new task instructions l' at arbitrary times during execution. Such dynamic inputs introduce out-of-distribution observation-instruction pairs $(l'|o_t, q_t)$, posing significant challenges for generalization. We identify two core sub-problems: (i) instruction grounding—aligning the policy with the latest instruction l' , and (ii) execution-aware switching—using execution feedback (e.g., physical contact) to decide whether to forward, rollback, or switch to a new behavior mode.

To address these challenges in task switching, we introduce two auxiliary supervision signals—the contact state and the behavior mode—as key latent indicators of task phase and execution feedback.

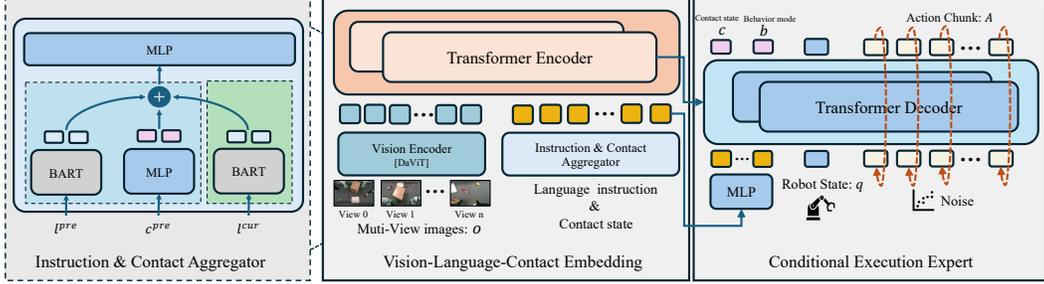


Figure 2: Overview of SwitchVLA. The framework consists of the Vision-Language-Contact Embedding module and the Conditional Execution Expert, which jointly fuse multimodal inputs to generate execution-aware and conditionally controlled actions.

3.2 Supervisory Signals for Instruction-Aware Control

Contact State. The contact state indicates physical interaction between the robot and objects. We define it as a binary variable $c_t \in \{0, 1\}$, where 0 denotes no contact and 1 denotes contact. It can be inferred through: tactile sensing, gripper open/close signals, heuristic motion or force thresholds, or vision-language parsing using pre-trained models. This binary state evolves over time and informs task phase progression, guiding the system’s next action strategy.

Behavior Mode. We define the behavior mode at timestep t as $b_t \in \{0 : \text{forward}, 1 : \text{rollback}, 2 : \text{advance}\}$, with each value corresponding to a distinct behavioral strategy: *forward* ($b_t = 0$) continues standard execution, *rollback* ($b_t = 1$) undoes previous actions upon detecting an intent mismatch while in contact, and *advance* ($b_t = 2$) transitions to a new subtask when the instruction updates and no physical interaction is present.

The labels of contact states and behavior modes can be weakly supervised from phase-aligned demonstrations or derived via behavior heuristics automatically parsed from execution feedback. Together, the contact state c_t and behavior mode b_t provide key supervision signals to address the task-switching challenge. The contact state offers real-time execution feedback to detect interaction phases, while the behavior mode encodes high-level task intent—whether to *forward*, *rollback*, or *advance*. These signals condition the policy to adaptively align with updated instructions and respond coherently under dynamic execution.

3.3 Architecture Overview

SwitchVLA establishes a unified architecture for robust and instruction-consistent task execution, as shown in Fig. 2. The architecture consists of two core components: (i) **Visual-Language-Contact (VLC) Embedding Module** encodes visual, language, and contact cues into unified representations. (ii) **Conditional Execution Expert** decodes behavior-aware actions conditioned on the current multimodal embedding. We build SwitchVLA upon Florence 2 [32].

3.3.1 VLC Embedding Module

To represent execution context with temporal and semantic richness, the Visual-Language-Contact (VLC) Embedding module fuses multi-view RGB observations, contact-aware execution cues, and paired task instructions into a unified token sequence. This embedding serves as the core state representation for downstream behavior selection.

Visual Encoder. We adopt a DaViT-based [33] backbone to encode multi-view RGB observations o into dense spatial tokens, projected for robust visual grounding.

Instruction and Contact Aggregator. By integrating historical and current contextual information, a rich conditioning signal is formed to enable behavior-aware action generation. Specifically, the previous instruction I^{pre} and contact state c^{pre} reflect past intent and interaction history, and the

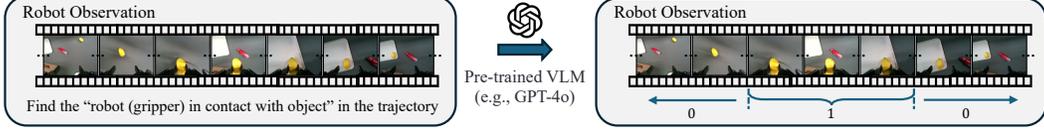


Figure 3: Identify and label time intervals of a specified event from trajectory data using a pre-trained VLM, such as GPT-4o. For example, with the prompt “Robot (gripper) in contact with object”, the model retrieves and labels the contact time intervals within the trajectory.

current instruction l^{cur} updates the semantic goal based on new task directives. These tokens are embedded using BART for language and an MLP for contact encoding, then concatenated to form the conditioning input token sequence.

Trajectory Parsing and Contact Annotation. In our implementation, we adopt a vision-language trajectory parsing approach: a pretrained VLM segments demonstrations into coarse, contact-aware phases based on wrist camera observations (Fig. 3), with gripper open-close signals incorporated to enhance the reliability of phase boundaries. These lightweight annotations enrich the tokenized input without modifying the core learning pipeline.

Transformer Fusion. All tokens are fused using Transformer encoder-decoder [34, 32], producing temporally and semantically rich embeddings.

3.3.2 Conditional Execution Expert

The Conditional Execution Expert module serves as a structured action decoder that integrates real-time contact signals with high-level behavioral intent to support adaptive and temporally coherent action generation. It simultaneously predicts three synchronized outputs in each timestep: the contact state $c_t \in \{0, 1\}$ reflecting whether the robot is in physical interaction with the environment, the behavior mode $b_t \in \{0, 1, 2\}$ indicating the current operational intent, forward execution, rollback of prior steps, or advance to a new subtask, and the action chunk $A_t = \{a_{t+k}\}_{k=1}^K$, which specifies a sequence of low-level actions spanning the next K timesteps.

The structured prediction design facilitates contact-aware behavioral adaptation, enabling reliable mode switching based on interaction signals. It also supports explicit behavioral reasoning over forward, rollback, and advance actions, as well as smooth temporal transitions through joint decoding of multi-step action chunks—reducing jitter and enhancing execution coherence.

