
Recovery-Aware HIL-SERL for Contact-Rich USB Insertion with Global State Encoding

Anonymous Author(s)

Affiliation

Address

email

Abstract

We study human-in-the-loop real-world reinforcement learning for USB pick-and-insert, a contact-rich manipulation task with tight tolerances and failure modes dominated by contact transitions. Building on HIL-SERL, we deploy an end-to-end on-robot training loop on a UR5e platform and introduce three system-level enhancements that target partial observability and recovery from non-terminal failures: (i) a practical high-rate compliant execution stack for stable in-contact exploration, (ii) global state encoding from an external camera to improve phase awareness under sparse rewards, and (iii) a human-gated **Recovery Buffer** to amplify recovery-critical experience. On real hardware, the global observer accelerates convergence to 100% success from 20k steps (3.5 h) to 12k steps (2 h), and the Recovery Buffer improves recovery success in two representative failure scenarios. We open-source our code and provide a project website with videos and additional details; the link will be released upon acceptance.

1 Introduction

USB connector insertion is a high-frequency and critical operation in electronics manufacturing, automated testing, and equipment maintenance. Improving its automation directly affects throughput, labor cost, and quality consistency. Despite standardized geometry and broad compatibility, USB insertion remains *contact-rich*: it involves tight clearances, strong contact constraints, and sensitivity to interaction forces. Consequently, small pose errors or inappropriate contact forces can cause insertion failure and may damage the connector or the robot end-effector.

Conventional automation pipelines typically rely on precision fixturing, manually tuned impedance/force controllers, or rule-based search strategies. In practice, however, deployment must cope with mounting offsets, compliance of cables and terminals, friction variability, and noisy observations, which makes rule-based solutions brittle and expensive to iterate. Reinforcement learning (RL) offers an attractive alternative by learning control policies through trial-and-error interaction, and has shown strong performance in continuous control and fine manipulation, with representative methods such as PPO and SAC [1, 2]. Yet, for contact-rich tasks, simulation training is highly sensitive to modeling errors in contact dynamics and friction and can suffer from sim-to-real transfer issues. Training directly in the real world removes the transfer gap, but it raises practical challenges: low sample efficiency, high data-collection cost, and non-trivial hardware risks during exploration.

Human-in-the-loop RL (HIL-RL) mitigates these challenges by combining human demonstrations and online intervention with on-robot fine-tuning. Recent progress such as HIL-SERL shows that high-performing real-robot policies can be learned within practical training time [3]. However, when instantiated for USB insertion on a physical UR5e system, we find that two practical bottlenecks remain under-addressed. **First**, the task is *partially observable and strongly temporal*: relying

37 predominantly on local proprioception and end-effector signals often provides insufficient modeling
38 capacity to disambiguate insertion phases and contact transitions, slowing learning and increasing
39 wall-clock cost. **Second**, robust deployment hinges on *learning from mistakes*: exploration yields
40 abundant non-terminal failures (e.g., jamming or misalignment), yet recovery-critical transitions are
41 typically rare and heterogeneous (many-class, few-shot) and are thus easily diluted in generic replay,
42 resulting in weak recovery behaviors and poor robustness to disturbances.

43 To address these bottlenecks, we study USB2.0 pick-and-insert on a UR5e manipulator and develop
44 an enhanced HIL-RL system that targets *execution safety*, *phase-aware state modeling*, and *recovery-*
45 *centric learning* while maintaining the HIL-SERL training paradigm. Our contributions are:

- 46 • **UR5e deployment with high-rate compliant execution.** We adapt the HIL-SERL baseline
47 to UR5e and implement an end-to-end on-robot training loop with high-frequency command
48 streaming and compliant execution based on PID force/torque regulation, improving stability
49 and feasibility for in-contact exploration.
- 50 • **Global stage-aware state encoding for improved task modeling.** We introduce a stage-
51 aware global observer that learns compact task-phase representations from an external
52 camera and fuses them into the policy input, reducing phase ambiguity and accelerating
53 convergence.
- 54 • **Human-gated Recovery Buffer for recovery-centric learning.** We add a dedicated
55 Recovery Buffer curated by human gating of recovery-relevant transitions and bias training
56 toward these samples, improving recovery robustness under typical failure states.

57 The remainder of this paper is organized as follows. Section 2 reviews related work on contact-rich
58 insertion, real-world reinforcement learning, and human-in-the-loop training. Section 3 presents
59 our framework and the three enhancements above. Section 4 evaluates the approach on a UR5e-
60 based USB2.0 pick-and-insert setup, reporting training efficiency, success rate, and robustness under
61 representative failure cases. Finally, Section 5 concludes with discussion and future directions.

