

TORL-VLA: Tactile Guided Online Reinforcement Learning for Contact-Rich Manipulation

422 A Appendix

423 This appendix provides implementation and experimental details omitted from the main text due
424 to space constraints. Appendix A.1 describes the latch-box manipulation benchmarks, the success
425 and failure criteria, and the evaluation protocol used for autonomous evaluation and online adapta-
426 tion. Appendix A.2 summarizes the robot system, tactile-to-wrench processing, demonstrations, and
427 online replay. Appendix A.3 details the baseline adaptation and the implementation of the wrench-
428 aware reference model and stage-specific online actor-critic refiners. Appendix A.4 presents addi-
429 tional diagnostics complementing the main results, including wrench and contact-state analysis and
430 intervention-censored critic diagnostics. Appendix A.5 reports deployment inference latency. All
431 task names, method names, and ablation names follow the notation in the main paper.

432 A.1 Task and Evaluation Protocol

433 A.1.1 Latch-Box Manipulation Tasks

434 The real-robot evaluation is conducted on a latch-box manipulation benchmark with three contact-
435 rich subtasks: **Coffee Cup**, **Latch**, and **Egg**. Each subtask contains one or more local contact-rich
436 bottleneck stages, where the stage estimator can activate the corresponding stage-specific online
437 refiner during execution. The workspace structure and subtask order are shared across methods,
438 and all trials follow the same task-specific reset protocol, with object initial poses allowed to vary
439 within predefined start regions. The three subtasks are selected to cover different contact bottlenecks.
440 Coffee Cup evaluates tight insertion into a narrow holder under partial occlusion. Latch evaluates
441 contact-dependent mechanical engagement, where the robot must grip, flip, and lock a latch. Egg
442 evaluates fragile-object grasping and placement, where both insufficient and excessive contact can
443 cause failure.

444 For full-task evaluation, the robot starts from the same latch-box initial configuration and executes
445 the task sequence in one autonomous rollout. A full-task rollout is counted as successful only if the
446 robot completes Coffee Cup, Latch, and Egg sequentially without human intervention, environment
447 reset, safety stop, object drop, or visible object damage.

448 A.1.2 Success and Failure Criteria

449 The subtasks are designed to expose contact states that may appear visually close to success while
450 differing in physical interaction, such as partial insertion, incomplete latch engagement, or exces-
451 sive grasping force. All trials are reset according to the same task-specific reset protocol before
452 evaluation.

453 The observable success and failure criteria are defined at the task-outcome level rather than by
454 internal model signals or intermediate wrench traces, ensuring consistent evaluation across methods.
455 Table 3 summarizes these criteria for all subtasks and the full-task evaluation. If any listed failure
456 event occurs, the trial is marked as a failure even if human intervention could have recovered the
457 task afterward.

458 A.1.3 Evaluation Metrics and Protocol

459 For both subtask and full-task success rates, each method is evaluated over 30 autonomous tri-
460 als. The same task objects, workspace structure, reset procedure, action representation, control fre-
461 quency, and receding-horizon execution protocol are used for all methods. During final autonomous
462 evaluation, no human intervention, environment reset, safety-stop recovery, or online parameter up-
463 date is allowed. Online parameter updates are performed only during the online data-collection

Table 3: **Observable success and failure criteria for autonomous evaluation.**

Task	Success condition	Failure condition
Coffee Cup	Pick cup; insert into holder; release; cup remains upright and stable.	Holder miss; rim hang-up; partial insertion; large tilt; drop; unsafe contact.
Latch	Grip latch; flip toward locking direction; press into locked state; latch remains secured after release.	Missed latch; slip during flipping; missed locking edge; incomplete lock; rebound; unsafe force.
Egg	Grasp egg; transfer to holder; release; egg remains stable without visible damage.	Slip; drop; holder collision; rolling out after release; excessive grasping force; visible damage.
Full Task	Complete Coffee Cup, Latch, and Egg sequentially in one autonomous rollout.	Any subtask failure; human intervention; reset; safety stop; object drop; object damage.

464 phase and are not performed while reporting the final success numbers. For reference-model ab-
 465 lations, variants are evaluated directly on the three contact-rich subtasks. For online-adaptation
 466 ablations, all variants use the same frozen wrench-aware reference model, checkpoint schedule, and
 467 matched online adaptation duration and interaction budget.

468 A.2 System, Sensing, and Data

469 A.2.1 Robot and Sensor Setup

470 The robot platform consists of a Piper robotic arm equipped with a Pika gripper from Agile Robots.
 471 The visual input consists of three camera views: a global view captured by an Intel RealSense
 472 D435i camera, a fisheye view captured by the camera integrated in the Pika gripper, and a wrist
 473 view captured by an Intel RealSense D405 camera mounted on the Pika gripper. The proprioceptive
 474 state includes the robot joint positions and the gripper state. Tactile sensing is provided by two Tacta
 475 fingertip tactile sensors from EnsuringTech, mounted on the inner contact surfaces of the gripper
 476 fingers. Each sensor has a compact size of 21 mm \times 18.5 mm \times 4.5 mm, and each tactile pad
 477 outputs a raw 6 \times 8 tactile array.