3.4 Training and Inference

3.4.1 Training with Behavior-Specific Conditioning

Given labeled expert trajectories $\tau_i = \{(l^i | o_t^i, a_t^i, c_t^i, b_t^i)\}_{t=0}^T$, we design behavior-specific supervision to encourage the model to adapt its action generation based on both task instructions and physical interaction signals, as shown in Fig. 4. The behaviors are defined as follows:

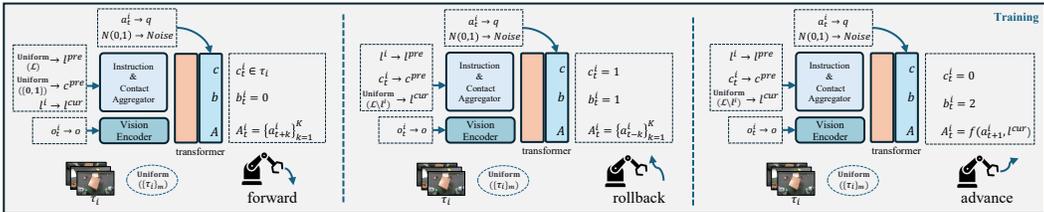


Figure 4: Illustration of the training pipelines for forward, rollback, and advance behaviors. Dynamic task transitions are achieved through policy modulation based on the predicted behavior mode, allowing the system to adapt to changing instructions and feedback.

- **forward** ($b_t = 0$): Predict the future action chunk $\{a_{t+1}^i, \dots, a_{t+K}^i\}$ under the matched instruction $l^{\text{cur}} = l^i$. Robustness to prior context is encouraged by randomly sampling different $l^{\text{pre}}, c^{\text{pre}}$ during training.
- **rollback** ($b_t = 1$): Generate reversed actions $\{a_{t-1}^i, \dots, a_{t-K}^i\}$ when presented with a mismatched instruction $l^{\text{cur}} \neq l^i$ and active contact $c_t^i = 1$. This encourages the model to recover gracefully from semantic mismatches through physical feedback.
- **advance** ($b_t = 2$): Interpolate linearly in joint space between the current action a_t^i and a canonical start pose a_0^{normal} , defined as the mean of initial trajectories for instruction l^{cur} . This interpolation is represented by $f(a_t^i, l^{\text{cur}})$, enabling rapid subtask switching under the assumption of no contact.

Action sequences are optimized using flow-matching loss [9], which promotes smooth and dynamically feasible trajectory generation. Contact states c_t^i and behavior modes b_t^i are supervised via standard classification objectives, ensuring accurate perception of both environmental interaction and high-level intent.

3.4.2 Inference with Conditional Switching

During execution, the policy operates on three key inputs: the current visual observation o_t , a pair of instructions ($l^{\text{pre}}, l^{\text{cur}}$), and the previously propagated contact state c^{pre} . At each timestep, the model jointly predicts (i) the new contact state c_t (to be used as c^{pre} in the next step), (ii) the behavior mode $b_t \in \{0, 1, 2\}$, and (iii) the corresponding action chunk A_t .

After each action chunk is executed, the model updates $l^{\text{pre}} \leftarrow l^{\text{cur}}$ and propagates the predicted contact state c_t to the next step. New instructions can be issued at any timestep, upon which the model re-evaluates b_t and dynamically adjusts behavior.

Depending on the predicted behavior mode, the system executes one of the following strategies: **forward** ($b_t = 0$): continue standard action prediction based on the current goal; **rollback** ($b_t = 1$): reverse previous actions to recover from errors under active contact; **advance** ($b_t = 2$): jump toward the new subgoal while updating the historical instruction as $l^{\text{pre}} \leftarrow l^{\text{cur}}$.

This conditional switching mechanism enables real-time, context-aware task transitions that couple semantic intent with physical interaction, supporting robust and adaptive task execution.

4 Experiments

4.1 Task Protocol

Task switching occurs when the execution of *Task A* is interrupted upon receipt of a new instruction for *Task B*. Based on this, we perform two types of task switching experiments: pairwise evaluation and long sequence evaluation. To facilitate a comprehensive evaluation of the model’s capabilities, in pairwise experiments, we send new instruction at different execution phases: early (pre-contact), mid (in-contact), and late (post-action). For long sequence experiments, we focus on mid-phase switching, as it offers a more rigorous evaluation of the method performance. A task switching is considered successful only if *Task A* enters its designated execution phase without failure, and when new instruction is triggered, *Task A* behaves as expected formulated in Section 3, and then *Task B* is subsequently completed. Failure conditions include the inability to place an object (e.g., due to dropping) at the target position, non-compliance with a new instruction, or an unexpected execution halt. Each evaluation is repeated over 12 trials, with results averaged to ensure statistical reliability.

Due to page constraints, extended experimental details, including setup, training procedures, model comparisons, and ablation results, are provided in Appendix (with its overview in Section A).

Table 1: Average success rates (%) on individual task (“No Switch”) and pairwise switching experiments in the LIBERO-Goal simulation.

Method	No Switch	Early Switch	Mid Switch	Late Switch
π_0 [9]	92.3	40.7	8.3	10.2
OpenVLA-OFT [11]	98.0	40.6	11.1	13.0
SwitchVLA	93.0	93.5	50.9	68.7

Table 2: Accumulated average success rate (%) of long sequence switch performance on both simulation and real experiments. Note that while individual tasks are not restricted to the same one, we ensure the use of distinct task pairs appeared in the long sequence.

Method	Task Sequence Length					
	A→B	A→B→C	A→...→D	A→...→E	A→...→F	
Simulation	π_0 [9]	0.0	0.0	0.0	0.0	0.0
	OpenVLA-OFT [11]	0.0	0.0	0.0	0.0	0.0
	SwitchVLA	100.0	83.3	83.3	75.0	50.0
Real	MT-ACT [2]	0.0	0.0	0.0	0.0	0.0
	Diffusion Policy [3]	0.0	0.0	0.0	0.0	0.0
	π_0 [9]	0.0	0.0	0.0	0.0	0.0
	SwitchVLA	95.6	83.3	79.2	58.3	54.2

4.2 Simulation Experiments

Experiment Setup. We conduct experiments on the LIBERO [35] platform using the LIBERO-Goal suite, which emphasizes task diversity within a single table manipulation environment. We select 8 tasks for multi-task training, each with 50 expert trajectories. For pairwise evaluation, we sample 9 representative task combinations from this set, while for long sequence evaluation, 6 consecutive tasks are chosen to assess sequential task-switching performance. We predict absolute joint space to ensure proper advance calculation.

Evaluation Results. We compare SwitchVLA with π_0 [9] and OpenVLA-OFT [11]. The overall pairwise switching success rates are summarized in Table 1. We first compare single-task performance (“No Switch”), where SwitchVLA achieves comparable success rates to these baseline methods. In pairwise task switching experiments, we observe that the baseline methods achieve a success rate of approximately 40% for early-phase switching, while demonstrating notably lower performance in mid- and late-phase switching scenarios ($\sim 10\%$). In contrast, SwitchVLA exhibits robust performance across all switching phases, outperforming competing methods by a substantial margin. Notably, simulation tasks terminate with distinct end poses (without resetting to home position), introducing additional challenges for late-switch scenarios. Leveraging our advance mechanism, the robot autonomously transitions to subsequent tasks’ starting positions, enabling natural task completion. For more challenging mid-phase long-sequence evaluations, SwitchVLA demonstrates robustness by steadily handling 6 consecutive task-switching scenarios—achieving competitive performance where other methods fail in the initial A→B transition, as shown in Table 2.