62 2 Related Work

63 USB insertion falls into *contact-rich* assembly (e.g., peg-in-hole and connector mating), where
64 success depends on tight tolerances, frictional contact transitions, and safe interaction forces. Rather
65 than enumerating methods, we summarize the central trade-offs that motivate our design along three
66 axes consistent with our system (Section 3): (i) compliant execution for safe contact interaction, (ii)
67 observation modeling for phase ambiguity, and (iii) recovery-centric data mechanisms for learning
68 from mistakes.

69 **Compliant assembly and perception-driven insertion.** Industrial insertion is often solved with
70 precise fixturing plus impedance/admittance control and handcrafted search heuristics. These ap-
71 proaches can be reliable when uncertainty is bounded, but performance often degrades under mounting
72 offsets, friction variation, and unmodeled compliance. Perception- and servoing-based pipelines im-
73 prove robustness by grounding insertion in geometric feedback, especially when socket pose changes,
74 but they typically require careful sensing, calibration, and pipeline engineering [4, 5]. Learning-based
75 approaches aim to replace or complement heuristics by extracting contact behaviors from interaction
76 data; prior work shows RL can acquire effective insertion strategies, particularly when paired with
77 compliant/force regulation [6–8]. A recurring lesson is that robust contact-rich performance depends
78 as much on system design—including control interfaces, safety constraints, and reset cost—as on the
79 learning algorithm itself [9]. Our high-rate compliant controller is aligned with this system-first view,
80 providing the execution foundation required for safe real-world exploration.

81 **Real-world RL and human-in-the-loop training.** Real-world RL avoids sim-to-real mismatch
82 for contact dynamics, but it makes safe exploration, fast iteration, and robust infrastructure essential.
83 System efforts such as SERL provide an end-to-end software stack for sample-efficient real-robot RL
84 and demonstrate strong performance on contact-rich tasks within short training budgets [10]. HIL-
85 SERL builds on this direction by incorporating human supervision and data collection strategies to
86 further improve efficiency and safety [3]. Complementary work constrains in-contact exploration via
87 shared control and interaction templates, often improving stability but requiring task-specific interface

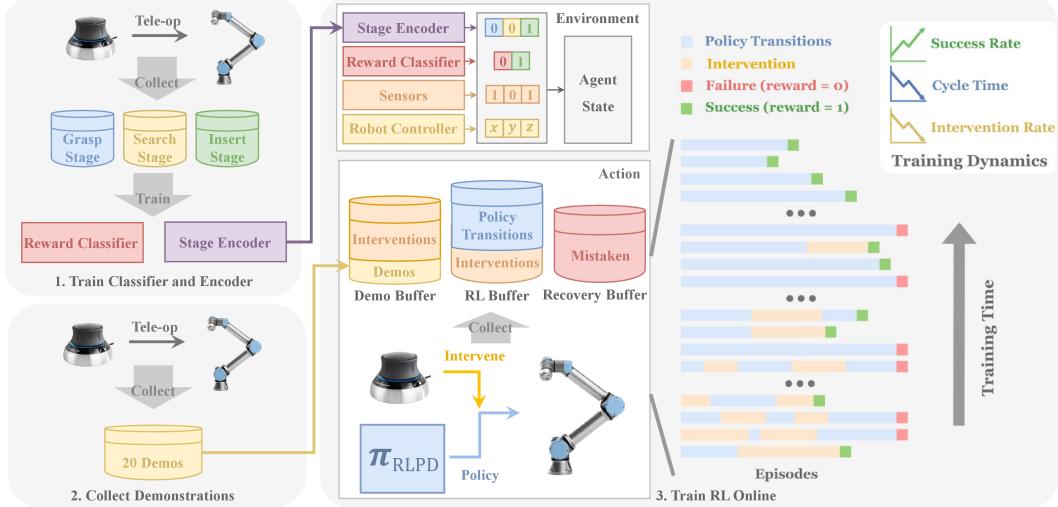


Figure 1: **System overview.** We adopt HIL-SERL as the baseline and augment it with three orthogonal, system-level improvements tailored for real-world USB2.0 pick-and-insert on UR5e. **Execution:** a high-rate compliant controller (RTDE, up to 500 Hz) enables smooth and safe in-contact exploration. **Observation:** a stage-aware global observer encodes external camera images into a compact task-phase representation and fuses it into the policy state. **Data curation:** a human-gated Recovery Buffer explicitly preserves recovery-critical transitions and is sampled together with the Demo Buffer and RL Buffer to strengthen learning from mistakes.

assumptions [11]. In USB insertion, distinct temporal phases (approach, contact, alignment, insertion) amplify partial observability; this motivates augmenting policy inputs with global, phase-aware representations rather than relying solely on local proprioceptive and end-effector signals.

