

# ReCoVLA: VLM-Guided Reward Compilation for Failure Recovery in Vision-Language-Action Policies

Anonymous Author(s)

Affiliation

Address

email

1     **Abstract:** Vision-language-action (VLA) policies provide strong priors for  
2     language-conditioned manipulation, but remain brittle in off-nominal states re-  
3     quiring targeted recovery. We propose ReCoVLA—a failure-conditioned residual  
4     recovery framework that keeps a pretrained VLA policy frozen, uses an exter-  
5     nal vision-language model (VLM) to infer the failure mode and recovery stage,  
6     and compiles a structured reward from task-relevant components. Rather than  
7     using the VLM to generate actions or rewards directly, ReCoVLA uses it as a se-  
8     mantic reward selector: it predicts a recovery descriptor and reward mask for in-  
9     simulation residual-policy training, followed by zero-shot sim-to-real deployment  
10    of the trained recovery policies. This decouples high-level failure understanding  
11    from low-level corrective control to support different VLAs. Experiments across  
12    short-horizon, long-horizon, and contact-rich manipulation tasks show that Re-  
13    CoVLA outperforms the tested baselines on average. In simulation, our reward  
14    compiler improves average success from 36.7% for the fine-tuned  $\pi_{0.5}$  baseline  
15    to 66.7%. In physical zero-shot sim-to-real experiments, ReCoVLA achieves the  
16    best average performance, with 61.7% success.

17    **Keywords:** Failure recovery, VLA policies, Residual reinforcement learning

## 18    1 Introduction

19    Large vision-language-action (VLA) policies, especially those with flow matching [1, 2], have be-  
20    come more capable of mapping language instructions and visual observations directly to robot ac-  
21    tions. Their scale gives them broad semantic knowledge and strong nominal manipulation skills, but  
22    deployment still exposes a practical weakness: imitation-trained policies can enter action-induced  
23    states that differ from the training distribution [3]. For instance, when the robot misplaces an object,  
24    loses a grasp, or reaches an off-distribution state, one policy that solves the nominal task may lack  
25    the failure-specific corrective behavior needed for recovery [4, 5].

26    One direct solution is to fine-tune the VLA policy on additional recovery data. However, collecting  
27    failure-recovery demonstrations is costly and the fine tuning may forget previously acquired capa-  
28    bilities [6, 7, 8]. Reinforcement learning (RL) offers an alternative, but introduces two challenges  
29    for VLA recovery. First, the reward structure must match the current failure mode: a generic task-  
30    level reward is often too sparse [9], whereas activating every hand-designed reward can introduce  
31    conflicting objectives or alter the effective optimization problem [10]. Second, applying standard  
32    policy-gradient RL to flow-matching VLA models is difficult due to lack of likelihood availability.

33    The recent VLA  $\pi_{0.6}$  [11] addresses the action-likelihood issue by avoiding direct policy-gradient  
34    optimization through offline RL. However, it still requires training a separate value function, relies  
35    on costly human interventions, and does not directly solve the reward-design problem for failure  
36    recovery. These challenges motivate a recovery system that keeps the base VLA intact while adapt-  
37    ing the reward structure to a diverse set of failure modes. As shown in Figure 1, ReCoVLA is a

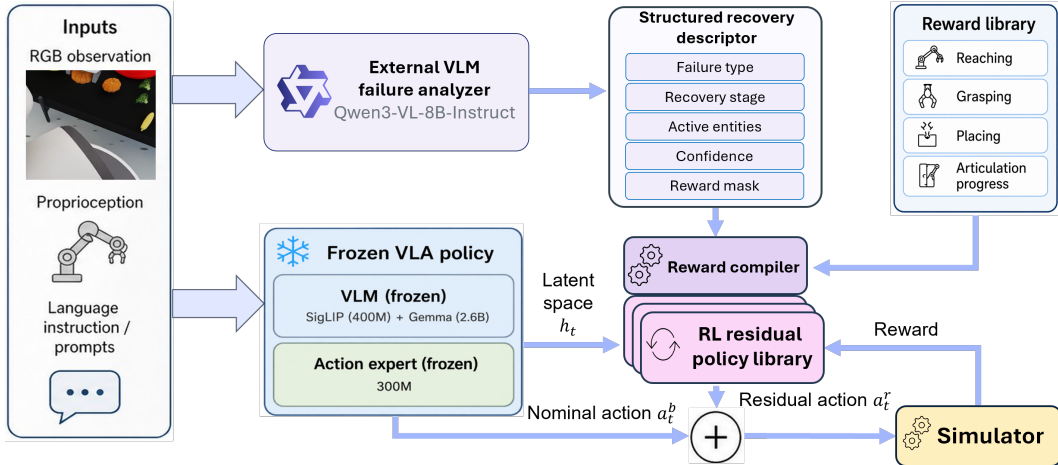


Figure 1: Overview of failure-conditioned residual VLA recovery. The frozen VLA policy maps robot inputs to a nominal action  $a_t^b$  and latent feature  $h_t$ . In parallel, an external Qwen3-VL-8B-Instruct VLM analyzes the RGB observation and prompt stream to produce a structured recovery descriptor containing the failure type, recovery stage, active entities, confidence, and reward mask. The reward compiler combines this descriptor with the reward library to generate the reward for residual-RL training. The residual policy observes  $h_t$ , receives the compiled reward during training, and outputs a corrective action  $a_t^r$ , which is added to the nominal VLA action before execution.

38 failure-conditioned residual recovery framework that uses an external VLM to analyze off-nominal  
 39 states and compile a structured reward for residual RL. The VLM does not directly generate robot  
 40 actions or free-form rewards. Instead, it produces a structured recovery descriptor containing the  
 41 failure type, recovery stage, active entities, confidence, and reward mask. A deterministic reward  
 42 compiler then grounds the entities, selects reward-library components, and inserts stage gates so that  
 43 reward terms are active only when their recovery preconditions are satisfied. An adaptively selected  
 44 one of residual policies is then trained over frozen VLA latents, preserving the nominal behavior of  
 45 the base VLA while learning corrective control for failure states.

46 We evaluate this design on three Fetch manipulation tasks selected to cover different recovery  
 47 regimes: long-horizon vegetable sorting, short-horizon soda-can disposal, and contact-rich toolbox  
 48 organization. The experiments show that reward structure is critical for reliable recovery, and the  
 49 proposed stage-gated reward compiler achieves the strongest overall performance across baselines.  
 50 In simulation, ReCoVLA improves average success by 30.0 percentage points over the fine-tuned  
 51  $\pi_{0.5}$  baseline. In physical Fetch experiments, it also achieves the highest average success, outper-  
 52 forming baselines by 18.3 percentage points. These results suggest that semantic failure identifica-  
 53 tion is useful not only for deciding when to recover, but also for determining which reward terms  
 54 should guide each stage of the recovery behavior. In summary, our contributions are:

- 55 1. We introduce a VLM-guided reward compiler that converts structured failure descriptors  
 56 into executable residual-RL rewards by grounding entities, selecting reward components,  
 57 and inserting stage gates, instead of using VLM for direct action or reward generation.
- 58 2. We formulate VLA failure recovery as a mixture of residual experts learning over VLA  
 59 latents, preserving nominal behavior while learning corrective control only for detected  
 60 off-nominal states.
- 61 3. We evaluate the framework across long-horizon sorting, short-horizon disposal, and  
 62 contact-rich organization tasks, demonstrating in simulation and on a physical robot that  
 63 ReCoVLA outperforms all baselines.

## 64 2 Related Works

65 **State-of-the-art VLA policies.** VIMA [12] and PaLM-E [13] established early VLA frameworks  
66 that robot behavior can be conditioned on rich multimodal prompts. RT-1/2 [14, 15] demonstrated  
67 that transformer-based VLA can scale to hundreds of real-world tasks. RoboCat [16], RT-X [17],  
68 RoboFlamingo [18], and GR-1 [19] explored broader pretraining and adaptation. More recently,  
69 open VLA policies [20, 1, 21, 22, 23, 24, 25, 2] highlight several emerging themes like efficient  
70 adaptation and improved reasoning. Rather than proposing new base VLA, our work focuses on  
71 robust recovery when the nominal policy enters failure states.

72 **RL-based refinement.** Residual reinforcement learning [26], residual policy learning [27], RL  
73 token [28], and residual feedback learning [29] show that they can repair imperfect controllers more  
74 efficiently than pure RL. Several works [30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40] study how to adapt  
75 pretrained policies with online RL. ReCoVLA is close in spirit to residual refinement, but differs by  
76 learning the residual in the latent space of a frozen VLA and by using an external VLM to condition  
77 the recovery reward structure on the failure mode.