478 Policy observations are queried and logged at 20 Hz. All modalities are timestamp-aligned during
 479 data collection. For higher-rate tactile and proprioceptive streams, the latest measurement before
 480 the query timestamp is used. For visual streams, the most recent synchronized image triplet is
 481 used. All methods that use tactile-derived wrench feedback share the same tactile preprocessing,
 482 normalization statistics, and temporal window.

483 A.2.2 Tactile-to-Wrench Processing

484 The main paper defines the tactile-derived wrench observation as $\mathbf{w}_t = [\mathbf{w}_t^L, \mathbf{w}_t^R] \in \mathbb{R}^{12}$. Here
 485 we clarify how this signal is obtained from the tactile sensor interface and how temporal wrench
 486 sequences are represented in implementation. The vendor SDK exposes two levels of tactile readout.
 487 The low-level interface returns 48 tactile points arranged on a 6 \times 8 pad, with a 3D force vector at
 488 each point. The high-level SDK interface provides a calibrated 6-DOF fingertip wrench estimate
 489 $[f_x, f_y, f_z, \tau_x, \tau_y, \tau_z]$. We denote this SDK-estimated fingertip wrench as $\tilde{\mathbf{w}}_t^i \in \mathbb{R}^6$ for each finger
 490 $i \in \{L, R\}$. The left and right fingertip outputs are concatenated in a fixed order to form the 12D
 491 wrench $\tilde{\mathbf{w}}_t = [\tilde{\mathbf{w}}_t^L, \tilde{\mathbf{w}}_t^R] \in \mathbb{R}^{12}$.

492 This vendor-provided wrench abstraction decouples the policy input from the raw tactile array lay-
 493 out. It also provides a compact force–torque representation that serves as a unified downstream
 494 interface, rather than requiring the policy to process sensor-specific raw tactile layouts.

495 For the current measured wrench used in the online refinement context, we use the SDK-estimated
 496 12D wrench, i.e., $\mathbf{w}_t = \tilde{\mathbf{w}}_t$. For temporal contact reasoning, the reference model uses $J = 10$

497 wrench samples over the most recent 2 seconds, corresponding to the recent wrench sequence
 498 $\mathbf{W}_t = \{\mathbf{w}_{t-J+1}, \dots, \mathbf{w}_t\}$ in the main paper. In implementation, before feeding this sequence into
 499 the wrench-history encoder, we apply a current-wrench-centered transform to the SDK-estimated
 500 wrench samples. Specifically, the encoder input is formed as $\{\tilde{\mathbf{w}}_{t-J+1} - \tilde{\mathbf{w}}_t, \dots, \tilde{\mathbf{w}}_{t-1} - \tilde{\mathbf{w}}_t, \mathbf{0}\}$.
 501 Similarly, the future wrench target used in joint action-wrench prediction is represented by residuals
 502 $\tilde{\mathbf{w}}_{t+k} - \tilde{\mathbf{w}}_t$. This implementation keeps the temporal structure described in the main paper while
 503 representing wrench sequences by relative wrench changes around the current contact state.

504 The same tactile-to-wrench preprocessing is used for all methods that consume wrench feedback,
 505 so differences between baselines and TORL-VLA come from the policy architecture and online
 506 adaptation mechanism rather than from different tactile preprocessing.

507 A.2.3 Training Data for Reference Model and Online Adaptation

508 Table 4 summarizes the main real-robot data budget used for the wrench-aware reference model
 509 and online actor-critic refinement. Full-task demonstrations are used to train the wrench-aware
 510 VLA reference model, which predicts reference action chunks and future wrench sequences. Main
 511 online interaction data are collected around contact-rich bottleneck stages for stage-specific online
 512 adaptation. Warmup settings are reported in Table 5.

Table 4: **Real-robot data budget for reference-model training and online adaptation.** Main online interaction time denotes accumulated robot execution time during online data collection, excluding human reset, object repositioning, and environment preparation. The Egg subtask contains two contact-rich bottleneck stages: grasping and placement.

Data source	Stage / Scope	Scale
Demonstrations	Full task	539 trajectories
Main online interaction	Coffee Cup	23.5 min
	Latch	40.3 min
	Egg: grasp	29.2 min
	Egg: place	27.8 min

513 A.3 Implementation Details

514 A.3.1 Baseline Adaptation

515 We compare all methods under the same robot platform, action representation, training demonstra-
 516 tions, execution protocol, and evaluation procedure. We do not compare against numbers reported
 517 in prior work, since robot embodiments, tactile sensors, task definitions, and data distributions differ
 518 across papers. Instead, we adapt the baselines to the same $\pi_{0.5}$ backbone and, where applicable, the
 519 same dual-fingertip tactile-derived wrench interface used by TORL-VLA.