Learning from demonstrations, feedback, and recovery. To reduce exploration burden, learning-from-demonstrations methods integrate expert trajectories into RL through replay and auxiliary objectives [12–14]. Human feedback has also been used to shape behavior online, ranging from evaluative reinforcement (e.g., TAMER) to corrective feedback schemes for deep policies [15, 16]. In contact-rich insertion, however, demonstration-augmented RL can still under-emphasize “near-failure” recovery, and generic replay may not allocate sufficient capacity to sparse recovery transitions. Safety-oriented formulations such as Recovery RL separate task learning from recovery regions/policies to avoid constraint violations [17], but may add algorithmic complexity and require explicit constraint specification. In contrast, our recovery-aware design is lightweight and data-centric: we curate recovery-critical transitions via human gating and bias sampling toward them, strengthening practical robustness without substantially increasing system complexity.

3 Method

Our goal is fast, efficient, and safe real-world RL for contact-rich USB insertion. We start from HIL-SERL [3], which combines sparse reward learning, human interventions, and off-policy training with prior data. We then introduce three targeted enhancements that address the two bottlenecks in Section 2: (i) reliable high-rate compliant execution for safe exploration, (ii) improved task modeling under partial observability via global phase encoding, and (iii) recovery-centric sampling to amplify learning from mistakes.

3.1 HIL-SERL Baseline

We build on Human-in-the-Loop Sample-Efficient Robotic Reinforcement Learning (HIL-SERL) [3], a real-world, vision-based RL system designed to achieve practical training times while maintaining safe on-robot exploration. HIL-SERL uses a sparse reward defined by a learned binary success classifier: the robot is first teleoperated to collect positive and negative examples, which are used to

114 train a reward detector offline; a small set of human demonstrations is then collected and inserted
115 into an offline *Demo Buffer* to bootstrap policy learning [3]. During online training, the classifier
116 provides the sparse reward signal, while a human operator can intervene via teleoperation (e.g., a
117 SpaceMouse) to correct unsafe or unproductive behaviors; interventions are typically more frequent
118 early in training and are reduced as success rate and cycle time improve [3].

119 At the system level, HIL-SERL follows an asynchronous actor–learner design. The actor executes
120 the latest policy on the robot, collects transitions, and streams experience to replay, while the learner
121 samples minibatches and periodically pushes updated parameters back to the actor [3]. Experience
122 is stored in two buffers: the Demo Buffer containing offline human trajectories and an *RL Buffer*
123 containing online interaction data. The learner samples *evenly* from these buffers and optimizes
124 the policy with RL with Prior Data (RLPD), instantiated with an off-policy backbone (SAC) [3].
125 Importantly, intervention segments are treated as privileged corrections: human-controlled transitions
126 are stored in both buffers, while policy transitions around interventions are stored in the RL Buffer,
127 increasing the utility of corrective data for off-policy learning [3]. We adopt this pipeline as our
128 baseline and introduce the following adaptations for USB insertion.

129 3.2 High-Rate Compliant Controller on UR5e

130 Contact-rich insertion requires not only sample-efficient learning but also *safe* and *responsive* execu-
131 tion: the control stack must react quickly to contact transitions while remaining compliant to reduce
132 impact forces and hardware risk. To support real-world HIL-RL on UR5e, we implement a high-rate
133 compliant control interface over RTDE (up to 500 Hz). Rather than purely position-based servoing,
134 we employ a PID force/torque regulation layer, which trades some absolute tracking precision for a
135 larger safety margin and stronger disturbance rejection during in-contact exploration. In addition, we
136 drive control via *relative* pose errors to reduce the chance that a policy exploits absolute workspace
137 coordinates.

138 Concretely, we maintain a global target end-effector pose $(\hat{\mathbf{x}}, \hat{\varphi})$, where $\mathbf{x} \in \mathbb{R}^3$ and $\varphi \in SO(3)$. A
139 lightweight HTTP server asynchronously receives action commands from the agent and updates only
140 this target pose without blocking. A separate high-rate control loop reads the current end-effector
141 pose (\mathbf{x}, φ) and computes pose errors

$$\delta\mathbf{x} = \hat{\mathbf{x}} - \mathbf{x}, \quad \delta\varphi = \hat{\varphi} \cdot \varphi^{-1}, \quad (1)$$

142 where rotational error is computed using $SO(3)$ group operations for geometric consistency. We then
143 generate the desired wrench via a PD law with discrete-time differencing,

$$\mathbf{F} = k_p \delta\mathbf{x} + k_d \Delta(\delta\mathbf{x}), \quad \mathbf{M} = k_p \delta\varphi + k_d \Delta(\delta\varphi), \quad (2)$$

144 and apply saturation before sending commands to the robot to enforce safety limits. Since the policy
145 operates at 10 Hz, we interpolate each action into sub-steps and update the target pose at 500 Hz,
146 yielding smoother trajectories and more compliant contact behavior. This controller completes the
147 execution layer needed for stable and safe real-world training.