78 **Failure recovery.** Prior work has studied damage adaptation and trial-and-error recovery [41], as  
79 well as RL recovery and safety monitors that guide the nominal policy near risky states [42]. Recent  
80 work studies language or VLM-guided recovery through demonstrations, language guidance, or  
81 executable actions [43, 44]. Instead of relying on manually designed monitors, resets, or task-  
82 specific recovery behaviors, ReCoVLA uses a VLM to produce a structured failure descriptor and  
83 compile the residual-RL reward to train mixture of residual policies.

## 84 3 Method

85 **Overview.** We propose a failure-conditioned residual recovery framework with three stages. First,  
86 we execute the frozen base VLA in simulation and use an external VLM to analyze the full failed  
87 rollouts, producing a catalog of simulator-observed failure categories and recovery descriptors. Sec-  
88 ond, a deterministic reward compiler constructs a reward from object-state reward components and  
89 trains one of residual policies in simulation for each recoverable failure category. Third, during real-  
90 world deployment, the VLM monitors a history of physical observations. When it detects a known  
91 failure category, the system dispatches the corresponding trained residual policy; when the category  
92 is unknown, no residual policy is called and the base VLA continues execution. The VLM is used  
93 for semantic failure identification and recovery-policy selection, not for direct action generation.

94 **Evaluation variants.** We consider six evaluation variants M1–M6. The first group M1–M4 uses  
95 the fine-tuned  $\pi_{0.5}$  policy as the base model. M1 is the frozen  $\pi_{0.5}$  base VLA without residual  
96 recovery. M2 adds a residual policy trained with a task-level reward, in which all task-relevant  
97 reward components are active simultaneously. M3 is an ablation of ReCoVLA that uses the VLM-  
98 detected failure category to activate only failure-relevant reward components, but applies no stage-  
99 aware gates. M4 is our full proposed method ReCoVLA, which uses the same failure-relevant  
100 reward components as M3, but adds stage-aware gates that activate later-stage terms only after their  
101 preconditions are satisfied, reducing premature reward activation and potential reward hacking. The  
102 second group M5–M6 evaluates whether the same recovery design can be applied to a different VLA  
103 backbone. M5 is the fine-tuned OpenVLA base policy without residual recovery. M6 applies the  
104 proposed stage-gated recovery design on top of OpenVLA.

### 105 3.1 Base VLA and residual control

106 We consider a language-conditioned manipulation task with instruction  $\ell$ . At time  $t$ , the frozen  
107 base VLA maps the robot observation,  $o_t$ , and  $\ell$  to a nominal action and an internal latent feature:  
108  $(a_t^b, h_t) = F_{\text{VLA}}(o_t, \ell)$ , where  $a_t^b \in \mathcal{A}$  is the base action and  $h_t$  is the VLA latent representation  
109 used by the residual recovery module. The base VLA is kept frozen throughout residual training

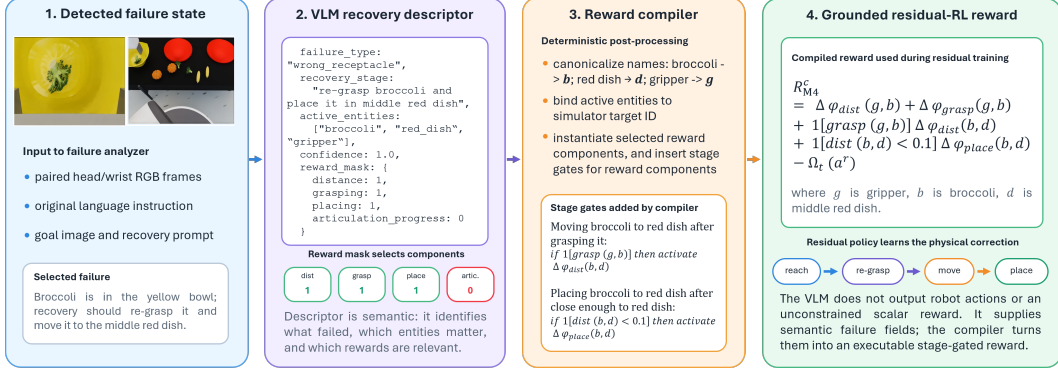


Figure 2: Example reward-compilation trace. The VLM analyzes the failed rollout and produces a structured recovery descriptor containing the failure type, recovery stage, active entities, confidence, and reward mask. The reward compiler then canonicalizes the entities, binds them to simulator object IDs, selects the corresponding reward-library components, and inserts stage-aware gates.

110 and deployment. A residual policy outputs an additive correction  $a_t^r \in \mathcal{A}$  from the VLA latent  
 111 representation:  $a_t^r \sim \pi_\theta^r(\cdot | h_t)$ . When residual recovery is active, the executed action is:

$$a_t = \Pi_{\mathcal{A}}(a_t^b + \beta_t a_t^r), \quad (1)$$

112 where  $\Pi_{\mathcal{A}}$  projects onto the valid action set and  $\beta_t \in \{0, 1\}$  is a binary deployment-time recovery  
 113 activation variable. This action composition preserves the nominal VLA behavior while allowing  
 114 the learned residual policy to provide targeted corrections in detected failure states.

### 115 3.2 Failure catalog construction

116 We construct the set of recoverable failure categories using simulator rollouts from the no-recovery  
 117 base policy M1. Let  $\tau_i^{M1} = \{(o_t, s_t, a_t^b)\}_{t=0}^{T_i}$  denote a rollout generated by the frozen base VLA  
 118 policy, where  $o_t$  is the visual observation,  $s_t$  is the robot state, and  $a_t^b$  is the base-policy action. For  
 119 each failed rollout, the external VLM failure analyzer,  $D_\phi$ , takes as input the full visual rollout, the  
 120 language instruction  $\ell_i$ , and a recovery-analysis prompt  $q$ , and outputs

$$\xi_i = D_\phi(\tau_i^{M1}, \ell_i, q) = (c_i, z_i, E_i, \rho_i, m_i). \quad (2)$$

121 Here  $c_i \in \mathcal{C}$  is the failure category,  $z_i \in \mathcal{Z}$  is the recovery stage,  $E_i \subseteq \mathcal{E}$  is the set of active  
 122 entities (e.g., "broccoli", "red\_dish"),  $\rho_i \in [0, 1]$  is the confidence score for failure states, and  $m_i \in$   
 123  $\{0, 1\}^K$  is a binary mask over the reward-component library. The high-confidence categories form  
 124 the simulator-trained failure catalog:  $\mathcal{C}_{\text{train}} = \{c_i : \rho_i \geq \tau_{\text{cat}}\}$ . For each  $c \in \mathcal{C}_{\text{train}}$ , the descriptor  
 125 provides the active entities  $E_c$  and the failure-relevant reward mask  $m_c$  for the reward compiler.

### 126 3.3 Reward variants and stage-gated compilation

127 The reward compiler uses a library of  $K$  object-state reward potentials  $\mathcal{R} = \{\varphi_k\}_{k=1}^K$ . Each po-  
 128 tential measures progress for a primitive recovery concept such as reaching, grasping, placing, or  
 129 articulation. For component  $k$ , we use potential-difference:

$$\Delta\varphi_k(s_t, s_{t+1}; E) = \varphi_k(s_{t+1}; E) - \varphi_k(s_t; E), \quad (3)$$

130 where  $E$  denotes the entities used to instantiate the component. The residual regularization term is

$$\Omega_t(a^r) = \lambda_1 \|a_t^r\|_2^2 + \lambda_2 \|a_t^r - a_{t-1}^r\|_2^2. \quad (4)$$

131 All reward terms are computed from simulator object states. The reward variants are defined as  
 132 follows. For M2, the task-level residual baseline activates all task-relevant components and entities  
 133 for instruction  $\ell$ . For example, in the soda-can disposal task, all reward components related to

134 picking and placing soda cans are enabled. Let  $m_\ell \in \{0, 1\}^K$  denote this task-level mask and  $E_\ell$   
 135 the set of task-level entities. Its reward is

$$R_{M2}^\ell(s_t, a_t, s_{t+1}) = \sum_{k=1}^K m_\ell^{(k)} \Delta\varphi_k(s_t, s_{t+1}; E_\ell) - \Omega_t(a^r). \quad (5)$$

136 For M3, the VLM-detected failure category activates only the failure-relevant components selected  
 137 by  $m_c$ , but all selected terms are active simultaneously:

$$R_{M3}^c(s_t, a_t, s_{t+1}) = \sum_{k=1}^K m_c^{(k)} \Delta\varphi_k(s_t, s_{t+1}; E_c) - \Omega_t(a^r). \quad (6)$$

138 M4 uses the same failure mask  $m_c$  as M3, but adds a binary stage-aware gate  $g_{c,k}(s_t; E_c) \in \{0, 1\}$   
 139 for each selected component:

$$R_{M4}^c(s_t, a_t, s_{t+1}) = \sum_{k=1}^K m_c^{(k)} g_{c,k}(s_t; E_c) \Delta\varphi_k(s_t, s_{t+1}; E_c) - \Omega_t(a^r). \quad (7)$$

140 The gates encode recovery-stage preconditions; M4 differs from M3 only by the stage-aware gates.  
 141 For example, in pick-and-place recovery, the placing component is inactive until the object is  
 142 grasped, preventing the residual policy from optimizing target motion before it has re-established  
 143 control of the object. M6 uses the same reward as M4, but replaces  $\pi_{0.5}$  with OpenVLA.