4.3 Real-World Experiments

Experiment Setup. Real-world experiments are conducted on two dual-armed Franka Emika Panda robot workstations, each executing four unique tasks. Data collection is performed via human-teleoperated demonstrations, covering foundational skills such as pushing, opening, and pick-and-place. Each task comprises 200 trajectories, recorded from two wrist-mounted cameras, one third-person RGB camera, and the robot’s proprioceptive states. For evaluation, objects are randomly positioned within a predefined workspace region.

Evaluation Results. We evaluate SwitchVLA against three representative manipulation policies: MT-ACT [2], Diffusion Policy (DP) [3], and π_0 [9, 36]. π_0 is a re-implementation based on the original paper. The results are summarized in Table 3. Consistent with simulation results, SwitchVLA matches baseline performance in “No Switch” evaluations while maintaining robust success rates

Table 3: Average success rates (%) on real-world experiments.

Method	Workstation 1				Workstation 2			
	No Switch	Early Switch	Mid Switch	Late Switch	No Switch	Early Switch	Mid Switch	Late Switch
MT-ACT [2]	56.3	0.0	0.0	4.9	50.0	0.0	0.0	0.0
DP [3]	93.8	13.2	4.8	13.2	83.3	49.3	0.0	34.7
π_0 [9]	97.9	50.1	0.0	32.0	100.0	75.0	0.0	64.6
SwitchVLA	95.9	99.3	95.1	75.0	100.0	95.1	96.5	94.4

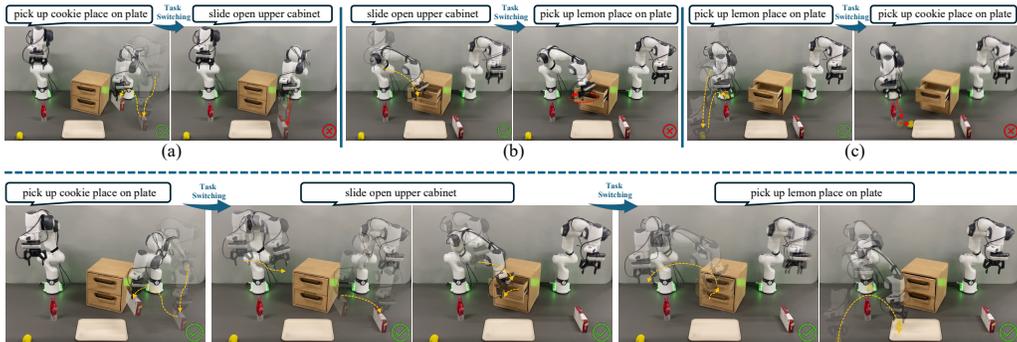


Figure 5: **Top:** Performance of π_0 [9] under pairwise task switching. (a), (b), and (c) each illustrate a unique task transition. Sudden switches during execution lead to erratic behaviors. **Bottom:** SwitchVLA enables smooth and consistent and instruction-aligned task transitions.

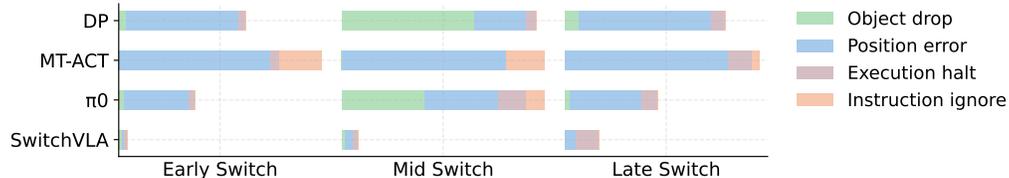


Figure 6: The proportion of failure causes during pairwise task switching in real-world experiments. Each bar’s length represents the ratio of all failure occurrences across total trials.

throughout all three switching phases. Similarly, long sequence experiments exhibit steady performance patterns as shown in Table 2. We also present real-world experiment scenes in Fig. 5, compared π_0 with SwitchVLA. In real-world experiments, we observe distinct failure patterns across switching phases, as shown in Fig. 6. During *Early Switch*, failures often result from advancing to incorrect positions; MT-ACT additionally tends to ignore updated instructions. *Mid Switch* is more challenging, with object drops and instruction neglect becoming common under active contact. In *Late Switch*, execution halts and misplacements occur frequently due to incomplete subtask transitions. These patterns reflect the increasing complexity of dynamic switching. SwitchVLA effectively mitigates such failures through its instruction-conditioned, execution-aware policy.

5 Conclusion

We propose **SwitchVLA**, a unified execution-aware framework for Vision-Language-Action (VLA) robots. By conditioning action generation on both task language and fine-grained execution signals, SwitchVLA enables seamless transitions across *forward*, *rollback*, and *advance* behaviors—without relying on modular planners or handcrafted switching logic. Extensive experiments in both simulated and real-world environments demonstrate that SwitchVLA significantly outperforms prior VLA baselines in instruction adherence, task recovery, and switching robustness. We position SwitchVLA as a more generalizable solution for instruction-conditioned control—capable of unifying diverse switching behaviors within a single policy framework. We believe this offers a new perspective for designing interaction-aware VLA systems in open-ended environments.

Limitations

While SwitchVLA demonstrates strong performance in dynamic task switching, several limitations remain: First, although our unified behavior-switching paradigm improves policy design, it relies on contact-based phase segmentation as a proxy for task progress. This may limit generalization to tasks involving abstract semantics or non-contact transitions. Incorporating richer temporal or semantic signals—such as intent prediction or temporal grounding—could address this issue. Second, the current use of three discrete behavior modes (forward, rollback, advance) may constrain flexibility in scenarios involving ambiguous or drifting human intent. Learning continuous or adaptive behavior representations could provide greater expressiveness. Finally, while our model is trained with weakly supervised paired demonstrations, adapting to noisy feedback, reinforcement refinement, or open-world settings remains an open challenge.

By introducing our approach and integrating it with high-level planning, we believe these limitations can be gradually addressed, extending the applicability of execution-aware VLA models in open-ended, interactive environments.