148 3.3 Global State Encoding via a Stage-Aware Observer

149 HIL-SERL learns with a sparse reward defined by a success classifier. While this design is general, it
150 can be inefficient for long-horizon, multi-phase tasks, where informative feedback arrives mainly at
151 task completion and credit assignment across contact transitions becomes difficult. To improve task
152 modeling without changing the sparse-reward paradigm, we introduce a *stage-aware global observer*
153 that encodes external camera images into a compact representation of task phase (Fig. 2).

154 **Stage classifier training.** In parallel with training the binary reward classifier, we train an auxiliary
155 stage classifier using global images. For USB insertion, we discretize the task into four stages: *grasp*,
156 *approach*, *hole-search*, and *insertion*. The classifier encodes the global observation with a ResNet-18
157 backbone [18] and predicts the stage label via an MLP head optimized with cross-entropy. We use
158 the penultimate feature vector as a compact stage representation that summarizes global task context.

159 **Policy integration (stop-gradient).** During RL, the stage encoder is frozen and used only for
160 inference: it maps each global image to a 256-dimensional stage representation, which is concatenated
161 to the agent observation. We match this dimension to the baseline visual feature size (256) to keep
162 the representation balanced and avoid over-weighting any single modality.

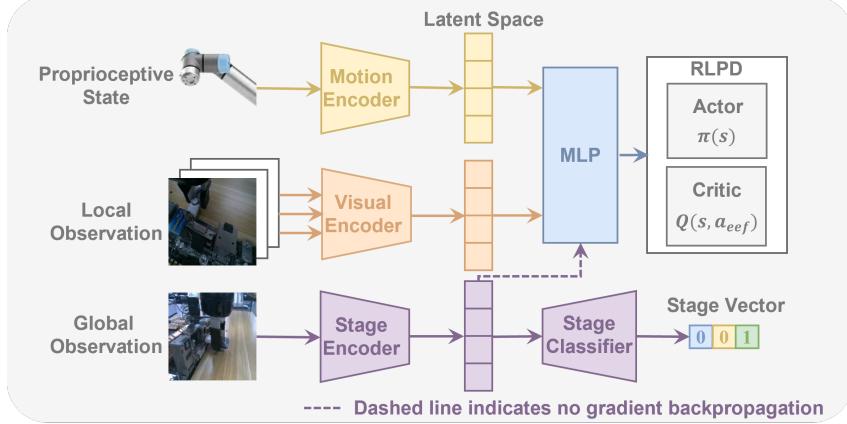


Figure 2: **Global stage-aware state encoding.** We augment the policy observation with a global stage representation extracted from an external camera. A ResNet-18 encoder and an MLP head are trained with cross-entropy to classify four phases (grasp, approach, hole-search, insertion). During RL, we freeze the encoder (stop-gradient) and concatenate its 256-d representation to the policy input to improve phase awareness and temporal task modeling.

163 **Practical impact.** Adding the global observer improves phase awareness and temporal modeling.
 164 In our UR5e USB insertion setup, this global encoding yields substantially faster convergence (about
 165 1.7 \times in wall-clock training time) and improves final success (Section 4).

166 3.4 Recovery Buffer for Learning from Mistakes

167 Motivated by how humans learn, we explicitly emphasize mistakes: in contact-rich manipulation,
 168 recovery from errors often provides more informative supervision than trajectories that simply
 169 succeed. HIL-SERL already highlights corrective behavior by prioritizing intervention segments and
 170 storing human-takeover transitions in both the Demo Buffer and RL Buffer [3]. However, in USB
 171 insertion we find that mistakes are *many-class, few-shot*: there are many distinct error types, but each
 172 type appears with limited frequency. Simply mixing these scarce recovery transitions into generic
 173 replay therefore dilutes their sampling probability, limiting systematic learning of robust recovery
 174 behaviors.