144 **Reward compiler.** The reward compiler acts as a deterministic interpreter that maps a schema-  
 145 valid VLM descriptor into an executable residual-RL reward. As illustrated for the sorting task in  
 146 Figure 2, the VLM detects a wrong-receptacle placement failure from paired rollout frames and out-  
 147 puts a structured descriptor  $\xi_c = (c, z, E_c, \rho, m_c)$ . The compiler first verifies that the failure category  
 148  $c$  is in the recoverable catalog and that the confidence  $\rho$  and reward mask  $m_c$  satisfy the required  
 149 thresholds. It then canonicalizes the VLM-predicted active entities  $E_c$  using the task specification  
 150 and the simulator’s object map, assigning each entity a semantic role (e.g., *gripper*, *target*, or *artic-*  
 151 *ulated part*). Consequently, the active-entity set serves as a candidate pool rather than a rigid argu-  
 152 ment list. Each reward component declares an argument signature to which the compiler binds only  
 153 the required roles. In our sorting example, the descriptor identifies the gripper, the misplaced broc-  
 154 coli, the source yellow bowl, and the target dish; thus, the grasping component is instantiated with  
 155 (gripper, broccoli), while the placing component is instantiated with (broccoli, target dish). Next,  
 156 the compiler attaches fixed stage-gate templates determined by the failure type and recovery stage  
 157 ( $z$ ), such as grasp-before-place or contact-before-articulation. This pipeline underscores a critical  
 158 design choice: the VLM avoids generating unstable, free-form scalar rewards or thresholds. Instead,  
 159 it provides structured semantic fields that the compiler safely grounds into the stage-gated reward  
 160 used in Eq. 7. To validate the stability and reproducibility of these VLM-detected states, we provide  
 161 a confusion matrix alongside comprehensive compiler implementation details in Appendix A.3.

162 **RL residual policy library.** Residual policies are trained entirely in simulation. For M2, we train  
 163 a task-level residual policy using  $R_{M2}^\ell$ . For M3, M4, and M6, we train failure-category residual  
 164 policies using the corresponding rewards in Eqs. (6) and (7) with PPO [45]. The actor observes the  
 165 VLA latent feature  $h_t$ . Simulator state is used to compute the object-state rewards and to train the  
 166 critic, but it is not provided to the deployed actor. The reward compiler uses object-state reward  
 167 potentials from a small library  $\mathcal{R} = \{\varphi_k\}_{k=1}^K$ . In our experiments, the primitive components are  
 168 distance progress, grasp state, placement progress, and articulation closing progress:

$$\begin{aligned} \varphi_{\text{dist}}(a, b) &= \max\left(-1, 1 - \frac{\|p(a) - p(b)\|_2}{d_{\text{init}}(a, b)}\right), \\ \varphi_{\text{grasp}}(\text{gripper}, \text{obj}) &= q_{\text{grasp}}(\text{gripper}, \text{obj}) \in [0, 1], \\ \varphi_{\text{place}}(\text{obj}, \text{target}) &= \max\left(-1, 1 - \frac{\|p(\text{obj}) - p(\text{target})\|_2}{d_{\text{init}}(\text{obj}, \text{target})}\right), \\ \varphi_{\text{close}}(\text{obj}) &= \max\left(-1, 1 - \frac{|\theta(\text{obj}) - \theta_{\text{closed}}|}{|\theta_{\text{init}}(\text{obj}) - \theta_{\text{closed}}|}\right). \end{aligned} \quad (8)$$



Figure 3: Physical experiments setup. Columns show the three evaluation tasks: organizing the toolbox; sorting vegetables; and soda-can disposal. Rows show the initial experiment setup, example task progress, a detected failure state from the VLM failure analyzer, and the recovered state after residual correction. The examples illustrate that the recovery descriptor must be task-specific and stage-specific: recovering a dropped cable or box, correcting an object-category placement error, and completing soda-can disposal require different active entities and reward components.

169 Where  $p(\cdot)$  denotes Cartesian position,  $\theta(\cdot)$  denotes the relevant articulation state,  $\theta_{\text{closed}}$  is the  
 170 closed target angle, and  $q_{\text{grasp}}$  is the simulator grasp-quality or grasp-indicator signal. The normal-  
 171 izers  $d_{\text{init}}(\cdot, \cdot)$  and  $\theta_{\text{init}}(\cdot)$  are measured when recovery is triggered.

172 For a residual method with reward  $R^u$ , where  $u$  denotes either the task instruction  $\ell$  for M2 or the  
 173 failure category  $c$  for M3/M4/M6, the training objective is:

$$\theta_u^* = \arg \max_{\theta_u} \mathbb{E} \left[ \sum_{t=0}^{T_u-1} \gamma^t R^u(s_t, a_t, s_{t+1}) \right], \quad (9)$$

174 with actions composed as in Eq. (1). After training, the learned residual policies are frozen and  
 175 stored in a recovery policy library. More details of the training can be found in Appendix A.4.

## 176 4 Experimental Results

### 177 4.1 Experiment setup

178 We evaluate all six variants M1–M6 on a Fetch robot in both simulation and physical experiments.  
 179 The benchmark includes three tabletop tasks spanning distinct recovery regimes: contact-rich tool-  
 180 box organization, long-horizon vegetable sorting, and short-horizon soda-can disposal. Each task is  
 181 specified by a language instruction and evaluated over 20 trials. We report both binary success and  
 182 a normalized goal-fulfillment quality score (Q-score), which assigns partial credit for intermediate  
 183 recovery progress. Task prompts, Q-score rubrics, and additional experimental details are provided  
 184 in Appendix A.2. Figure 3 shows representative task setups, detected failure states, and recovered  
 185 states.

186 Across simulation rollouts, we observe several recurring recoverable failures. In the toolbox task,  
 187 common failures include a cable or box left on the ground, an object grasped but not inserted into  
 188 the toolbox, and a toolbox lid that remains open after object placement. In the sorting task, failures  
 189 include vegetables placed in the wrong receptacle, objects remaining on the table, and incomplete

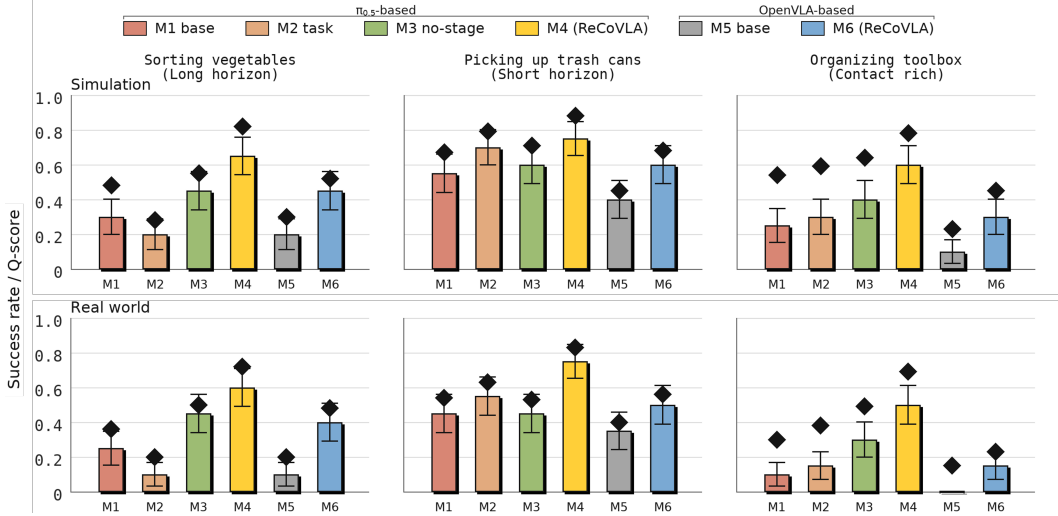


Figure 4: Simulation and physical experiments over 20 trials per method and task. The top and bottom rows report simulation and physical results, respectively. Bars show success rate, error bars show binomial standard error, and black diamonds show Q-score.

190 category-specific placement. In the soda-can disposal task, failures include a can remaining on the  
 191 table, a dropped can near the bin, and a can held by the gripper but not released into the trash can.  
 192 A complete library of trained residual policies is provided in Appendix A.2.