520 The $\pi_{0.5}$ baseline directly executes the pretrained VLA policy without tactile-derived wrench input.
 521 TA-VLA and ForceVLA are reimplemented on the same $\pi_{0.5}$ backbone and adapted to the dual-
 522 fingertip wrench observation. This keeps the comparison focused on whether tactile- or force-aware
 523 offline VLA designs can improve contact-rich manipulation under the same hardware and data set-
 524 ting. TORL-VLA w/o RL directly executes the learned wrench-aware reference model, without
 525 online RL refinement. RLT uses reference-guided online refinement, but it does not incorporate
 526 wrench input during online refinement and does not include the intervention-censored critic. TORL-
 527 VLA uses the wrench-aware reference model together with wrench-conditioned online refinement
 528 and intervention-censored value learning.

529 All online-adaptation variants use the same frozen wrench-aware reference model and the same
 530 online adaptation duration, while each variant collects its own online robot data. This ensures that
 531 the comparison between RLT, TORL-VLA, and online ablations focuses on the effect of wrench-
 532 context conditioning and intervention-censored value learning under a matched online interaction
 533 budget.

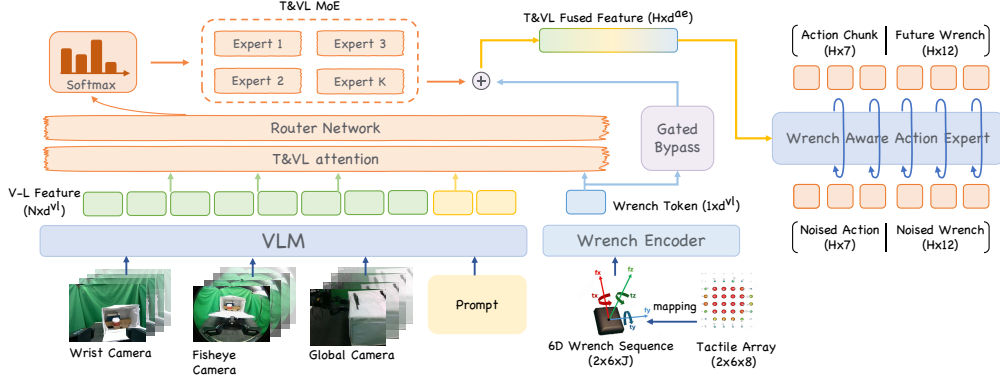


Figure 7: Implementation details of the wrench-aware VLA reference model.

534 A.3.2 Wrench-Aware Reference Model

535 This subsection specifies the implementation details of Stage I, which are omitted from the main
 536 text for space. At each timestep t , the wrench-aware VLA reference model takes as input three
 537 camera views (wrist, fisheye, and global), a language instruction, robot proprioception, and tactile
 538 readings from two piezoelectric pads mounted on the inner surfaces of the gripper fingers. As
 539 described in Appendix A.2.2, the tactile readings are converted into the current measured wrench
 540 \mathbf{w}_t , and the recent wrench window is represented in a current-wrench-centered form before being
 541 fed into the wrench-history encoder. The resulting history sequence is compressed into a fixed-size
 542 wrench token $\mathbf{z}_t^{\text{tac}} \in \mathbb{R}^{1 \times d_{v1}}$ by a lightweight MLP encoder. In parallel, the three camera images,
 543 language instruction, and proprioception are encoded by the VLM backbone into visual-language
 544 prefix hidden states $\mathbf{P}_t \in \mathbb{R}^{N_p \times d_{v1}}$.

545 Rather than injecting the wrench token early into the VLM, we adopt a late fusion strategy. The
 546 wrench token is appended to \mathbf{P}_t after visual-language encoding, and a self-attention layer is ap-
 547 plied to the concatenated sequence $[\mathbf{P}_t; \mathbf{z}_t^{\text{tac}}]$ to produce the fused multimodal representation $\mathbf{H}_t \in$
 548 $\mathbb{R}^{(N_p+1) \times d_{v1}}$. This preserves the pretrained visual-semantic features while enabling cross-modal
 549 interaction with tactile-derived wrench feedback. The fused representation is then passed through
 550 a soft-routing Mixture-of-Experts module, where a router network computes per-expert weights
 551 $\boldsymbol{\rho}_t = \text{Softmax}(R_{\text{tac}}(\mathbf{H}_t))$, and the routed feature is obtained as $\mathbf{F}_t = \sum_{m=1}^M \boldsymbol{\rho}_{t,m} \odot E_m(\mathbf{H}_t)$. This
 552 allows the model to dynamically route features according to the current contact phase.