175 We address this with a dedicated **Recovery Buffer** as an explicit “mistake notebook.” At the end
 176 of each training episode, a supervisor decides whether the collected transitions contain meaningful
 177 recovery behavior (i.e., transitions that correct an error state); if so, the transitions are added to the
 178 Recovery Buffer. During learning, we sample from the Demo Buffer, RL Buffer, and Recovery Buffer
 179 with a fixed ratio of 4:4:2, increasing exposure to recovery-critical experience without changing the
 180 underlying RL objective. This lightweight data curation improves recovery behaviors and robustness
 181 in our experiments (Section 4).

182 4 Experiments

183 We evaluate our system on real-world USB2.0 pick-and-insert with UR5e. The experiments are
 184 designed to answer two questions aligned with our method: **(Q1)** does the global stage-aware
 185 observer improve sample efficiency and final performance under sparse reward? **(Q2)** does the
 186 human-gated Recovery Buffer improve recovery behaviors under typical failure states? Unless
 187 otherwise stated, all methods share the same execution stack (Section 3.2) and the HIL-SERL training
 188 pipeline (Section 3.1).

189 4.1 Platform, Task, and Data Collection

190 **Hardware setup.** Our platform is a UR5e manipulator equipped with a wrist force/torque sensor
 191 and a Robotiq two-finger gripper. We use three Intel RealSense D405 cameras: two wrist-mounted

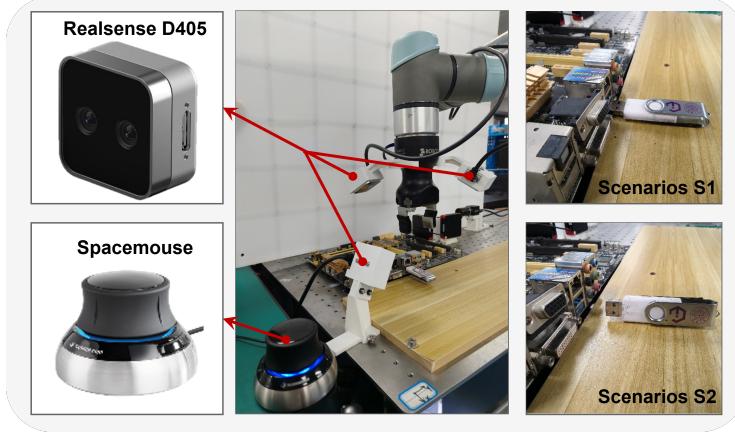


Figure 3: **Platform and failure scenarios.** **Left:** UR5e USB insertion setup with RealSense D405 cameras and SpaceMouse teleoperation. **Right:** recovery test initializations: **S1** (top) hard-contact under the port; **S2** (bottom) misoriented USB pickup (sideways).

192 cameras and one side-view camera for global observation. A SpaceMouse is used for teleoperation
 193 during data collection and online interventions. The task object is a standard USB2.0 port on a
 194 computer motherboard and a white USB flash drive. The policy outputs actions at 10 Hz, while
 195 the low-level controller executes at up to 500 Hz via RTDE with compliant force/torque regulation
 196 (Section 3.2). Figure 3 (left) provides an overview of the workcell and sensing/teleoperation setup.

197 **Reward classifier and optional stage classifier.** Following HIL-SERL [3], we first collect tele-
 198 operated rollouts to train a binary success classifier that provides sparse rewards. We collect 200
 199 successful and 600 failed episodes and train a ResNet-10 reward classifier. When enabling global
 200 state encoding (Section 3.3), we additionally annotate the collected data with four stages (grasp,
 201 approach, hole-search, insertion) and train a stage classifier with a ResNet-18 encoder and an MLP
 202 head. The input image resolution is 224×224 , and the stage classifier reaches $> 98\%$ validation
 203 accuracy within 2–3 minutes.

204 **Demonstrations and online interventions.** We collect 20 human demonstrations to initialize the
 205 Demo Buffer. During online training, a human operator provides interventions to maintain safe
 206 exploration and stabilize early learning. Interventions are scheduled in a curriculum-like manner:
 207 initially the operator may complete the full insertion; assistance is then progressively reduced (e.g.,
 208 only assisting grasp+approach, then only grasp), and finally the policy operates autonomously. We
 209 log intervention steps and episode outcomes throughout training.

210 **Evaluation protocol and metrics.** We evaluate policies at fixed training checkpoints using 20
 211 independent trials per checkpoint. We report success rate and average insertion time. For learning
 212 dynamics, we plot training curves including (i) success rate over training steps and (ii) intervention
 213 steps over time. We also report wall-clock time to reach 100% success under our evaluation protocol.