References

- [1] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- [2] H. Bharadhwaj, J. Vakil, M. Sharma, A. Gupta, S. Tulsiani, and V. Kumar. Roboagent: Towards sample efficient robot manipulation with semantic augmentations and action chunking. *arxiv*, 2023.
- [3] 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*, 2023.
- [4] A. Goyal, J. Xu, Y. Guo, V. Blukis, Y.-W. Chao, and D. Fox. Rvt: Robotic view transformer for 3d object manipulation. *CoRL*, 2023.
- [5] J. Zhang, Y. Guo, X. Chen, Y.-J. Wang, Y. Hu, C. Shi, and J. Chen. Hirt: Enhancing robotic control with hierarchical robot transformers. *arXiv preprint arXiv:2410.05273*, 2024.
- [6] Z. Fu, T. Z. Zhao, and C. Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. In *Conference on Robot Learning (CoRL)*, 2024.
- [7] C.-L. Cheang, G. Chen, Y. Jing, T. Kong, H. Li, Y. Li, Y. Liu, H. Wu, J. Xu, Y. Yang, et al. Gr-2: A generative video-language-action model with web-scale knowledge for robot manipulation. *arXiv preprint arXiv:2410.06158*, 2024.
- [8] 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.
- [9] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, et al. *pi_0*: A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- [10] Z. Wu, Y. Zhou, X. Xu, Z. Wang, and H. Yan. Momanipvla: Transferring vision-language-action models for general mobile manipulation. *arXiv preprint arXiv:2503.13446*, 2025.
- [11] M. J. Kim, C. Finn, and P. Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025.

- [12] H. Yuan, C. Zhang, H. Wang, F. Xie, P. Cai, H. Dong, and Z. Lu. Plan4MC: Skill reinforcement learning and planning for open-world Minecraft tasks. *arXiv preprint arXiv:2303.16563*, 2023.
- [13] K. Namasivayam, H. Singh, V. Bindal, A. Tuli, V. Agrawal, R. Jain, P. Singla, and R. Paul. Learning neuro-symbolic programs for language guided robot manipulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
- [14] W. Huang, C. Wang, Y. Li, R. Zhang, and L. Fei-Fei. Rekep: Spatio-temporal reasoning of relational keypoint constraints for robotic manipulation. *arXiv preprint arXiv:2409.01652*, 2024.
- [15] P. Liu, Y. Orru, C. Paxton, N. M. M. Shafiullah, and L. Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024.
- [16] L. X. Shi, B. Ichter, M. Equi, L. Ke, K. Pertsch, Q. Vuong, J. Tanner, A. Walling, H. Wang, N. Fusai, et al. Hi robot: Open-ended instruction following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*, 2025.
- [17] Y. Dai, J. Lee, N. Fazeli, and J. Chai. Racer: Rich language-guided failure recovery policies for imitation learning. *arXiv preprint arXiv:2409.14674*, 2024.
- [18] S. Huang, Y. Liao, S. Feng, S. Jiang, S. Liu, H. Li, M. Yao, and G. Ren. Adversarial data collection: Human-collaborative perturbations for efficient and robust robotic imitation learning. *arXiv preprint arXiv:2503.11646*, 2025.
- [19] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [20] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- [21] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- [22] X. Li, M. Liu, H. Zhang, C. Yu, J. Xu, H. Wu, C. Cheang, Y. Jing, W. Zhang, H. Liu, et al. Vision-language foundation models as effective robot imitators. *arXiv preprint arXiv:2311.01378*, 2023.
- [23] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, C. Xu, J. Luo, T. Kreiman, Y. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.
- [24] S. Haldar, Z. Peng, and L. Pinto. Baku: An efficient transformer for multi-task policy learning. *arXiv preprint arXiv:2406.07539*, 2024.
- [25] S. Nasiriany, S. Kirmani, T. Ding, L. Smith, Y. Zhu, D. Driess, D. Sadigh, and T. Xiao. Rt-affordance: Affordances are versatile intermediate representations for robot manipulation. *arXiv preprint arXiv:2411.02704*, 2024.
- [26] AgiBot-World-Contributors, Q. Bu, J. Cai, L. Chen, X. Cui, Y. Ding, S. Feng, S. Gao, X. He, X. Hu, X. Huang, S. Jiang, Y. Jiang, C. Jing, H. Li, J. Li, C. Liu, Y. Liu, Y. Lu, J. Luo, P. Luo, Y. Mu, Y. Niu, Y. Pan, J. Pang, Y. Qiao, G. Ren, C. Ruan, J. Shan, Y. Shen, C. Shi, M. Shi, M. Shi, C. Sima, J. Song, H. Wang, W. Wang, D. Wei, C. Xie, G. Xu, J. Yan, C. Yang, L. Yang, S. Yang, M. Yao, J. Zeng, C. Zhang, Q. Zhang, B. Zhao, C. Zhao, J. Zhao, and J. Zhu. Agibot world colosseum: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025.

- [27] 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.
- [28] Y. Hu, F. Lin, T. Zhang, L. Yi, and Y. Gao. Look before you leap: Unveiling the power of gpt-4v in robotic vision-language planning. *arXiv preprint arXiv:2311.17842*, 2023.
- [29] A. Xiao, N. Janaka, T. Hu, A. Gupta, K. Li, C. Yu, and D. Hsu. Robi butler: Remote multimodal interactions with household robot assistant. *arXiv preprint arXiv:2409.20548*, 2024.
- [30] L. X. Shi, Z. Hu, T. Z. Zhao, A. Sharma, K. Pertsch, J. Luo, S. Levine, and C. Finn. Yell at your robot: Improving on-the-fly from language corrections. *arXiv preprint arXiv: 2403.12910*, 2024.
- [31] S. Belkhale, T. Ding, T. Xiao, P. Sermanet, Q. Vuong, J. Tompson, Y. Chebotar, D. Dwibedi, and D. Sadigh. Rt-h: Action hierarchies using language. *arXiv preprint arXiv:2403.01823*, 2024.
- [32] B. Xiao, H. Wu, W. Xu, X. Dai, H. Hu, Y. Lu, M. Zeng, C. Liu, and L. Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks (2023). URL <https://arxiv.org/abs/2311.06242>, 2023.
- [33] M. Ding, B. Xiao, N. Codella, P. Luo, J. Wang, and L. Yuan. Davit: Dual attention vision transformers. In *European conference on computer vision*, pages 74–92. Springer, 2022.
- [34] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 2017.
- [35] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *arXiv preprint arXiv:2306.03310*, 2023.
- [36] A. Z. Ren. open-pi-zero. <https://github.com/allenzren/open-pi-zero>, 2024.
- [37] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [38] OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2025.

Appendix

A Overview

This appendix provides comprehensive technical details of SwitchVLA, organized into several key aspects. Section B presents implementation specifics, including: (1) our dynamic condition-action data sampling strategy (Section B.1), (2) model architecture details (Section B.2), and (3) the contact state labeling methodology (Section B.3). Finally, we document the hyperparameter configurations used during training in Section B.4.

We present additional experimental validation in Section C, beginning with an ablation study of model components. We then compare SwitchVLA against baseline approaches that require specialized switching-behavior datasets, demonstrating our method’s superior data efficiency.