553 To retain fine-grained physical information that may be diluted by the MoE transformations, we
 554 add a zero-initialized wrench-token bypass. The routed feature and the original wrench token are
 555 projected to the action horizon and combined as $\mathbf{G}_t = \Pi_{\text{tac}}(\mathbf{F}_t) + \alpha P_{\text{tac}}^H(\mathbf{z}_t^{\text{tac}}) \in \mathbb{R}^{H \times d_{\text{ae}}}$, where
 556 Π_{tac} maps the routed feature to the action horizon, P_{tac}^H projects and expands the wrench token to
 557 horizon-aligned guidance tokens, and α is a learnable scalar initialized to zero. Thus, the bypass
 558 contributes nothing at the start of training and grows adaptively as the model learns to use wrench
 559 feedback.

560 Within the flow-matching action expert, the noisy joint generation variable $\mathbf{x}_{t,s}^{\text{fm}} \in \mathbb{R}^{H \times (d_a+12)}$
 561 is processed at interpolation time $s \in [0, 1]$ to yield action-side hidden states $\mathbf{S}_{t,s}$. The wrench-
 562 conditioned guidance signal is added as a residual, producing $\tilde{\mathbf{S}}_{t,s} = \mathbf{S}_{t,s} + \mathbf{G}_t$. A joint output
 563 head then predicts a single velocity field over the concatenated action-wrench target, which is split
 564 into action and wrench dimensions as $\hat{\mathbf{u}}_{t,s} = [\hat{\mathbf{u}}_{t,s}^a, \hat{\mathbf{u}}_{t,s}^w]$. Solving the learned flow from $s = 0$ to
 565 $s = 1$ yields the joint prediction: a reference action chunk $\hat{\mathbf{A}}_{t:t+H-1} \in \mathbb{R}^{H \times d_a}$ that serves as the
 566 motion prior for the online actor, and a future wrench sequence $\hat{\mathbf{W}}_{t:t+H-1} \in \mathbb{R}^{H \times 12}$ that provides
 567 the expected contact cue for wrench-conditioned action refinement in Stage II. Following the pre-
 568 processing described in Appendix A.2.2, the future wrench target used in training is represented by
 569 current-centered residuals, i.e., $\tilde{\mathbf{w}}_{t+k} - \tilde{\mathbf{w}}_t$.

570 A.3.3 Online Reference Adaptation

571 This subsection adds implementation details that are not fully specified in the main text. During
 572 Stage II adaptation, the wrench-aware VLA is kept frozen and only the lightweight stage-specific
 573 actor-critic refiners are updated.

574 **Reference interface.** The online refiner uses the same executable reference interface and context
 575 definition as in the main text. In implementation, the frozen wrench-aware VLA predicts a 50-step
 576 action sequence and a 50-step future wrench sequence, and the first $K = 10$ steps are used for online
 577 refinement and execution. At 20 Hz, this corresponds to a 0.5 s replanning interval.

578 To obtain the compact VLA token z_t^{vla} , we append a learned special token to the final hidden-state
 579 sequence of the frozen wrench-aware VLA and pass the augmented sequence through a lightweight
 580 token encoder. The output at the special-token position is used as z_t^{vla} . The token encoder is trained
 581 before online adaptation and is frozen during online actor-critic updates.

582 **Stage estimator routing.** As shown in Figure 8, the stage estimator is a lightweight chunk-level pre-
 583 diction module rather than an action policy. Its purpose is to decide how each VLA reference chunk
 584 should be executed during Stage II adaptation. At each VLA query, the frozen wrench-aware VLA
 585 first produces the reference action chunk and the policy-side features used by the online module.
 586 The stage estimator reuses these features, including the compact VLA token context, proprioception
 587 context, and wrench context, to predict a discrete stage label and its confidence. The predicted stage
 588 label is then mapped to a runtime route before executing the next K -step action chunk.

589 The route mapping separates full-task execution into
 590 base execution and local contact-window refine-
 591 ment. Labels outside predefined contact windows are
 592 mapped to the base route, where the frozen VLA re-
 593 ference chunk is executed directly. Only labels corre-
 594 sponding to contact-critical windows activate the as-
 595 sociated stage-specific actor-critic refiner. Therefore,
 596 the stage estimator does not assign an entire subtask to
 597 an online refiner. Instead, it only selects short win-
 598 dows in which contact-aware refinement is needed.
 599 For example, approach and inter-stage transport are
 600 usually handled by the frozen VLA, while online re-
 601 finement is activated around cup insertion, latch en-
 602 gagement/locking, egg grasping, and egg placement.

603 To avoid unstable switching caused by isolated predic-
 604 tion noise, we apply a simple route gate before execu-
 605 tion. A predicted route is accepted only when the stage
 606 confidence exceeds 0.9 and the same mapped route is
 607 predicted for two consecutive chunks. Otherwise, the
 608 system keeps the previous stable route or falls back to
 609 the base route. Because the stage estimator reuses fea-
 610 tures already produced at the VLA query, this routing
 611 step does not require an additional VLA forward pass.