214 **4.2 Global Observer: Convergence and Final Performance**

215 To isolate the effect of the stage-aware global observer, we compare **Baseline** (HIL-SERL with
 216 standard observations) vs. **+Global** (Baseline plus global stage encoding; Section 3.3). Both variants
 217 use the same reward classifier, demonstrations, intervention schedule, and training setup.

218 **Checkpoint evaluation.** We evaluate both methods at 6k, 10k, and 12k training steps. Table 1
 219 reports success rate and average insertion time over 20 trials per checkpoint. Across all checkpoints,
 220 **+Global** achieves strictly better performance: it reaches 70% success at 6k steps, while the baseline
 221 achieves 0% (no autonomous capability under our protocol). At 10k and 12k steps, **+Global** maintains
 222 higher success and consistently shorter insertion time, indicating that phase-aware global encoding
 223 not only improves completion but also reduces unproductive contact behaviors.

Table 1: **Global observer evaluation.** Success rate and average insertion time over 20 trials.

Method	6k	10k	12k
Baseline: Success rate (%) \uparrow	0	75	90
+Global: Success rate (%) \uparrow	70	95	100
Baseline: Time (s) \downarrow	30.00	14.60	12.70
+Global: Time (s) \downarrow	15.15	6.75	5.75

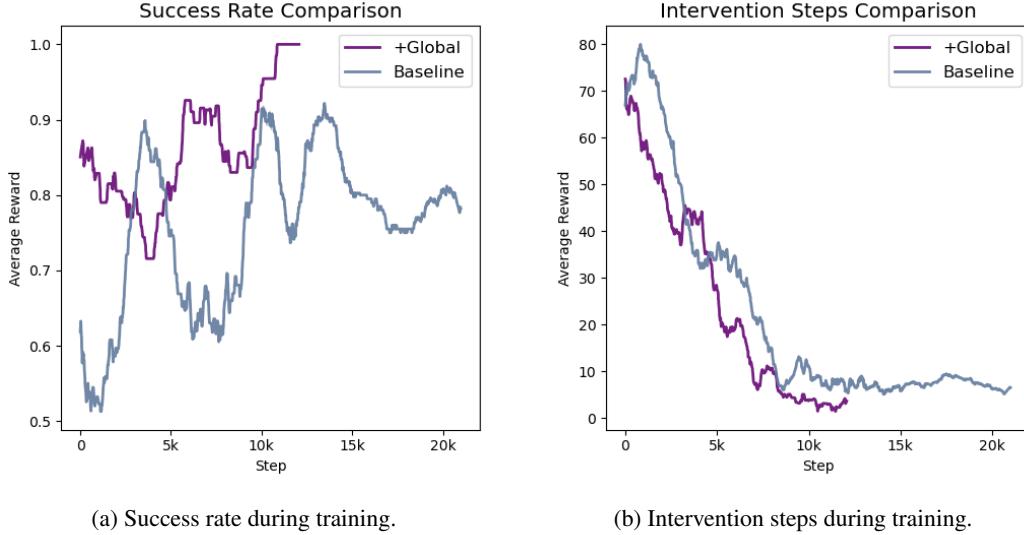


Figure 4: **Learning curves.** Baseline vs. +Global. Global encoding accelerates success and reduces required intervention.

224 **Learning dynamics and wall-clock convergence.** Figure 4 provides complementary evidence
 225 from the training process. Throughout learning, **+Global** achieves higher success while requiring
 226 fewer human intervention steps, indicating improved sample efficiency under the same sparse reward
 227 and intervention protocol. Concretely, **+Global** reaches 100% success with near-zero intervention
 228 by 12k steps (about 2 h). In contrast, the baseline requires about 20k steps (3.5 h) and still exhibits
 229 unstable success and non-trivial intervention demand. Overall, global stage encoding accelerates
 230 wall-clock convergence by more than 1.7 \times in our real-world setting, consistent with the hypothesis
 231 that reducing phase ambiguity improves credit assignment and stabilizes contact-rich learning.

232 4.3 Recovery Buffer: Error Recovery Under Typical Failure States

233 We evaluate whether the Recovery Buffer improves robustness to error states that arise in contact-rich
 234 insertion. We compare **Baseline** vs. **+Recovery** (Baseline plus Recovery Buffer; Section 3.4). During
 235 training, **+Recovery** adds human-selected recovery-relevant transitions at episode end and samples
 236 Demo/RL/Recovery buffers with a 4:4:2 ratio.