193 The mixture of residual policies library is deployed zero-shot on the physical robot. The robot exe-  
 194 cutes the frozen base VLA by default. In parallel, the VLM maintains a sparse history of five RGB  
 195 observations:  $H_t = \{I_{t-4\Delta}, I_{t-3\Delta}, I_{t-2\Delta}, I_{t-\Delta}, I_t\}$ ,  $\Delta = 15s$ . Given this history, the instruction,  
 196 and the recovery prompt, the VLM works as failure state detector, and residual recovery is activated  
 197 only when the VLM detects a failure category that was observed and trained in simulation:

$$\beta_t = \mathbf{1} [\rho_t \geq \tau_{\text{deploy}} \wedge c_t \in \mathcal{C}_{\text{train}}]. \quad (10)$$

198 If  $\beta_t = 1$ , the dispatcher invokes the trained residual policy associated with the detected category.  
 199 If  $c_t \notin \mathcal{C}_{\text{train}}$ , no residual policy is called; the system continues with the base VLA, and the trial is  
 200 counted as unrecovered if the base policy cannot recover from that state.

## 201 4.2 Simulation and physical results

202 Figure 4 summarizes the simulation and physical results. We observe that the full stage-gated reward  
 203 compiler, M4, achieves the strongest overall performance. In simulation, M4 improves average  
 204 success from 36.7% for the no-recovery  $\pi_{0.5}$  baseline M1 to 66.7%, with average Q-score increasing  
 205 from 0.56 to 0.83. In physical experiments, M4 reaches 61.7% average success and 0.75 average  
 206 Q-score, outperforming the baselines on each task by 18.3 points in success and 0.21 in Q-score.

207 The ablations indicate that recovery depends strongly on reward structure. M2 uses a task-level  
 208 residual reward, while M3 activates VLM-selected failure-relevant reward components without  
 209 stage-aware gates. Both improve some tasks, but neither matches M4. In simulation, M3 reaches  
 210 48.3% average success, 18.4 percentage points below M4, and has lower Q-scores on every task.  
 211 Thus, selecting relevant reward terms is useful but insufficient: the main gain comes from stage-  
 212 aware gates that suppress reward components before their recovery preconditions are satisfied.

213 Task-level results support the same conclusion. In simulation, M4 achieves the best result on all  
 214 three tasks, improving toolbox organization from 25% to 60% success and vegetable sorting from  
 215 30% to 65%. These gains are largest in tasks requiring staged recovery, such as re-grasping, object  
 216 placement, and lid closing. The same design also benefits another VLA backbone: M6 improves  
 217 OpenVLA average success from 23.3% for M5 to 45.0%. In physical experiments, M4 again per-  
 218 forms best on all tasks, reaching 50%, 60%, and 75% success on toolbox organization, vegetable

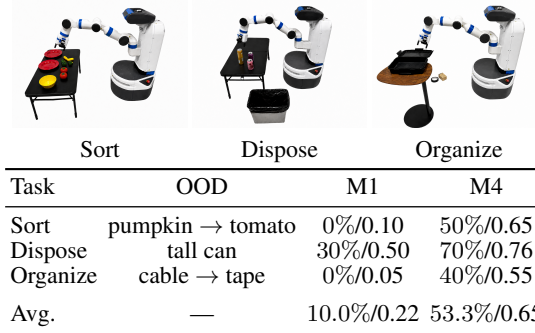
219 sorting, and soda-can disposal, respectively. The soda-can task is the most favorable to generic  
 220 recovery, with M2 reaching 55% success, but M4 still achieves the highest success and Q-score.  
 221 Overall, M1 and M5 show that base VLAs remain brittle after failures, while the M3–M4 gap iso-  
 222 lates the value of stage-gated reward compilation over simply activating failure-relevant rewards.

### 223 4.3 Out of distribution stress test

224 Figure 5 reports an out-of-distribution (OOD) evaluation under object substitutions that preserve  
 225 task structure but change visual appearance or geometry. The no-recovery base policy M1 drops  
 226 to 10.0% average success and 0.22 Q-score, including zero successful trials for two tasks. Un-  
 227 der the same OOD changes, M4 achieves 53.3% average success and 0.65 Q-score, with a 43.3  
 228 percentage-point success gain and a 0.43 Q-score gain over M1. The soda-can variant remains the  
 229 most robust at 70% success and 0.76 Q-score for M4, while the contact-rich toolbox variant is harder.

230

231 To validate the advantage of M4 over re-  
 232 cent recovery and VLA+RL methods, we  
 233 additionally evaluate against RLinf [46]  
 234 and RACER [43], and test M4 on  
 235 Behavior-1K [47] using a different mo-  
 236 bile humanoid manipulator, the Galaxea  
 237 R1 Pro. More details are provided in Ap-  
 238 pendix A.1.



Task	OOD	M1	M4
Sort	pumpkin → tomato	0%/0.10	50%/0.65
Dispose	tall can	30%/0.50	70%/0.76
Organize	cable → tape	0%/0.05	40%/0.55
Avg.	—	10.0%/0.22	53.3%/0.65

## 239 5 Conclusion

240 We presented a failure-conditioned reward compilation framework for residual VLA recovery. In-  
 241 stead of fine-tuning the full base policy or training a residual policy with a fixed reward, ReCoVLA  
 242 keeps the pretrained VLA frozen, uses an external VLM to infer the current failure mode and re-  
 243 covery stage, and compiles a reward from selected task-grounded components. This design lets  
 244 the residual policy receive feedback matched to the current recovery subproblem while preserv-  
 245 ing the nominal competence of the base policy. Across three manipulation tasks, the experiments  
 246 show that reward selection is central to reliable recovery. The proposed ReCoVLA is the strongest  
 247 method overall, increasing average simulation success from 36.7% for the fine-tuned  $\pi_{0.5}$  baseline  
 248 to 66.7%. The same stage-gated recovery mechanism also improves OpenVLA in simulation and  
 249 physical experiments. In physical experiments, ReCoVLA achieves the best success and Q-score on  
 250 all three tasks and outperforms the other tested recovery methods on average. These results suggest  
 251 that semantic failure understanding should shape not only the choice to recover, but also the reward  
 252 structure used to learn the recovery behavior.

## 253 6 Limitations

254 ReCoVLA requires each recoverable failure category to be reproducible in simulation before de-  
 255 ployment. When a new failure mode is observed on the real robot, we can add it to the recovery  
 256 catalog by reproducing the corresponding state distribution in simulation, compiling a stage-gated  
 257 reward, and training a new residual policy. However, the current system does not synthesize or  
 258 train a new recovery policy online during a physical trial. The empirical evaluation is limited in  
 259 scale. We evaluate three tabletop Fetch manipulation tasks with 20 trials per task and method, cov-  
 260 ering long-horizon sorting, short-horizon disposal, and contact-rich toolbox organization, but not  
 261 the full diversity of manipulation failures. Residual policies are trained in simulation and deployed  
 262 zero-shot, so performance can still be affected by sim-to-real mismatch, perception errors, contact  
 263 variation, and VLM failure-classification mistakes. Future work should automate failure-catalog  
 264 expansion, learn or generate reward gates, and evaluate broader physical baselines.

## References

- [1] 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.
- [2] P. Intelligence, K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, et al.  $\pi_{0.5}$ : a vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.
- [3] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 627–635, 2011.
- [4] Y. Dai, J. Lee, N. Fazeli, and J. Chai. RACER: Rich language-guided failure recovery policies for imitation learning. In *IEEE International Conference on Robotics and Automation*, 2025.
- [5] S. Pan, Y. Xu, R. Xu, Z. Zhou, S. Wu, and Z. Yu. Self-correcting robot manipulation via gaussian-splatted foresight. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025.
- [6] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [7] N. Hounsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. de Laroussilhe, A. Gesmundo, M. Atariyan, and S. Gelly. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning*, pages 2790–2799, 2019.
- [8] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017.
- [9] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba. Hindsight experience replay. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [10] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [11] P. Intelligence, A. Amin, R. Aniceto, A. Balakrishna, K. Black, K. Conley, G. Connors, J. Darpinian, K. Dhabalia, J. DiCarlo, et al.  $\pi_{0.6}$ : a VLA that learns from experience. *arXiv preprint arXiv:2511.14759*, 2025.
- [12] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu, and L. Fan. VIMA: Robot manipulation with multimodal prompts. In *International Conference on Machine Learning*, pages 14975–15022. PMLR, 2023.
- [13] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, et al. PaLM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- [14] 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.