612 Figure 9 compares the estimated stages with ground-truth annotations over a full-task rollout. The
 613 stage estimator achieves an absolute semantic accuracy of 95.5%. The remaining mismatches are
 614 concentrated near the continuous boundaries between adjacent stages. This is expected because the
 615 physical transition between two neighboring phases is gradual rather than instantaneous, and manual
 616 stage annotations at such boundaries can vary slightly. Direct inspection of the aligned camera
 617 observations shows that the physical manipulation states corresponding to the ground truth (GT)
 618 and estimation results (ER) are visually and functionally similar across stable intervals. Thus, these
 619 boundary-level discrepancies mainly reflect annotation ambiguity around phase transitions rather

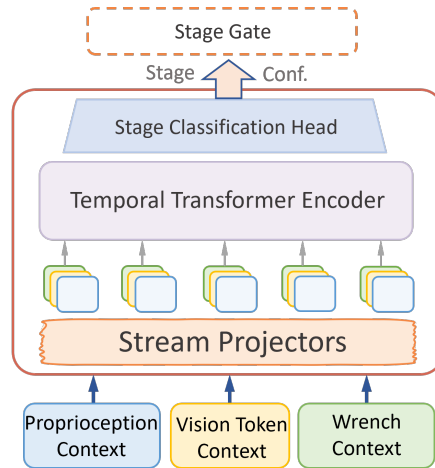


Figure 8: **Stage estimator architecture.** Policy-side features are projected, temporally aggregated, and used to predict a stage label and confidence for route filtering.

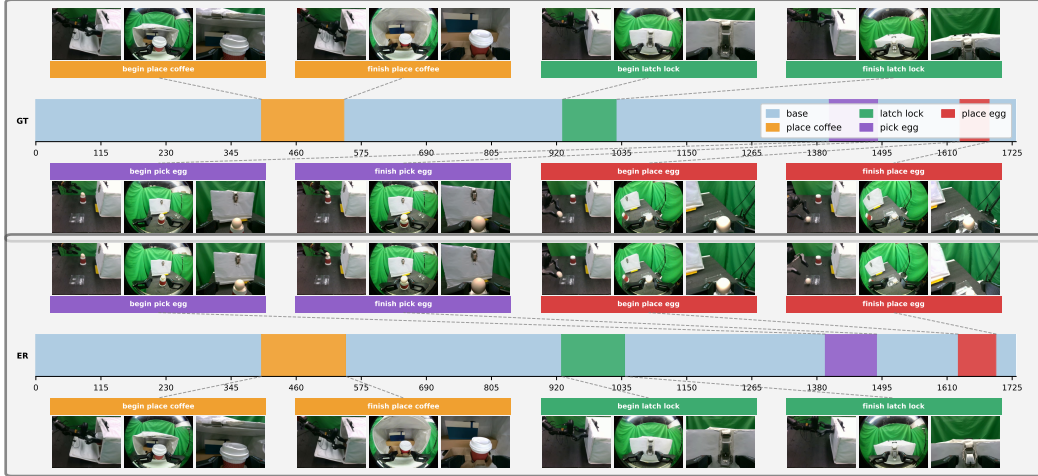


Figure 9: **Detailed alignment of predicted stages against ground truth over full-task execution.** The visualization trace shows chunk-level stage estimates aligned with sequential camera observations. Minor discrepancies mainly appear at continuous phase boundaries, while the physical manipulation states in GT and ER remain consistent across stable intervals.

620 than incorrect recognition of stable task stages. In deployment, the confidence and temporal-stability
 621 gate further reduces the chance that such boundary fluctuations lead to unstable actor switching.
 622 Moreover, the stage-specific online actors operate on continuous wrench-conditioned contexts, making
 623 the refinement robust to small timing shifts around the beginning or end of a contact window.

624 When a chunk is routed to a contact window, the corresponding stage-specific actor-critic refiner is
 625 used. The actor embeds the input groups in c_t , concatenates the features, and directly outputs the
 626 refined executable chunk:

$$A_t^\phi = \pi_{\phi_{gt}}(c_t). \quad (5)$$

627 The actor does not predict an explicit residual added to A_t^{ref} . Action chunks are normalized inter-
 628 nally for training and converted back to executable robot commands before deployment.

629 **Replay buffer and control-source labels.** Before main online adaptation, each stage-specific re-
 630 finer is initialized with warmup data, where execution mainly follows the VLA reference and human
 631 intervention is applied near failure or unsafe contact. After warmup, the corresponding actor snap-
 632 shot is used in closed-loop execution.

633 Online data are stored at the executed K -step chunk level. Each replay entry contains A_t^{ref} , the
 634 executed chunk A_t^{data} , robot proprioception, the current measured wrench, recent measured wrench
 635 history, the aligned predicted wrench cue \hat{W}_t^{ref} , rewards, terminal flag, next context, and per-step
 636 control-source labels. For each low-level step j , $\sigma_{t,j} \in \{\text{policy}, \text{human}\}$ records whether the action
 637 was generated by the learned policy or by human intervention. These labels are not actor observa-
 638 tions; they are used only for behavior regularization, intervention-boundary target construction, and
 639 diagnostics.