237 **Failure scenarios.** We design two representative recovery scenarios, illustrated in Figure 3 (right).
 238 **(S1) Hard-contact under the port** (top-right): the USB starts beneath the port in hard contact with
 239 the board; successful completion requires retreating to clear contact and re-approaching for insertion.
 240 **(S2) Misoriented USB pickup** (bottom-right): the USB is initially placed sideways (tilted on its side)
 241 rather than flat on the table, making direct grasp-and-insert infeasible. To succeed, the agent must
 242 first recover the object to a graspable pose—by using the gripper to push the USB into the correct
 243 orientation—and only then execute the standard pick-and-insert sequence.

244 **Protocol and metrics.** For each scenario, we run 20 trials per method starting from the specified
 245 failure state. We report recovery success rate and average recovery time.

Table 2: **Recovery evaluation.** Recovery success rate and time (20 trials).

Method	S1: hard-contact		S2: misoriented pickup	
	Success (%) ↑	Time (s) ↓	Success (%) ↑	Time (s) ↓
Baseline	55	10.16	0	30
+Recovery	90	8.63	70	10.05

246 **Results.** Table 2 summarizes recovery performance under the two designed failure initializations
247 (Figure 3, right). Adding the **Recovery Buffer** improves both *recovery success* and *efficiency*. In **S1**
248 (hard contact under the port), success rises from 55% to 90% (+35 points) while recovery time drops
249 from 10.16 s to 8.63 s, indicating more reliable retreat–re-align behavior under contact constraints.
250 In **S2** (misoriented USB pickup), the **Baseline** fails to handle the sideways initialization, whereas
251 **+Recovery** learns an explicit error-recovery routine (reorienting the USB by pushing it with the
252 gripper before grasping) and achieves 70% success with a 10.05 s average recovery time. Overall, the
253 gains across both scenarios empirically support our design rationale: prioritizing recovery-relevant
254 transitions in a dedicated buffer strengthens mistake-driven learning without changing the underlying
255 RL objective.

256 5 Conclusion

257 This paper presents a recovery-aware human-in-the-loop real-world RL system for contact-rich
258 USB2.0 pick-and-insert on UR5e. Building on HIL-SERL, we introduced three complementary
259 enhancements targeting execution, observation, and data curation: a high-rate compliant control stack
260 for safe in-contact exploration, a stage-aware global observer to reduce phase ambiguity under sparse
261 classifier rewards, and a human-gated Recovery Buffer to amplify learning from mistake-recovery
262 transitions. Experiments on a real UR5e platform show that global stage encoding improves both
263 success rate and insertion time at matched training steps and accelerates end-to-end training to 100%
264 success from 20k steps (3.5 h) to 12k steps (2 h). Moreover, the Recovery Buffer improves recovery
265 robustness in representative failure states, yielding higher recovery success and faster completion
266 time.

267 Our study highlights that, for contact-rich real-world RL, system design choices—particularly
268 compliant execution, phase-aware state modeling, and recovery-centric data organization—can be as
269 important as the underlying off-policy learner. Limitations include reliance on human supervision
270 for interventions and recovery gating, and evaluation on a single connector family and workcell
271 configuration. Future work will explore reducing human workload via learned gating, extending
272 to broader connector types and tolerances, and integrating richer representations for long-horizon
273 contact sequences. To support reproducibility, we open-source our code and provide a project website
274 with videos and additional details; the link will be released upon acceptance.

275 References

- 276 [1] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
277 policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- 278 [2] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
279 maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the*
280 *35th International Conference on Machine Learning (ICML)*, volume 80 of *Proceedings of*
281 *Machine Learning Research*, pages 1861–1870, 2018. URL <https://proceedings.mlr.press/v80/haarnoja18b.html>.
- 283 [3] Jianlan Luo et al. Precise and dexterous robotic manipulation via human-in-the-loop reinforce-
284 ment learning. *Science Robotics*, 10(105):eads5033, 2025. doi: 10.1126/scirobotics.ads5033.
285 URL <https://www.science.org/doi/abs/10.1126/scirobotics.ads5033>.
- 286 [4] Rasmus Haugaard, Simon Langaa, Christian Sloth, and Anders Glent Buch. Fast and robust
287 peg-in-hole insertion with continuous visual servoing. In *Proceedings of The 5th Conference on*