- 308 [15] B. Zitkovich, T. Yu, S. Xu, P. Xu, T. Xiao, F. Xia, J. Wu, P. Wohlhart, S. Welker, A. Wahid,  
309 et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. In  
310 *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023.
- 311 [16] K. Bousmalis, G. Vezzani, D. Rao, C. Devin, A. X. Lee, et al. RoboCat: A self-improving  
312 generalist agent for robotic manipulation. *Transactions on Machine Learning Research*, 2024.  
313 URL <https://arxiv.org/abs/2306.11706>.
- 314 [17] A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta,  
315 A. Mandlekar, A. Jain, et al. Open X-embodiment: Robotic learning datasets and RT-X models:  
316 Open X-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and  
317 Automation (ICRA)*, pages 6892–6903. IEEE, 2024.
- 318 [18] X. Li, M. Liu, H. Zhang, C. Yu, J. Xu, H. Wu, C. Cheang, Y. Jing, W. Zhang, H. Liu, et al.  
319 Vision-language foundation models as effective robot imitators. In *International Conference  
320 on Learning Representations*, volume 2024, pages 26703–26721, 2024.
- 321 [19] H. Wu, Y. Jing, C. Cheang, G. Chen, J. Xu, X. Li, M. Liu, H. Li, and T. Kong. Unleashing large-  
322 scale video generative pre-training for visual robot manipulation. In *International Conference  
323 on Learning Representations*, volume 2024, pages 10641–10662, 2024.
- 324 [20] O. M. Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, et al. Octo: An open-source  
325 generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. URL [https://arxiv.org/  
326 abs/2405.12213](https://arxiv.org/abs/2405.12213).
- 327 [21] P. An, Y. Guo, X. Li, J. Liu, M. Liu, Z. Wang, et al. RoboMamba: Efficient vision-language-  
328 action model for robotic reasoning and manipulation. *Advances in Neural Information Pro-  
329 cessing Systems*, 2024.
- 330 [22] Q. Li, Y. Liang, Z. Wang, L. Luo, X. Chen, M. Liao, F. Wei, Y. Deng, S. Xu, Y. Zhang, et al.  
331 CogACT: A foundational vision-language-action model for synergizing cognition and action  
332 in robotic manipulation. *arXiv preprint arXiv:2411.19650*, 2024.
- 333 [23] D. Qu, H. Song, Q. Chen, Y. Yao, X. Ye, Y. Ding, et al. SpatialVLA: Exploring spatial  
334 representations for visual-language-action model. *arXiv preprint arXiv:2501.15830*, 2025.  
335 URL <https://arxiv.org/abs/2501.15830>.
- 336 [24] J. Wen, Y. Zhu, J. Li, M. Zhu, Z. Tang, K. Wu, Z. Xu, N. Liu, R. Cheng, C. Shen, et al.  
337 TinyVLA: Towards fast, data-efficient vision-language-action models for robotic manipula-  
338 tion. *IEEE Robotics and Automation Letters*, 2025.
- 339 [25] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, et al. OpenVLA:  
340 An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. URL  
341 <https://arxiv.org/abs/2406.09246>.
- 342 [26] T. Johannink, S. Bahl, A. Nair, J. Luo, A. Kumar, M. Loskyll, J. A. Ojea, E. Solowjow, and  
343 S. Levine. Residual reinforcement learning for robot control. In *2019 international conference  
344 on robotics and automation (ICRA)*, pages 6023–6029. IEEE, 2019.
- 345 [27] T. Silver, K. Allen, J. Tenenbaum, and L. Kaelbling. Residual policy learning. *arXiv preprint  
346 arXiv:1812.06298*, 2018.
- 347 [28] C. Xu, J. T. Springenberg, M. Equi, A. Amin, A. Esmail, S. Levine, and L. Ke. RL token: Boot-  
348 strapping online RL with vision-language-action models. *arXiv preprint arXiv:2604.23073*,  
349 2026.
- 350 [29] A. Ranjbar, N. A. Vien, H. Ziesche, J. Boedecker, and G. Neumann. Residual feedback learning  
351 for contact-rich manipulation tasks with uncertainty. In *2021 IEEE/RSJ International confer-  
352 ence on intelligent robots and systems (IROS)*, pages 2383–2390. IEEE, 2021.

- 353 [30] H. Bharadhwaj, R. Mottaghi, A. Gupta, and S. Tulsiani. Track2Act: Predicting point tracks  
354 from internet videos enables generalizable robot manipulation. In *European Conference on*  
355 *Computer Vision*, pages 306–324. Springer, 2024.
- 356 [31] L. Ankile, A. Simeonov, I. Shenfeld, M. Torne, and P. Agrawal. From imitation to refinement-  
357 residual RL for precise assembly. In *2025 IEEE International Conference on Robotics and*  
358 *Automation (ICRA)*, pages 01–08. IEEE, 2025.
- 359 [32] K. Pertsch, Y. Lee, and J. Lim. Accelerating reinforcement learning with learned skill priors.  
360 In *Conference on robot learning*, pages 188–204. PMLR, 2021.
- 361 [33] P. J. Ball, L. Smith, I. Kostrikov, and S. Levine. Efficient online reinforcement learning with  
362 offline data. In *International Conference on Machine Learning*, pages 1577–1594. PMLR,  
363 2023.
- 364 [34] J. Luo, Z. Hu, C. Xu, Y.L. Tan, J. Berg, A. Sharma, S. Schaal, C. Finn, A. Gupta, and S. Levine.  
365 SERL: A software suite for sample-efficient robotic reinforcement learning. In *2024 IEEE*  
366 *International Conference on Robotics and Automation (ICRA)*, pages 16961–16969. IEEE,  
367 2024.
- 368 [35] J. Luo, C. Xu, J. Wu, and S. Levine. Precise and dexterous robotic manipulation via human-  
369 in-the-loop reinforcement learning. *Science Robotics*, 10(105):eads5033, 2025.
- 370 [36] Z. Zhou, A. Peng, Q. Li, S. Levine, and A. Kumar. Efficient online reinforcement learning fine-  
371 tuning need not retain offline data. In *International Conference on Learning Representations*,  
372 volume 2025, pages 32343–32368, 2025.
- 373 [37] A. Ren, J. Lidard, L. Ankile, A. Simeonov, P. Agrawal, A. Majumdar, B. Burchfiel, H. Dai,  
374 and M. Simchowitz. Diffusion policy optimization. In *International Conference on*  
375 *Learning Representations*, volume 2025, pages 77288–77329, 2025.
- 376 [38] H. Jiang and Z. Yang. Adaptive diffusion policy optimization for robotic manipulation. *arXiv*  
377 *preprint arXiv:2505.08376*, 2025.
- 378 [39] X. Yuan, T. Mu, S. Tao, Y. Fang, M. Zhang, and H. Su. Policy decorator: Model-agnostic  
379 online refinement for large policy model. *arXiv preprint arXiv:2412.13630*, 2024. doi:10.  
380 48550/arXiv.2412.13630. URL <https://arxiv.org/abs/2412.13630>.
- 381 [40] Z. Song, G. Ouyang, M. Li, Y. Ji, C. Wang, Z. Xu, Z. Zhang, X. Zhang, Q. Jiang, F. Ji, et al.  
382 ManipLVM-R1: Reinforcement learning for reasoning in embodied manipulation with large  
383 vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
384 volume 40, pages 18558–18566, 2026.
- 385 [41] A. Cully, J. Clune, D. Tarapore, and J.-B. Mouret. Robots that can adapt like animals. *Nature*,  
386 521(7553):503–507, 2015.
- 387 [42] B. Thananjeyan, A. Balakrishna, S. Nair, M. Luo, K. Srinivasan, M. Hwang, J. E. Gonzalez,  
388 J. Ibarz, C. Finn, and K. Goldberg. Recovery RL: Safe reinforcement learning with learned  
389 recovery zones. *IEEE Robotics and Automation Letters*, 6(3):4915–4922, 2021.
- 390 [43] Y. Dai, J. Lee, N. Fazeli, and J. Chai. Racer: Rich language-guided failure recovery policies  
391 for imitation learning. *arXiv preprint arXiv:2409.14674*, 2024.
- 392 [44] Z. Lin, J. Duan, H. Fang, D. Fox, R. Krishna, C. Tan, and B. Wen. Failsafe: Reasoning and  
393 recovery from failures in vision-language-action models. *arXiv preprint arXiv:2510.01642*,  
394 2025.
- 395 [45] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
396 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- 397 [46] C. Yu, Y. Wang, Z. Guo, H. Lin, S. Xu, H. Zang, Q. Zhang, Y. Wu, C. Zhu, J. Hu, et al. Rlinf:  
398 Flexible and efficient large-scale reinforcement learning via macro-to-micro flow transforma-  
399 tion. *arXiv preprint arXiv:2509.15965*, 2025.
- 400 [47] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine,  
401 W. Ai, B. Martinez, H. Yin, M. Lingelbach, M. Hwang, A. Hiranaka, S. Garlanka, A. Ay-  
402 din, S. Lee, J. Sun, M. Anvari, M. Sharma, D. Bansal, S. Hunter, K.-Y. Kim, A. Lou, C. R.  
403 Matthews, I. Villa-Renteria, J. H. Tang, C. Tang, F. Xia, Y. Li, S. Savarese, H. Gweon, C. K.  
404 Liu, J. Wu, and L. Fei-Fei. BEHAVIOR-1K: A human-centered, embodied AI benchmark with  
405 1,000 everyday activities and realistic simulation. *arXiv preprint arXiv:2403.09227*, 2024.