640 A.3.4 Training Details and Hyperparameters

641 Table 5 summarizes the key hyperparameters used for reference generation and online adaptation.
 642 We report the settings that define the reference interface, execution schedule, online update cadence,
 643 intervention-boundary construction, and stage-specific actor-critic refiners. Since the Egg subtask
 644 contains two contact-rich bottleneck stages, we use separate actor-critic refiners for egg grasping
 645 and egg placement.

Table 5: Key hyperparameters for reference generation and online adaptation.

		Online learning / AC	Value
Reference / execution	Value	Discount factor γ	0.99
Reference horizon H	50	Update-to-data ratio	5
Executed horizon K	10	Critic update frequency	every step
Action / wrench dimension	7 / 12	Actor update frequency	every 2 steps
Robot command representation	delta chunk	Target-network EMA τ	0.005
Control frequency	20 Hz	Actor snapshot interval	100 steps
Replanning interval	0.5 s	Checkpoint interval	1000 steps
Wrench history length	10	Intervention cost c_{int}	1.0
Wrench history window	2.0 s	Human-control threshold ρ_{int}	0.5
Future wrench horizon	50	Warmup replay chunks	300
Future wrench loss λ_w	0.3	Warmup gradient updates	5k
VLA token dimension z_t^{vla}	2048	Hidden dimension	256
		Coffee Cup AC layers	3 / 3
		Egg-grasp AC layers	3 / 3
		Egg-place AC layers	3 / 3
		Latch AC layers	4 / 4

646 Online learning uses replay-buffer updates. For each newly added replay transition, the learner
 647 performs five gradient updates on average; the critic is updated at every learner step, while the actor
 648 is updated every two learner steps. Checkpoints are saved every 1000 learner steps.

649 These settings specify the TORL-VLA implementation used in our online adaptation experiments,
 650 including the execution horizon, control frequency, warmup schedule, online update cadence, and
 651 stage-specific actor-critic architecture.

652 A.4 Additional Diagnostics

653 A.4.1 Wrench and Contact-State Analysis

654 This section analyzes tactile-derived wrench signals during contact-rich execution. The goal is not to
 655 introduce an additional benchmark metric, but to show how measured wrench and predicted future
 656 wrench behave in the task phases visualized below.

657 **Measured contact evolution.** Figures 12 and 10 show representative real-execution segments from
 658 the latch-locking and cup-insertion tasks. The wrench curves shown in this analysis are measured
 659 from the right-fingertip tactile sensor. The latch sequence in Figures 12 covers the process from the
 660 initial latch grasp to latch flipping. During this process, the commanded action changes smoothly,
 661 while the measured fingertip wrench exhibits clear temporal variation. In the cup task shown in
 662 Figures 10, the gripper must maintain a normal force f_z to stably hold the cup, while also applying
 663 an appropriate downward force f_y during insertion to place the cup into the holder. These examples
 664 show that the contact state in our tasks is not a simple binary contact/no-contact variable, but evolves
 665 continuously with the interaction among the robot, the manipulated object, and the environment.
 666 Beyond characterizing normal contact evolution, the wrench signal also provides cues for identifying
 667 deviations from the expected execution. For example, grasping the latch too tightly may introduce
 668 excessive contact forces that prevent the latch from being flipped. Similarly, during cup insertion,
 669 a slight positional offset may be difficult to detect from visual observations alone, but it can be
 670 revealed by abnormal variations in f_y .

671 **Future-wrench prediction.** Figure 13 evaluates whether the reference model can predict future
 672 contact trends. Each panel pairs the current visual observation and the observation after 50 action
 673 steps with the predicted and measured fingertip wrench norms. The examples cover latch lock-
 674 ing, cup insertion, and egg grasping, thereby testing the prediction head under mechanism con-
 675 tact, environmental contact, and delicate object contact. These curves suggest that the reference