Section D details our comprehensive task selection methodology. We begin by defining the fundamental single-task benchmarks, which serve as the basis for our evaluation. Building upon these atomic tasks, we then present:

- Carefully constructed pairwise switching tasks to evaluate basic switching capabilities.
- Strategically sampled long sequence switching tasks to assess complex, multi-step transitions.

The complete experimental evaluation is documented in Section E, structured as:

- Implementation setup and configuration details (Section E.1)
- Quantitative results for single-task and pairwise switching scenarios (Section E.2).
- Extended analysis of long-sequence task switching performance (Section E.3).

B Implementation Details

B.1 Data Sampling Strategy

To facilitate task switching learning without collecting additional demonstration data, we introduce a novel training paradigm that dynamically adapts action allocation according to both task instructions and physical interaction signals, shown in Section 3.4. The detailed data sampling strategy for training SwitchVLA is presented in Algorithm 1. Importantly, while we demonstrate this framework with contact states and three distinct behaviors, the proposed sampling strategy is fundamentally extensible—other condition-action combinations can be implemented.

B.2 Model Architecture Details

SwitchVLA adopts Florence-2-base [32] architecture as its backbone with several key modifications. Following the original Florence-2-base design, we maintain a consistent token dimension of 768 throughout the model.

Vision-Language-Contact Encoder. Our Vision-Language-Contact (VLC) model consists of three key components. For visual processing, we adopt Florence-2-base’s default image tokenizer, DaViT [33], which takes $3 \times 224 \times 224$ input images and produces 50 output tokens per image. For language processing, we utilize the default BART encoder [37], processing all instructions to a fixed length of 20 tokens through padding. The contact state is encoded using a single-layer MLP, producing a 1-dimensional token. We construct the input representation by concatenating the previous instruction tokens (l^{pre}), current instruction tokens (l^{cur}), and previous contact state token (c^{pre}), forming a combined condition vector of length 41. This condition vector is then concatenated with the image tokens and processed through the transformer encoder, enabling multi-modal integration of visual, linguistic, and contact signals.

Algorithm 1 Data Sampling Strategy in SwitchVLA Training

Input: Given one sampled observation o_t from an episode at timestep t , full trajectory actions A , robot state q_t , robot action a_t (with $q_t = a_t$), chunk size K , task label l , contact label c_t , all tasks \mathcal{L} , all contact state $\mathcal{C} = \{c_0, c_1\}$, task mode $\mathcal{M} = \{\text{current}, \text{previous}\}$, predefined probability to sample current task mode P_{cur} , predefined three behaviors $\mathcal{B} = \{b_0, b_1, b_2\}$, function that calculates the trajectory actions from current robot state to a specific task robot home state f .

- 1: $L_{\text{mode}} \sim \text{Categorical}([P_{\text{cur}}, 1 - P_{\text{cur}}])$ ▷ Randomly assign obs to be cur or pre task
- 2: **if** $L_{\text{mode}} = \text{current}$ **then** ▷ Forward behavior
- 3: $l^{\text{cur}} \leftarrow l$
- 4: $l^{\text{pre}} \sim \text{Uniform}(\mathcal{L})$ ▷ Sample a task that can be the same as the current task
- 5: $c^{\text{pre}} \sim \text{Uniform}(\{c_0, c_1\})$
- 6: $A_t \leftarrow a_{t+1:t+K}$
- 7: $b_t \leftarrow b_0$
- 8: **else**
- 9: $l^{\text{cur}} \sim \text{Uniform}(\mathcal{L} \setminus l)$ ▷ Sample a different task from the current task
- 10: $l^{\text{pre}} \leftarrow l$
- 11: $c^{\text{pre}} \leftarrow c_t$
- 12: **if** $c^{\text{pre}} = c_0$ **then** ▷ Rollback behavior
- 13: $A_t \leftarrow a_{t-1:t-K}$ ▷ Revert actions
- 14: $b_t \leftarrow b_1$
- 15: **else** ▷ Advance behavior
- 16: $A_t \leftarrow f(a_t, l^{\text{cur}})$ ▷ Get trajectory from current pos to next task’s home pos
- 17: $b_t \leftarrow b_2$
- 18: **end if**
- 19: **end if**
- 20: **Return:** One data sample $\leftarrow \{\text{‘observation’}: o_t, \text{‘action’}: A_t, \text{‘condition’}: \{\text{‘task_cur’} : l^{\text{cur}}, \text{‘task_pre’}: l^{\text{pre}}, \text{‘stage_pre’}: c^{\text{pre}}, \text{‘behavior’}: b_t\}\}$

Condition Execution Expert. Our action expert builds upon Florence-2-base’s transformer decoder, processing inputs from both the VLC encoder and a combined representation of aggregated conditions, robot proprioception, and action noise tokens (with chunk size K for diffusion). The 41-dimensional condition tokens are first compressed to 5 dimensions via a single-layer MLP, while robot joint positions are encoded into a 1-dimensional token through another MLP. These components are concatenated with the action noise tokens, forming the expert’s input (total dimension $5 + 1 + K$). The decoder outputs are processed through two separate 2-layer MLP projectors with SiLU activation: the first 5 tokens are decoded to predict contact state and behavior mode. For action generation, we employ diffusion flow-matching following [36]. Our multi-task learning objective combines: (1) flow-matching MSE loss, (2) binary cross-entropy (BCE) loss for contact state prediction (weight 10^{-2}), and (3) cross-entropy (CE) loss for behavior mode classification (weight 10^{-4}).

B.3 Contact State Label Acquisition

Multiple approaches exist for acquiring robot-object contact information, including manual human annotation, force/torque sensing, vision-language models (VLMs), and gripper state monitoring. In our framework, we employ an automated pipeline for contact state labeling: First, we extract the current task’s meta-data to identify (1) the active robotic arm and (2) relevant target objects. We then process the corresponding wrist-mounted camera image from the identified arm through VLMs to generate per-observation contact labels for entire trajectories.

For each trajectory, we generate contact state labels by sampling images at 1/3 Hz and processing them through GPT-4o [38] to detect gripper-object contact. The visual analysis is guided by structured text prompts. To enhance prediction reliability, we integrate the gripper’s opening status from proprioceptive observations, using this mechanical constraint to both refine contact predictions and eliminate false positives. An example of the prompt is shown below.

Table 4: Contact state label accuracy (%) for real-world workstation 1 and workstation 2. Refer to Sec D.1 for task ID associations.

Workstation 1				Workstation 2				Average
R1.1	R1.2	R1.3	R1.4	R2.1	R2.2	R2.3	R2.4	
91.7	88.0	80.8	60.6	79.5	84.2	73.9	86.7	80.7

System Prompt:

You are an expert in analyzing images from robotic arm grasping operations. You will receive an image from the hand-eye camera. Your task is to analyze this image according to the specific conditions provided by the user and determine if those conditions are met. Return ‘1’ if the conditions are satisfied and ‘0’ if they are not. Do not provide any additional explanation.