406 **A Appendix: Additional Experimental Details**

407 **A.1 Additional results**

408 **Test on Behavior-1k challenge.** To evaluate whether the recovery design can be transferred  
 409 beyond the Fetch tabletop setup, we additionally implemented ReCoVLA on three tasks in the  
 410 Behavior-1K Challenge<sup>1</sup>: sorting vegetables, bringing in wood, and preparing a lunch box. These  
 411 tasks use a different simulated robot platform (R1 Pro) and require recovery from object drops,  
 412 constrained bimanual manipulation, and placement errors. Figure 6 shows representative recovery  
 413 trajectories. Across these examples, the VLM descriptor identifies the failed object or manipulation  
 414 stage, while the compiled residual reward encourages the local correction needed to return the rollout  
 415 to a solvable state. Table 1 shows the performance of ReCoVLA and Openpi Comet<sup>2</sup>, demonstrating  
 416 the advantage of our method.

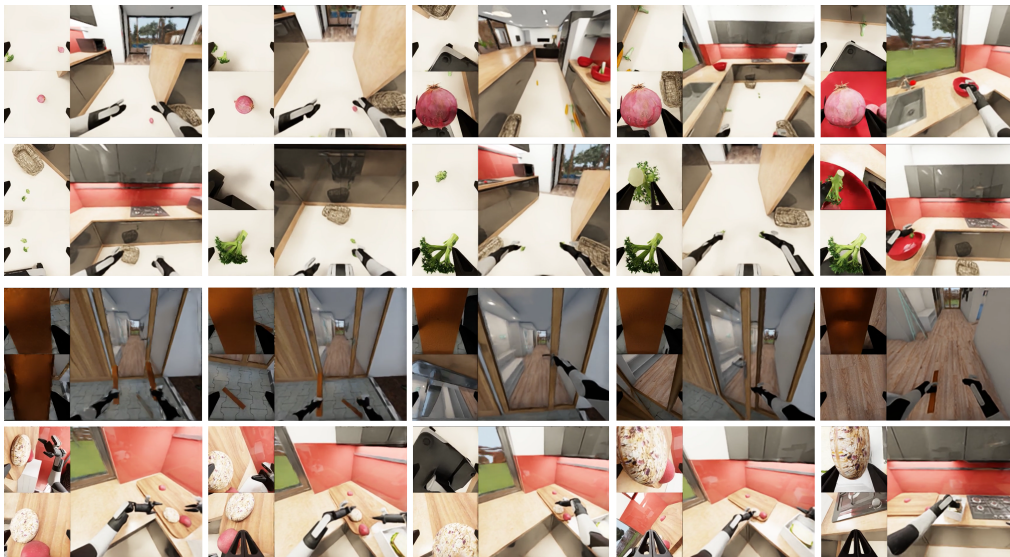


Figure 6: Failure recovery examples on Behavior-1K Challenge tasks. The top two rows show recovery from dropped onion and broccoli failures: the robot re-grasps the object from the ground and places it into the red bowl. The third row shows a door-opening failure in which both hands are occupied by wood sticks; the learned recovery drops one stick, opens the door, and continues the task. The final row shows recovery for picking up a bagel and placing it into the lunch box.

Table 1: Behavior-1K challenge test results. Each method is evaluated over 20 trials per task. Success rate and Q-score are reported in  $[0, 1]$ , where higher is better.

Task	Method	Trials	Success rate	Q-score
Sorting vegetables on the table	ReCoVLA	20	<b>0.20</b>	<b>0.51</b>
Sorting vegetables on the table	OpenPI Comet	20	0.05	0.26
Bringing in wood	ReCoVLA	20	<b>0.65</b>	<b>0.79</b>
Bringing in wood	OpenPI Comet	20	0.50	0.62
Preparing a lunch box	ReCoVLA	20	<b>0.25</b>	<b>0.49</b>
Preparing a lunch box	OpenPI Comet	20	0.10	0.33
Average over three tasks	ReCoVLA	–	<b>0.37</b>	<b>0.60</b>
Average over three tasks	OpenPI Comet	–	0.22	0.40

<sup>1</sup>B1K Challenge: <https://behavior.stanford.edu/index.html>

<sup>2</sup>A top-ranking VLA model in the B1K Challenge: <https://github.com/mli0603/openpi-comet>

Table 2: Additional comparison with existing baselines. Each method is evaluated over 20 trials per task. Success rate and Q-score are reported in  $[0, 1]$ , where higher is better.

Task	Method	Trials	Success rate	Q-score
Sorting vegetables	ReCoVLA	20	<b>0.65</b>	<b>0.82</b>
Sorting vegetables	RLinf	20	0.00	0.12
Sorting vegetables	RACER	20	0.35	0.49
Soda-can disposal	ReCoVLA	20	<b>0.75</b>	<b>0.88</b>
Soda-can disposal	RLinf	20	0.35	0.40
Soda-can disposal	RACER	20	0.65	0.74
Organizing toolbox	ReCoVLA	20	<b>0.60</b>	<b>0.78</b>
Organizing toolbox	RLinf	20	0.00	0.05
Organizing toolbox	RACER	20	0.30	0.56
Average over three tasks	ReCoVLA	–	<b>0.67</b>	<b>0.83</b>
Average over three tasks	RLinf	–	0.12	0.19
Average over three tasks	RACER	–	0.43	0.60

417 **Comparison with additional baselines.** We also compare ReCoVLA with two recent recovery  
 418 or VLA+RL baselines, RLinf [46] and RACER [43], on the same three simulation tasks. Table 2  
 419 reports 20 trials per method and task. Averaged over the three tasks, ReCoVLA obtains 0.67 success  
 420 and 0.83 Q-score, compared with 0.12/0.19 for RLinf and 0.43/0.60 for RACER. RACER performs  
 421 reasonably on the short-horizon soda-can task, but remains lower than ReCoVLA on all three tasks.  
 422 RLinf is more sensitive to the sparse and staged nature of the recovery objectives, particularly in  
 423 vegetable sorting and toolbox organization. These comparisons support the main result that failure-  
 424 conditioned, stage-gated reward compilation is important for robust recovery rather than merely  
 425 adding an RL correction to the base VLA.

## 426 A.2 Additional experiment setup details

427 **Additional physical experiment details.** All physical experiments use a Fetch mobile manipula-  
 428 tor with the task prompts listed in Table 3. In the vegetable-sorting task, we randomize vegetable  
 429 positions on the table while keeping the yellow bowl and red dishes fixed. In the soda-can disposal  
 430 task, we randomize soda-can positions and appearance while keeping the can size fixed. In the  
 431 toolbox organization task, we randomize the cable and camera-box positions on the ground. These  
 432 randomizations preserve the task semantics while changing the initial recovery geometry across  
 433 trials.

434 During physical experiments, we terminate a trial if the robot makes no measurable progress for a  
 435 3-minute window, defined as no increase in Q-score over that interval. We also terminate a trial if  
 436 the robot collides with the environment or an undesired object. The maximum execution time for  
 437 each task is set to twice the average duration of the corresponding human expert demonstrations.

438 **Q-score criteria.** For vegetable sorting, each pumpkin placed in the yellow bowl receives 0.15,  
 439 each broccoli piece placed in the middle red dish receives 0.10, and each corn placed in the right red  
 440 dish receives 0.15; full task completion receives a Q-score of 1.0. For soda-can disposal, each can  
 441 placed in the trash bin receives 0.30, with full completion assigned 1.0. For toolbox organization,  
 442 placing the cable or box in the toolbox receives 0.30 each, and closing the lid receives 0.40; full  
 443 completion receives 1.0.

444 **Simulation environment.** The simulation environments mirror the three layouts used in the phys-  
 445 ical Fetch experiments. The toolbox environment contains a table-mounted toolbox and ground  
 446 objects that must be picked up, inserted into the toolbox, and followed by lid closure. The sort-  
 447 ing environment contains red dishes, a yellow bowl, and multiple vegetable categories that must be  
 448 placed into their specified containers. The trash-can environment contains soda cans on the table

Table 3: Evaluation tasks and language prompts.

Task	Purpose	Language prompt
Organizing toolbox	Contact rich	Pick up the cable and box on the ground and place them in the toolbox on the table, and close the lid of the toolbox.
Sorting vegetables	Long horizon	Sort the vegetables into the red dishes and yellow bowl on the table: put all three broccoli pieces into one red dish; put two corns into the other red dish; and put both pumpkin pieces into the yellow bowl.
Picking up trash cans	Short horizon	Put the three cans of soda from the table inside the trash can on the left of the table.

449 and a trash can placed beside the table. These layouts instantiate the same recovery regimes shown  
 450 in Figure 3: contact-rich recovery, long-horizon object-category sorting, and short-horizon disposal.