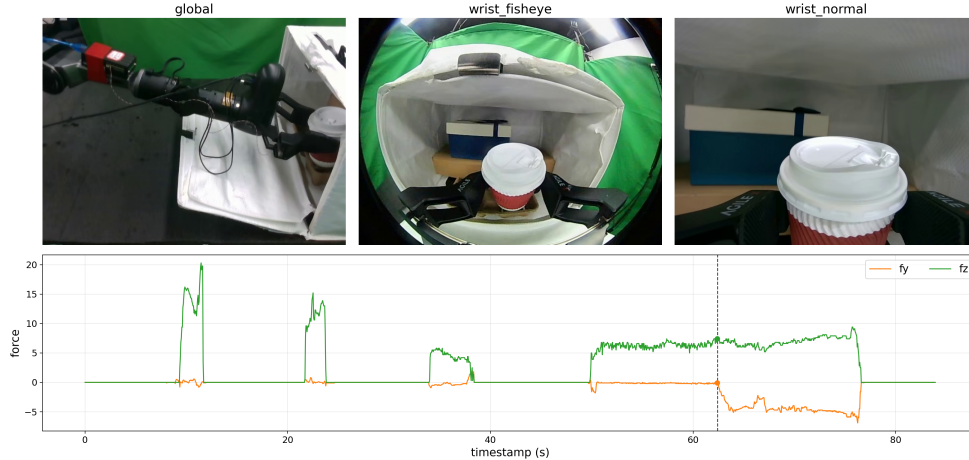


Figure 10: Measured contact evolution in the cup-insertion task. The wrench curves are measured from the right-fingertip tactile sensor. The gripper maintains the normal force f_z to stably hold the cup while applying an appropriate downward force f_y during insertion.

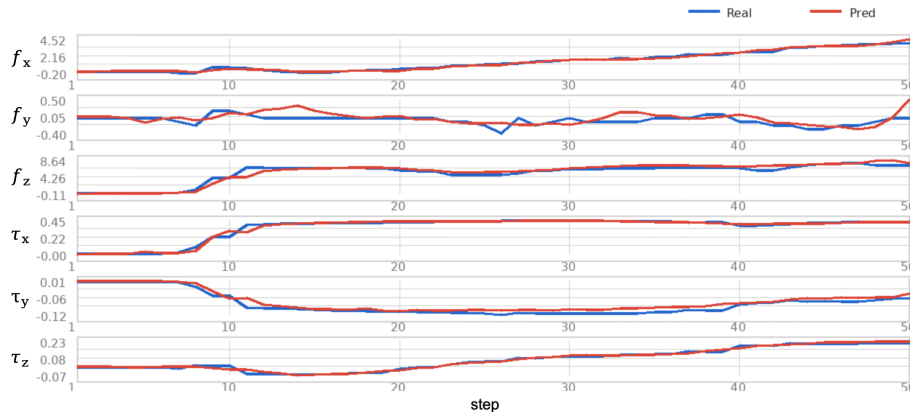


Figure 11: Detailed six-dimensional future-wrench prediction on a latch-locking segment. The figure compares the predicted and measured six-dimensional wrench trajectories for the left-fingertip tactile sensor, including the three force components and three torque components. This view complements the norm-based summaries in Figure 13 by showing whether the model captures the component-wise contact evolution.

676 model anticipates the contact consequences of its proposed future action. Figure 11 further expands
 677 one latch-locking segment into the left-fingertip components, showing the predicted and measured
 678 force–torque trajectories along all six dimensions.

679 **Beyond history extrapolation.** A potential concern is that the future-wrench head may simply ex-
 680 trapolate the trend of the input wrench history. To probe this, we examine the examples in subpanels
 681 (a1), (b1), and (c1) of Figure 13. In these cases, the input wrench histories are equal to zero before
 682 substantial contact is established. A history-copying predictor would therefore remain near zero or
 683 continue the previous flat trend. Instead, the model predicts a future rise or change in the fingertip
 684 wrench while generating the corresponding reference action. This suggests that the reference model
 685 infers upcoming contact dynamics from the visual state and the shared action–wrench generation
 686 context, rather than merely extrapolating short-term wrench history.

687 A.4.2 Intervention-Censored Critic Diagnostics

688 This section provides process-level diagnostics for the intervention-censored critic on the Latch-
 689 Lock stage, i.e., the locking contact window within the Latch subtask. We use this stage because

690 it contains sufficient policy-to-human boundary contexts for stable analysis. These diagnostics are
 691 not additional final evaluation metrics; they are used only to inspect whether Q_{ic} captures the in-
 692 tended intervention-risk signal and whether the learned actor moves away from actions with lower
 693 intervention-censored value. Unless otherwise specified, all diagnostics are computed from the fi-
 694 nal Latch-Lock online replay of the full model, and critic values are evaluated using the final critic
 695 checkpoint.

696 **Boundary contexts and diagnostic scores.** For diagnostics, we define a clean policy-generated
 697 chunk as a mostly policy-generated chunk that is not followed by an intervention-dominated chunk.
 698 A policy-to-human boundary chunk is a mostly policy-generated chunk immediately preceding an
 699 intervention-dominated chunk, using the same threshold as the intervention-boundary mask in the
 700 main method. A boundary context denotes the actor context c_t associated with such a boundary
 701 chunk. These diagnostic labels are used only for offline analysis and are not provided to the actor as
 702 observations.

703 All critic values are reported as the minimum over the double-Q pair. For AUC-based boundary-risk
 704 ranking, we use $-Q(c, A)$ as the score: lower critic values indicate lower predicted value, and for
 705 Q_{ic} this is intended to correspond to higher intervention risk. Thus, a larger AUC means that a critic
 706 more reliably ranks boundary chunks above clean policy-generated chunks under this diagnostic
 707 score.