User Instruction:

Analyze the provided image and determine whether:

- The robotic gripper has grasped/touched sandwich or its edge (return 1)
- The gripper is visibly closed (return 1)
- Neither condition is met (return 0)

Finally, we apply 1D nearest-neighbor imputation to infer contact states for all observations in the trajectory. To further filter outliers, we employ a sliding window approach. Specifically, for a window size of K , we evaluate the contact results within the interval $[i - \lfloor K/2 \rfloor, i + \lfloor K/2 \rfloor]$ for each observation O_i . The contact state of O_i is then assigned the *mode* of the contact labels within this window, reducing noise and improving robustness.

The contact state information does not require high precision, as it serves primarily as coarse supervision for SwitchVLA to learn contact trends for relative prediction. For evaluation, we randomly sample 1,000 observations (including images and contact state predictions) from each real-world task for human verification of labeling accuracy. As shown in Table 4, the average accuracy across all tasks reaches 80.7%, which is sufficient for SwitchVLA contact state prediction. During training, SwitchVLA contact state prediction accuracy on evaluation dataset is 89.3%, 98.6%, and 93.5% respectively for workstation 1, workstation 2, and simulation.

B.4 Training Hyper-parameters

We train SwitchVLA with batch size of 512, using a constant learning rate scheduler which increases learning rate linearly between 0 and the initial learning rate with 5 warm-up steps. We train our models using the AdamW optimizer (with a learning rate of 5×10^{-5} with weight decay of 10^{-4} and Exponential Moving Average (EMA). During training, we predict 100 future actions (chunk size) and use first 30 actions for temporal ensembling to allow smooth and accurate manipulation, following implementations in ACT [21]. In real-world training and inference, we use 3-RGB images, and predict 16-dimension actions (7-DoF×2 & 2 grippers). In simulation, we use 2-RGB images and predict 9-dimension actions (7-DoF & 2 gripper tips). Both real-world and simulation share same other training settings.

C Additional Experiments

Ablating Model Components. We conduct an ablation study on SwitchVLA by examining two key components: contact state and behavior mode. Our real-world pairwise task-switching experiments reveal that removing these components reduces the model to a standard visual-language-action (VLA) imitation learning approach, which fails to determine appropriate task-switching timing during evaluation. As shown in Table 5, incorporating contact state information yields significant

Table 5: Average success rates (%) of additional teleoperated switching data experiment on two pairwise switching tasks “pick up plate and place on plate” & “pick up basket sandwich and place on plate”, and “pick up red gum and place on plate” & “push plate to customer”.

Method	Early Switch	Mid Switch	Late Switch
<i>SwitchVLA-base</i>	79.2	0.0	70.8
<i>SwitchVLA-base+contact</i>	87.5	95.8	100.0
<i>SwitchVLA-base+contact+behavior</i>	91.7	100.0	100.0

Table 6: Average success rates (%) of additional teleoperated switching data experiment on pairwise switching task “pick up plate and place on plate” & “pick up basket sandwich and place on plate”.

Method	Early Switch	Mid Switch	Late Switch
<i>SwitchVLA-base</i>	75.0	0.0	66.7
<i>SwitchVLA-base-add-data</i>	83.3	91.7	100.0
SwitchVLA	83.3	100.0	100.0

performance gains: early-switch accuracy improves from 79.2% to 87.5%, mid-switch from 0.0% to 95.8%, and late-switch from 70.8% to 100%. The addition of behavior mode modeling as an auxiliary loss further enhances performance, achieving 91.7% (early), 100.0% (mid), and 100.0% (late) success rates. These results demonstrate that both components are crucial for robust task-switching performance.

Comparison with additional behavior specific data To demonstrate the data efficiency and robustness of SwitchVLA, we additionally collect teleoperated task switching data and train the *SwitchVLA-base*. Specifically, we focus on a pairwise switch task—“Pick up plate place on plate” switch to “Pick up basket sandwich place on plate”. We collect 50 early switch data, 100 mid switch data, and 50 late switch data, combined with original data for both tasks (200 each) to train *SwitchVLA-base-add-data*. As shown in Table 6, we compare *SwitchVLA-base*, *SwitchVLA-base-add-data*, and SwitchVLA. The results demonstrate that additional training data improves the baseline model’s performance. Our analysis suggests that with sufficient additional behavior-specific data, the baseline model could achieve performance comparable to or even surpassing SwitchVLA. However, such an approach demands substantial data curation efforts, which are neither cost-effective nor scalable. In contrast, SwitchVLA effectively addresses these challenges without relying on extensive additional data, demonstrating superior efficiency and practicality.

D Evaluation Tasks

In this section, we provide more details of real-world and simulation task definitions and evaluation results. The overview of experiment setups for real-world and simulation is shown in Fig. 7.

D.1 Single Tasks

In workstation 1 and workstation 2, we setup the scenes and create 4 unique single tasks for each workstation, based on the single tasks, we combine two tasks and result in 12 pairwise switching tasks for each workstation. For LIBERO-Goal simulator, we select 8 unique single tasks and create 9 pairwise tasks. We describe all tasks below.

D.1.1 Real Franka Workstation 1

Task R1.1: pick up lemon place on plate. A basic pick-and-place task where the robot’s goal is to grasp a lemon from a specified location on table and place it onto the plate, evaluating its ability to handle small objects with precision.

Task R1.2: pick up coke bottle place on plate. This is a more challenging task where the robot

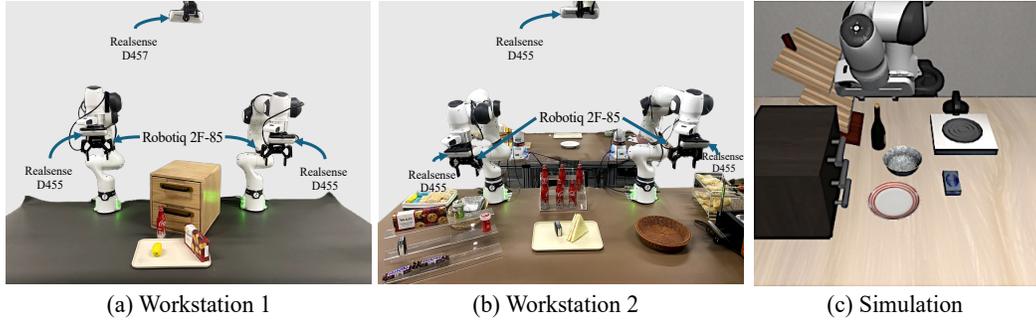


Figure 7: (a) Dual-armed Franka workstation 1 setup. (b) Dual-armed Franka workstation 2 setup. (c) LIBERO-Goal simulation setup.

needs to identify and lift a coke bottle from table and place it on a plate, assessing its capability to manipulate the small bottle cap with more precision.