451 During simulation rollouts, the external VLM failure analyzer is queried on the current RGB obser-  
 452 vation and recovery prompt. The resulting descriptor records the failure type, recovery stage, active  
 453 entities, confidence, and reward mask. In the toolbox task, typical detected failure states include a  
 454 cable or box left on the ground, an object grasped but not inserted into the toolbox, and a toolbox  
 455 lid that remains open after object placement. In the sorting task, typical failures include vegetables  
 456 placed in the wrong receptacle, remaining objects on the table, and incomplete category-specific  
 457 placement. In the trash-can task, typical failures include a can remaining on the table, a dropped  
 458 can near the bin, or a can held by the gripper but not yet released into the trash can. These detected  
 459 states determine which entities and reward-library components are activated by the reward compiler  
 460 during residual policy optimization.

### 461 A.3 Reward compiler implementation and VLM failure detector confusion matrix

462 **Implementation details of reward compiler.** Algorithm 1 gives the implementation-level reward  
 463 compilation procedure used in our experiments. Beyond the high-level description in the main  
 464 paper, the implementation performs conservative validation before any residual reward is produced.  
 465 It normalizes the VLM category and stage labels to the task vocabulary, checks that the reward mask  
 466 has length  $K$ , rejects descriptors with unresolved entities, and rejects descriptors whose selected  
 467 components cannot be bound to the required semantic roles. This rejection behavior is intentional:  
 468 if the descriptor is incomplete or inconsistent with the simulator object map, the system does not  
 469 synthesize a fallback reward. Instead, the dispatcher leaves the nominal VLA active or invokes no  
 470 recovery policy for that descriptor. The algorithm also makes the component-level binding explicit.  
 471 The active entity set  $E_c$  may contain extra objects observed by the VLM, but each reward potential  
 472  $\varphi_k$  receives only the entities required by its signature  $\sigma_k$ . This avoids, for example, passing a  
 473 receptacle entity to a grasp term or an end-effector entity to an articulation term. Stage gates are  
 474 retrieved only after role binding, so each gate is parameterized by the same grounded entities used  
 475 by the corresponding reward term.

476 **Quantitative results of experiments.** Table 4 provides the full numerical results corresponding  
 477 to the main-paper plots, including simulation, physical robot, and OOD settings. The averages  
 478 show that M4, the full ReCoVLA instantiation on  $\pi_{0.5}$ , is the strongest variant in both simulation  
 479 and physical deployment. In simulation, M4 improves the average result from 0.37/0.56 for the  
 480 no-recovery  $\pi_{0.5}$  baseline M1 to 0.67/0.83. On the physical robot, M4 improves the average from  
 481 0.27/0.40 to 0.62/0.75. The OOD rows show the same trend under object substitutions: M4 achieves  
 482 0.53/0.65 on average, while M1 reaches only 0.10/0.22. The OpenVLA comparison further indicates  
 483 that the recovery design is not tied to a single base VLA, with M6 improving over M5 in both  
 484 simulation and physical experiments.

485 **VLM failure-analysis accuracy.** Figure 7 reports a row-normalized failure-mode confusion ma-  
 486 trix for the external VLM detector used during physical rollouts. The matrix covers ten recoverable

---

**Algorithm 1** Deterministic reward compilation

---

**Require:** VLM descriptor  $\xi_c = (c, z, E_c, \rho, m_c)$ ; task specification  $\mathcal{T}$ ; simulator object map  $\mathcal{O}$ ; reward library  $\mathcal{R} = \{\varphi_k\}_{k=1}^K$ ; component signatures  $\{\sigma_k\}_{k=1}^K$ ; stage-gate templates  $\mathcal{G}$ ; recoverable catalog  $\mathcal{C}_{\text{train}}$ ; confidence threshold  $\tau$

**Ensure:** Executable reward  $R_{\text{M4}}^c(s_t, a_t, s_{t+1})$ , or rejection

- 1: Normalize category and stage labels:  
 $(\hat{c}, \hat{z}) \leftarrow \text{NORMALIZELABELS}(c, z, \mathcal{T})$
- 2: **if**  $\hat{c} \notin \mathcal{C}_{\text{train}}$  **or**  $\rho < \tau$  **or**  $m_c \notin \{0, 1\}^K$  **then**
- 3:     Reject descriptor and return no recovery reward.
- 4: **end if**
- 5: Canonicalize VLM entities:  
 $\hat{E}_c \leftarrow \text{RESOLVEENTITIES}(E_c, \mathcal{T}, \mathcal{O})$
- 6: **if**  $\hat{E}_c$  contains unresolved entities **then**
- 7:     Reject descriptor and return no recovery reward.
- 8: **end if**
- 9: Assign semantic roles:  
 $\mathcal{R}_c \leftarrow \text{ASSIGNROLES}(\hat{E}_c, \hat{c}, \hat{z}, \mathcal{T})$   
(e.g., end-effector, object, target, source, articulated part)
- 10: Initialize selected term set  $\mathcal{S}_c \leftarrow \emptyset$ .
- 11: **for**  $k = 1, \dots, K$  **do**
- 12:     **if**  $m_c^{(k)} = 1$  **then**
- 13:         Read the component signature  $\sigma_k$  for  $\varphi_k$ .
- 14:         Bind component-specific entities:  
 $\bar{E}_{c,k} \leftarrow \text{BINDROLES}(\mathcal{R}_c, \sigma_k)$
- 15:         **if**  $\bar{E}_{c,k}$  is missing any required role **then**
- 16:             Reject descriptor and return no recovery reward.
- 17:         **end if**
- 18:         Retrieve the fixed stage gate:  
 $g_{\hat{c},k}(\cdot) \leftarrow \mathcal{G}(\hat{c}, \hat{z}, k, \bar{E}_{c,k})$
- 19:         Add the gated potential-difference term:  
 $\mathcal{S}_c \leftarrow \mathcal{S}_c \cup \{g_{\hat{c},k}(s_t; \bar{E}_{c,k})\Delta\varphi_k(s_t, s_{t+1}; \bar{E}_{c,k})\}$
- 20:     **end if**
- 21: **end for**
- 22: **if**  $\mathcal{S}_c = \emptyset$  **then**
- 23:     Reject descriptor and return no recovery reward.
- 24: **end if**
- 25: Return the compiled reward:  
 $R_{\text{M4}}^c(s_t, a_t, s_{t+1}) = \sum_{\psi \in \mathcal{S}_c} \psi - \Omega_t(a_t^r)$

---

487 failure modes: five from vegetable sorting, two from soda-can disposal, and three from toolbox  
488 organization. Because the active task is known from the language instruction, the evaluation con-  
489 strains predictions to the corresponding task group; errors therefore measure confusion between  
490 semantically related failures within the same task, rather than impossible cross-task predictions. For  
491 vegetable sorting, the observed failure occurrences are 9 wrong-broccoli-receptacle, 8 wrong-corn-  
492 receptacle, 7 wrong-pumpkin-receptacle, 8 dropped-vegetable-on-table, and 8 dropped-vegetable-  
493 on-ground cases, totaling 40 detections with 87.5% aggregate accuracy, 87.4% macro accuracy, and  
494 1.0 percentage-point class-wise standard deviation. For soda-can disposal, the detector observes  
495 20 fallen-can-on-table and 18 dropped-can-on-ground cases, totaling 38 detections with 86.8% ag-  
496 gregate accuracy, 86.7% macro accuracy, and 3.3 percentage-point standard deviation. For toolbox  
497 organization, the detector observes 12 dropped-box, 14 dropped-cable, and 13 unclosed-lid cases,  
498 totaling 39 detections with 79.5% aggregate accuracy, 80.0% macro accuracy, and 8.5 percentage-  
499 point standard deviation. Across all tasks, the detector reaches 85.0% macro accuracy, 84.6% ag-  
500 gregate accuracy, and 6.0 percentage-point class-wise standard deviation. Most errors are local,  
501 such as confusion between wrong vegetable receptacles, table-versus-ground soda-can failures, or  
502 cable/box/lid failures in the toolbox task. This supports the use of the VLM as a semantic failure

Table 4: Failure-recovery experiment results in simulation, physical robot, and out-of-distribution (OOD) settings. Each entry reports success rate / Q-score. Higher is better.