- 288 *Robot Learning (CoRL)*, volume 164 of *Proceedings of Machine Learning Research*, pages 646–
 289 656. PMLR, 2021. URL <https://proceedings.mlr.press/v164/haugaard21a.html>.
- 290 [5] Tao Chen et al. Robust peg-in-hole assembly. In *Robotics: Science and Systems (RSS)*, 2021.
 291 doi: 10.15607/RSS.2021.XVII.089.
- 292 [6] Takuma Inoue, Giovanni De Magistris, Asim Munawar, and Ryuki Tachibana. Deep reinforce-
 293 ment learning for high precision assembly tasks. In *2017 IEEE/RSJ International Conference on*
 294 *Intelligent Robots and Systems (IROS)*, pages 819–825, 2017. doi: 10.1109/IROS.2017.8202244.
- 295 [7] Jianlan Luo, Eugen Solowjow, Chengtao Wen, Juan Aparicio Ojea, and Alice M. Agogino.
 296 Deep reinforcement learning for robotic assembly of mixed deformable and rigid objects. In
 297 *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages
 298 2062–2069, 2018. doi: 10.1109/IROS.2018.8594353.
- 299 [8] Cristian Camilo Beltran-Hernandez, Damien Petit, Ixchel G. Ramirez-Alpizar, Takayuki Nishi,
 300 Shinichi Kikuchi, Takamitsu Matsubara, and Kensuke Harada. Learning force control for
 301 contact-rich manipulation tasks with rigid position-controlled robots. *IEEE Robotics and*
 302 *Automation Letters*, 5(4):5709–5716, 2020. doi: 10.1109/LRA.2020.3010739.
- 303 [9] Iñigo Elguea-Aguinaco, Antonio Serrano-Muñoz, Dimitrios Chrysostomou, Ibai Inziarte-
 304 Hidalgo, Simon Bøgh, and Nestor Arana-Arexolaleiba. A review on reinforcement learning for
 305 contact-rich robotic manipulation tasks. *Robotics and Computer-Integrated Manufacturing*, 81:
 306 102517, 2023. doi: 10.1016/j.rcim.2022.102517.
- 307 [10] Jianlan Luo, Zheyuan Hu, Charles Xu, You Liang Tan, Jacob Berg, Archit Sharma, Stefan Schaal,
 308 Chelsea Finn, Abhishek Gupta, and Sergey Levine. Serl: A software suite for sample-efficient
 309 robotic reinforcement learning. *arXiv preprint arXiv:2401.16013*, 2024.
- 310 [11] Abhishek Padalkar, Gabriel Quéré, Antonin Raffin, João Silvério, and Freek Stulp. Shared
 311 control templates for in-contact reinforcement learning. *Autonomous Robots*, 48:279–301, 2024.
 312 doi: 10.1007/s10514-023-10164-6.
- 313 [12] Todd Hester, Mel Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Andrew
 314 Sendonaris, Gabriel Dulac-Arnold, Ian Osband, John Agapiou, et al. Deep q-learning from
 315 demonstrations. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 2018.
 316 URL <https://ojs.aaai.org/index.php/AAAI/article/view/11757>.
- 317 [13] Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Over-
 318 coming exploration in reinforcement learning with demonstrations. In *2018 IEEE International*
 319 *Conference on Robotics and Automation (ICRA)*, 2018. doi: 10.1109/ICRA.2018.8463162.
 320 Also available as arXiv:1709.10089.
- 321 [14] Bingyi Kang, Zequn Jie, and Jiashi Feng. Policy optimization with demonstrations. In
 322 *Proceedings of the 35th International Conference on Machine Learning (ICML)*, volume 80
 323 of *Proceedings of Machine Learning Research*, pages 2469–2478. PMLR, 2018. URL <https://proceedings.mlr.press/v80/kang18a.html>.
- 325 [15] W. Bradley Knox and Peter Stone. Interactively shaping agents via human reinforcement: The
 326 tamer framework. In *Proceedings of the Fifth International Conference on Knowledge Capture*
 327 (*K-CAP*), 2009. doi: 10.1145/1597735.1597738.
- 328 [16] Rodrigo Pérez-Dattari, Carlos Celegin, Javier Ruiz-del Solar, and Jens Kober. Interactive
 329 learning with corrective feedback for policies based on deep neural networks. *arXiv preprint*
 330 *arXiv:1810.00466*, 2018.
- 331 [17] Brijen Thananjeyan, Ashwin Balakrishna, Suraj Nair, Michael Luo, Krishnan Srinivasan, Minho
 332 Hwang, Joseph E. Gonzalez, Julian Ibarz, Chelsea Finn, and Ken Goldberg. Recovery rl: Safe
 333 reinforcement learning with learned recovery zones. *IEEE Robotics and Automation Letters*, 6
 334 (3):4915–4922, 2021. doi: 10.1109/LRA.2021.3070252.
- 335 [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for im-
 336 age recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
 337 *Recognition (CVPR)*, pages 770–778, 2016. doi: 10.1109/CVPR.2016.90.