Task R1.3: pick up cookie place on plate. This task involves carefully picking up a fragile cookie box without breaking it and placing it on a plate, testing the robot gripper’s precise manipulation skills.

Task R1.4: slide open upper cabinet. The robot’s target is to interact with an articulated cabinet door by sliding it open, evaluating its ability to handle constrained motion in a semi-structured environment.

D.1.2 Real Franka Workstation 2

Task R2.1: pick up red gum place on plate. The robot must locate and grasp a small red gum bottle from among multiple objects and place it on the plate, testing object discrimination and precise manipulation.

Task R2.2: pick up tissue place on plate. A challenging task where the robot’s objective is to pick up a pack of tissue and place it accurately on the plate, assessing its ability to handle deformable and flexible materials.

Task R2.3: pick up sandwich place on plate. The robot’s goal is to grasp a sandwich without deforming it and place it on the plate.

Task R2.4: push plate to customer. The robot needs to push a plate forward to simulate serving, testing its ability to perform controlled planar manipulation while maintaining object stability.

D.1.3 Simulation

We show simple description of tasks in LIBERO-Goal table manipulation environment here, refer to original paper [35] for more details.

Task S1: put the bowl on the plate. A nested placement task where the robot must position a bowl securely on a plate, assessing spatial awareness and stability awareness.

Task S2: put the cream cheese in the bowl. The robot’s objective is to pick up cream cheese and place into a bowl, evaluating its ability to perform precise manipulation.

Task S3: turn on the stove. The robot needs to interact with a stove knob to turn it on, testing its capability to operate mechanical controls with rotational motion.

Task S4: put the bowl on the stove. The robot’s target is to place a bowl onto a stove surface, requiring careful positioning and assess balanced manipulation.

Task S5: put the wine bottle on top of the cabinet. The robot must lift and place a wine bottle onto a high cabinet shelf, testing its ability to perform tasks requiring height adjustment and careful object placement.

Task S6: push the plate to the front of the stove. The robot’s goal is to slide a plate forward near the stove, evaluating controlled planar movement and spatial awareness in a kitchen environment.

Task S7: put the wine bottle on the rack. The robot needs to place the wine bottle into a rack, assessing precision in constrained placement tasks with orientation sensitivity.

Table 7: Selected pairwise switching tasks for real-world and simulation evaluation.

Workstation 1		Workstation 2		Simulation	
Pairwise Task ID	Task Pairs	Pairwise Task ID	Task Pairs	Pairwise Task ID	Task Pairs
PR1.1	R1.1→R1.2	PR2.1	R2.1→R2.2	PS1	S1→S3
PR1.2	R1.1→R1.3	PR2.2	R2.2→R2.1	PS2	S4→S1
PR1.3	R1.2→R1.1	PR2.3	R2.1→R2.3	PS3	S5→S6
PR1.4	R1.2→R1.3	PR2.4	R2.3→R2.1	PS4	S2→S6
PR1.5	R1.3→R1.1	PR2.5	R2.1→R2.4	PS5	S3→S5
PR1.6	R1.3→R1.2	PR2.6	R2.4→R2.1	PS6	S5→S3
PR1.7	R1.4→R1.1	PR2.7	R2.2→R2.3	PS7	S3→S2
PR1.8	R1.1→R1.4	PR2.8	R2.3→R2.2	PS8	S8→S3
PR1.9	R1.4→R1.2	PR2.9	R2.2→R2.4	PS9	S5→S6
PR1.10	R1.2→R1.4	PR2.10	R2.4→R2.2		
PR1.11	R1.4→R1.3	PR2.11	R2.3→R2.4		
PR1.12	R1.3→R1.4	PR2.12	R2.4→R2.3		

Table 8: Selected long sequence tasks for real-world and simulation evaluation.

	Long-Seq Task ID	Task Sequence
Workstation 1	LR1	R1.1→R1.3→R1.2→R1.1→R1.2→R1.4
Workstation 2	LR2	R2.1→R2.3→R2.4→R2.2→R2.3→R2.1
Simulation	LS1	S2→S5→S4→S1→S8→S3
	LS2	S1→S3→S2→S8→S3→S6
	LS3	S8→S3→S1→S3→S7→S4

Task S8: put the bowl on top of the cabinet. The robot’s objective is to lift and position a bowl on an elevated cabinet surface, testing its ability to handle objects and maintaining stability.

D.2 Pairwise Switching Tasks

We pair single tasks to form our pairwise switching tasks for real-world and simulation. We use all permutations for 4 tasks taken 2, *i.e.* ${}^4P_2 = 12$ for two each real-world workstation. In LIBERO-Goal simulation, we sample 9 pairs. Full details are shown in Table 7.

D.3 Long Sequence Switching Tasks

In Section 4, we show long sequence switching task performance in simulation environment for task LS1, the task selection is shown in Table 8, all other methods show 0.0% success rates, as they by chance fail at the beginning S2→S5 pairwise tasks. We samples two more long sequence—LS2 and LS3 tasks for comparison.

E Evaluation Results

In this section, we present experimental environment setup, as illustrated in Fig. 7, and evaluation details, including detailed results for polices on single and pairwise switching.

E.1 Experiment Setup

Workstations. We use 7-DoF Franka Emika plus RobotiQ 2F-85 grippers setup with three Intel RealSense D455 or D457 cameras, mounted at left arm wrist, right arm wrist and third-person top view. The data collection and inference frequency is 15Hz. All models employ 224 × 224 image size for both training and inference. During inference, all models are deployed on a single NVIDIA RTX 4090 GPU.

Table 9: Simulation joint position controller parameters.

Parameters	Value
output_min	1.0
output_max	1.0
kp	15,000.0
ramp_ratio	1.0

Table 10: Average success rates (%) on all single task experiments. We run 50 trials for each simulation evaluation, and 12 trials for real-world evaluation.

Method	Workstation 1				Workstation 2				Average (%)
	R1.1	R1.2	R1.3	R1.4	R2.1	R2.2	R2.3	R2.4	
MT-ACT [2]	75.0	83.3	66.7	0.0	100.0	16.7	83.3	0.0	53.1
DP [3]	83.3	91.7	100.0	100.0	100.0	91.7	83.3	58.3	88.5
π_0 [9]	100.0	91.7	100.0	100.0	100.0	100.0	100.0	100.0	99.0
SwitchVLA	100.0	91.7	91.7	100.0	100.0	100.0	100.0	100.0	97.9

Method	Simulation								Average (%)
	S1	S2	S3	S4	S5	S6	S7	S8	
OpenVLA-OFT [11]	100.0	100.0	100.0	94.0	94.0	100.0	98.0	98.0	98.0
π_0 [9]	98.0	96.0	100.0	96.0	98.0	90.0	76.0	88.0	92.3
SwitchVLA	100.0	62.0	100.0	100.0	100.0	82.0	100.0	100.0	93.0

Simulation. In LIBERO-Goal, we employ both wrist camera and third-person camera for training and inference. The image size is 128×128 . Each task contains 50 human-teleoperated demonstrations. We keep the same setting of each LIBERO-Goal task using its corresponding BDDL file, while using *JOINT_POSITION* control and changing some parameters in Robosuite [19] to keep proper joint position inference on simulation (default inference uses end effector *OSC_POSE* control). The configurations are shown in Table 9. Same to real-world setup, during inference, all models are deployed on a single NVIDIA RTX 4090 GPU.