Setting	Task / split	M1	M2	M3	M4 (Ours)	M5	M6
<i>Average over three in-distribution tasks</i>							
Simulation	Average	0.37 / 0.56	0.40 / 0.55	0.48 / 0.63	<b>0.67 / 0.83</b>	0.23 / 0.33	0.45 / 0.55
Physical robot	Average	0.27 / 0.40	0.27 / 0.40	0.40 / 0.51	<b>0.62 / 0.75</b>	0.15 / 0.25	0.35 / 0.42
<i>In-distribution task results</i>							
Simulation	Sorting vegetables	0.30 / 0.48	0.20 / 0.28	0.45 / 0.55	<b>0.65 / 0.82</b>	0.20 / 0.30	0.45 / 0.52
Simulation	Picking up trash cans	0.55 / 0.67	0.70 / 0.79	0.60 / 0.71	<b>0.75 / 0.88</b>	0.40 / 0.45	0.60 / 0.68
Simulation	Organizing toolbox	0.25 / 0.54	0.30 / 0.59	0.40 / 0.64	<b>0.60 / 0.78</b>	0.10 / 0.23	0.30 / 0.45
Physical robot	Sorting vegetables	0.25 / 0.36	0.10 / 0.20	0.45 / 0.50	<b>0.60 / 0.72</b>	0.10 / 0.20	0.40 / 0.48
Physical robot	Picking up trash cans	0.45 / 0.54	0.55 / 0.63	0.45 / 0.53	<b>0.75 / 0.83</b>	0.35 / 0.40	0.50 / 0.56
Physical robot	Organizing toolbox	0.10 / 0.30	0.15 / 0.38	0.30 / 0.49	<b>0.50 / 0.69</b>	0.00 / 0.15	0.15 / 0.23
<i>OOD stress tests</i>							
OOD (physical robot)	Sorting vegetables pumpkin → tomato	0.00 / 0.10	–	–	<b>0.50 / 0.65</b>	–	–
OOD (physical robot)	Picking up trash cans soda can → tall soda can	0.30 / 0.50	–	–	<b>0.70 / 0.76</b>	–	–
OOD (physical robot)	Organizing toolbox cable → tape	0.00 / 0.05	–	–	<b>0.40 / 0.55</b>	–	–

M1: fine-tuned  $\pi_{0.5}$  without recovery. M2: task-level residual reward. M3: equal-weight failure rewards. M4: proposed  $\pi_{0.5}$  VLM reward compiler. M5: fine-tuned OpenVLA without recovery. M6: OpenVLA with the proposed recovery method.

Bold entries denote the best result in each row. OOD tests compare M4 against the corresponding non-recovery baseline M1.

503 dispatcher, while also highlighting that detector mistakes can still select a suboptimal residual policy  
504 in visually similar recovery states.

#### 505 A.4 Training details

506 **Base VLA fine-tuning.** For each task, we collected 40 expert demonstrations: 20 in simulation  
507 and 20 on the physical robot. We use these demonstrations to fine-tune the  $\pi_{0.5}$  base model used  
508 by M1–M4. The fine-tuned policy is then used as the nominal controller for all variants built on  
509 the same backbone. During residual recovery training, the base VLA is frozen, so no gradients are  
510 propagated into the VLA encoder, language model, or action expert. The same process is applied to  
511 fine-tune the OpenVLA base. The main fine-tuning parameters are listed in Table 5.

512 **Residual policy training.** Residual policies are trained entirely in simulation with PPO, using the  
513 same simulator task layouts and recoverable failure categories described in Section A.2. For each  
514 recoverable failure category, rollouts are initialized from a restored OmniGibson scene state that  
515 reproduces the corresponding failure, such as an object in the wrong receptacle, a dropped soda can,  
516 a ground object near the toolbox, or an unclosed toolbox lid. The actor receives the frozen VLA  
517 latent feature  $h_t$  and outputs an additive residual action  $a_t^r$ . The executed action is the clipped sum  
518 of the base action and residual action, as defined in Eq. 1. After training, the residual actor is frozen  
519 and stored in the recovery policy library indexed by the VLM-detected failure category. PPO and  
520 residual-action parameters are listed in Table 5.

521 **Reward and regularization.** The reward used for each residual rollout is generated by the deter-  
522 ministic compiler in Algorithm 1. M2 uses a task-level reward mask, M3 uses the VLM-selected  
523 failure mask without stage gates, and M4/M6 use the same failure mask with stage-aware gates. The  
524 implementation also enforces per-joint action bounds from the simulator action space. In addition  
525 to the residual penalty  $\Omega_t(a_t^r)$ , the implementation applies smoothness penalties on action deltas  
526 and sign-flip reversals, with coefficients reported in Table 5. This regularization is important during  
527 sim-to-real transfer because it biases the residual policy toward small corrective motions around the  
528 nominal VLA behavior rather than replacing the base controller.

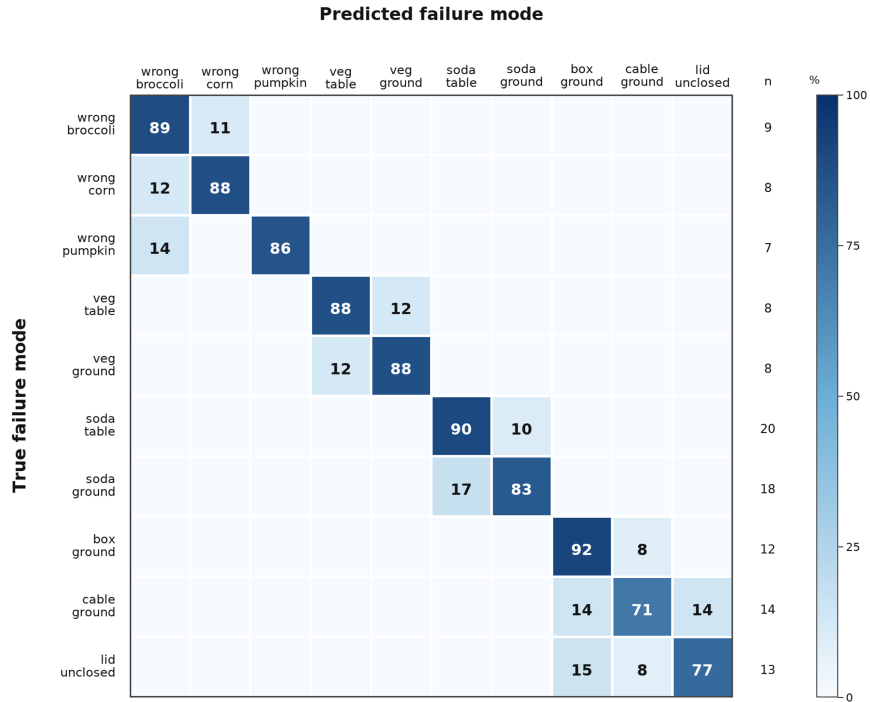


Figure 7: VLM failure detector confusion matrix. Rows are normalized by the true failure mode, and the support column reports the number of physical-rollout occurrences used for each row.

Table 5: Training parameters for base-VLA fine-tuning and residual PPO.

Stage	Parameter	Value
<i><math>\pi_{0.5}</math> fine-tuning</i>		
$\pi_{0.5}$	Expert demonstrations per task	40 total: 20 simulation, 20 physical robot
$\pi_{0.5}$	Base checkpoint	gs://openpi-assets/checkpoints/pi05_base/params
$\pi_{0.5}$	Action horizon	32
$\pi_{0.5}$	Training steps	20000 gradient steps
$\pi_{0.5}$	Batch size	4
$\pi_{0.5}$	Learning-rate schedule	Cosine decay
$\pi_{0.5}$	Peak learning rate	$2.5 \times 10^{-5}$
$\pi_{0.5}$	Warmup steps	$0.1 \times$ training steps, capped at 1000
$\pi_{0.5}$	Seed	42
$\pi_{0.5}$	Freeze rule	OpenPI $\pi_{0.5}$ model freeze filter
<i>Residual PPO</i>		
PPO	Policy input	Frozen $\pi_{0.5}$ VLA latent feature $h_t$
PPO	Actor/critic architecture	MLP, hidden sizes (4096, 1024, 128)
PPO	Activation	Tanh
PPO	Total timesteps per residual policy	300,000
PPO	Vectorized environments	1
PPO	Episode horizon	1000 steps
PPO	Simulator steps per residual action	1
PPO	PPO rollout length $n_{steps}$	100
PPO	Discount factor $\gamma$	0.99
PPO	GAE parameter	0.95
PPO	PPO clip range	0.2
PPO	PPO learning rate	$3 \times 10^{-4}$
PPO	Residual action clip	$[-1, 1]$ before scaling
PPO	Default residual scale	0.25
PPO	Mobile-base residual scale	0.03
PPO	Torso residual scale	0.01
PPO	Startup gripper-open duration	8 environment steps
PPO	Action-delta penalty coefficient	0.005
PPO	Action sign-flip penalty coefficient	0.01

529 **Deployment.** The trained residual policies are deployed zero-shot on the physical Fetch robot.  
530 The robot executes the frozen base VLA by default. At a fixed monitoring interval, the VLM failure  
531 detector receives the recent RGB observation history and task prompt. If the detector predicts a  
532 known failure category with confidence above the deployment threshold, the corresponding residual  
533 policy is activated; otherwise the robot continues with the base policy. The residual policy is there-  
534 fore used only as a targeted recovery controller for known failure states, not as a replacement for  
535 nominal task execution.