E.2 Evaluation details

We evaluate SwitchVLA against three prior manipulation policies and selected the more representative real-robot works: MT-ACT [2], Diffusion Policy (DP) [3], and π_0 [9, 36]. π_0 is a re-implementation based on the original paper. For simulation, we compare with state-of-the-art (SOTA) VLAs— π_0 [9] and OpenVLA-OFT [11] on simulation. We show all individual task results in Table 10, and pairwise results in Table 11. For single-task performance (“No Switch”), SwitchVLA demonstrates comparable performance to selected baseline methods. In pairwise task switching experiments, we observe that the baseline methods achieve higher for early- and late-switching, while showing notably lower performance in mid-switching tasks. In contrast, SwitchVLA exhibits robust performance across all switching phases, outperforming competing methods by a substantial margin.

E.3 Additional Long Sequence Evaluations

In Section 4, we evaluate mid-phase task-switching performance in simulation for task LS1 (see Table 8 for task selection). While all baseline methods achieve 0.0% success rates—failing immediately during the initial S2→S5 transition—SwitchVLA demonstrates robust performance. To further validate our approach, we test on two additional long-sequence tasks (LS2 and LS3). As shown in Table 12, SwitchVLA consistently outperforms baselines, achieving substantially higher success rates on these challenging consecutive switching tasks. This further highlights our method’s robustness in dynamic task-switching scenarios.

Table 11: Average success rate (%) over 12 trials on all pairwise task switching experiments.

Method	Switch Stage	Pairwise Switching Experiment Task IDs										Average (%)		
		Workstation 1												
		PR1.1	PR1.2	PR1.3	PR1.4	PR1.5	PR1.6	PR1.7	PR1.8	PR1.9	PR1.10	PR1.11	PR1.12	
MT-ACT	Early	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DP		0.0	58.3	0.0	41.7	0.0	0.0	58.3	0.0	0.0	0.0	0.0	0.0	13.2
π_0		41.7	16.7	16.7	41.7	58.3	75.0	0.0	91.7	16.7	83.3	58.3	100.0	50.1
SwitchVLA		100.0	100.0	100.0	91.7	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	99.3
MT-ACT	Mid	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DP		0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	58.3	0.0	4.8
π_0		0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SwitchVLA		100.0	83.3	100.0	100.0	100.0	100.0	91.7	83.3	100.0	91.7	100.0	91.7	95.1
MT-ACT	Late	41.7	0.0	16.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	4.9
DP		50.0	0.0	25.0	0.0	0.0	0.0	0.0	41.7	0.0	0.0	41.7	0.0	13.2
π_0		41.7	41.7	33.3	50.0	66.7	50.0	0.0	0.0	0.0	41.7	0.0	58.3	32.0
SwitchVLA		50.0	91.7	100.0	100.0	91.7	100.0	0.0	83.3	0.0	83.3	100.0	100	75.0
		Workstation 2												
		PR2.1	PR2.2	PR2.3	PR2.4	PR2.5	PR2.6	PR2.7	PR2.8	PR2.9	PR2.10	PR2.11	PR2.12	
MT-ACT	Early	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DP		33.3	0.0	100.0	58.3	0.0	0.0	100.0	100.0	0.0	100.0	0.0	100.0	49.3
π_0		100.0	0.0	100.0	0.0	100.0	0.0	100.0	100.0	100.0	100.0	100.0	100.0	75.0
SwitchVLA		100.0	100.0	58.3	100.0	100.0	100.0	83.3	100.0	100.0	100.0	100.0	100.0	95.1
MT-ACT	Mid	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DP		0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
π_0		0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SwitchVLA		100.0	100.0	83.3	100.0	100.0	100.0	100.0	83.3	100.0	91.7	100.0	100.0	96.5
MT-ACT	Late	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
DP		58.3	0.0	0.0	66.7	0.0	0.0	83.3	100.0	66.7	0.0	41.7	0.0	34.7
π_0		100.0	100.0	100.0	100.0	100.0	0.0	0.0	100.0	100.0	0.0	75.0	0.0	64.6
SwitchVLA		100.0	100.0	100.0	100.0	100.0	83.3	100.0	100.0	100.0	100.0	50.0	100.0	94.4
		Simulation												
		PS1	PS2	PS3	PS4	PS5	PS6	PS7	PS8	PS9				
OpenVLA-OFT	Early	91.7	33.3	25.0	66.7	33.3	33.3	33.3	25.0	25.0				
π_0		25.0	8.3	33.3	83.3	33.3	66.7	41.7	41.7	33.3				
SwitchVLA		100.0	91.7	75.0	100.0	91.7	100.0	83.3	100.0	100.0				
OpenVLA-OFT	Mid	41.7	0.0	0.0	0.0	0.0	0.0	0.0	58.3	0.0				
π_0		0.0	0.0	0.0	58.3	0.0	0.0	16.7	0.0	0.0				
SwitchVLA		100.0	83.3	33.3	8.3	33.3	91.7	41.7	66.7	0.0				
OpenVLA-OFT	Late	0.0	100.0	0.0	16.7	0.0	0.0	0.0	0.0	0.0				
π_0		0.0	0.0	0.0	16.7	0.0	8.3	0.0	66.7	0.0				
SwitchVLA		100.0	41.7	41.7	66.7	25.0	100.0	41.7	100.0	100.0				

Table 12: Accumulated average success rate (%) of long sequence switch performance on simulation experiments. Note that while individual tasks are not restricted to the same one, we ensure the use of distinct task pairs appeared in the long sequence.

Method		Task Sequence Length				
		A→B	A→B→C	A→...→D	A→...→E	A→...→F
LS1	π_0	0.0	0.0	0.0	0.0	0.0
	OpenVLA-OFT	0.0	0.0	0.0	0.0	0.0
	SwitchVLA	100.0	83.3	83.3	75.0	50.0
LS2	π_0	0.0	0.0	0.0	0.0	0.0
	OpenVLA-OFT	50.0	0.0	0.0	0.0	0.0
	SwitchVLA	100.0	66.7	66.7	58.3	0.0
LS3	π_0	0.0	0.0	0.0	0.0	0.0
	OpenVLA-OFT	58.3	8.3	0.0	0.0	0.0
	SwitchVLA	66.7	50.0	41.7	33.3	0.0