708 **Boundary-context distribution and actor behavior.** Table 6 reports two diagnostics on the final
 709 Latch-Lock replay. The left table shows how policy-to-human boundary contexts are distributed
 710 across online episodes, indicating that they are recurring contact-risk events rather than a few iso-
 711 lated outliers. The right table compares the logged pre-intervention action, the VLA reference
 712 action, and the final actor action on the same 165 boundary contexts. The final actor achieves
 713 higher average Q_{ic} value than both the logged action and the VLA reference, suggesting that the
 714 learned actor moves away from lower-IC-value actions on these boundary contexts. We call a
 715 boundary context *IC-hinge-active* when the reference-relative IC hinge in the actor loss is active,
 716 i.e., $\min_i Q_{ic}^{(i)}(c, A^\phi) < \min_i Q_{ic}^{(i)}(c, A^{\text{ref}})$.

Table 6: **Boundary-context distribution and actor behavior.** Left: episode-level distribution of
 policy-to-human boundary contexts in the final Latch-Lock replay. Right: action comparison on the
 same 165 boundary contexts using the final IC critic.

Metric	Value	Actor diagnostic	Value
Episodes with replay	162	Logged $Q_{ic}(c, A^{\text{data}})$	-1.003
Boundary contexts	165	Reference $Q_{ic}(c, A^{\text{ref}})$	-0.573
Episodes ≥ 1 boundary	125 / 162 = 77.2%	Final actor $Q_{ic}(c, A^\phi)$	-0.271
Episodes ≥ 2 boundaries	32 / 162 = 19.8%	ΔQ_{ic} : actor vs. reference	+0.301
Episodes ≥ 3 boundaries	8 / 162 = 4.9%	ΔQ_{ic} : actor vs. logged data	+0.731
Mean per episode	1.02	IC-hinge-active boundary contexts	21 / 165 = 12.73%
Median per episode	1		
90th percentile	2		
Maximum per episode	3		

717 **Boundary-risk separation.** Table 7 compares Q_{ic} and Q_{task} on clean policy-generated chunks and
 718 policy-to-human boundary chunks. Clean Mean and Boundary Mean denote the average critic values
 719 on the two groups, and Clean-Boundary is their difference. AUC is computed using $-Q(c, A)$ as the
 720 boundary-risk score. The IC critic separates boundary chunks from clean chunks much more clearly
 721 than the task critic, supporting the separation between task-return estimation and intervention-risk
 722 modeling.

723 **Actor checkpoint comparison.** Table 8 compares actor checkpoints saved at learner global steps
 724 10k, 20k, and 33k. To avoid comparing values from drifting critic checkpoints, we use the final
 725 Q_{ic} as a fixed diagnostic scorer for all actor checkpoints. The comparison is retrospective: all
 726 checkpoints are evaluated on the same final boundary-context set, although not all of these contexts

Table 7: **Boundary-risk separation by task and IC critics.** AUC is computed using $-Q(c, A)$ as the boundary-risk score; for Q_{ic} , lower values are intended to indicate higher intervention risk.

Critic / Action	Clean Mean	Boundary Mean	Clean-Boundary	AUC \uparrow
$Q_{ic}(c, A^{\text{data}})$	0.073	-1.003	1.076	0.999
$Q_{ic}(c, A^{\text{ref}})$	-0.007	-0.573	0.565	0.913
$Q_{ic}(c, A^{\phi})$	0.091	-0.271	0.363	0.806
$Q_{\text{task}}(c, A^{\text{data}})$	0.333	0.247	0.086	0.621
$Q_{\text{task}}(c, A^{\text{ref}})$	0.316	0.241	0.074	0.626
$Q_{\text{task}}(c, A^{\phi})$	0.351	0.272	0.079	0.618

727 were present in replay at earlier checkpoints. A boundary context is counted as IC-hinge-active at a
 728 checkpoint if the actor action at that checkpoint has lower IC value than the VLA reference under the
 729 same context. A 10k IC-hinge-active context is considered resolved if it is no longer IC-hinge-active
 730 under the final actor. The number of IC-hinge-active boundary contexts decreases from 89/165 at
 731 10k to 21/165 at the final 33k checkpoint, indicating that the actor progressively reduces boundary
 732 contexts where its action has lower IC value than the VLA reference.

Table 8: **Actor checkpoint comparison on boundary contexts.** The final IC critic is used as a fixed scorer on the same final boundary-context set. Checkpoints are indexed by learner global step. Margin statistics are computed only on the 89 contexts that are IC-hinge-active at the 10k checkpoint.

Retrospective actor diagnostic	Value
IC-hinge-active contexts at 10k global-step checkpoint	89 / 165 = 53.94%
IC-hinge-active contexts at 20k global-step checkpoint	66 / 165 = 40.00%
IC-hinge-active contexts at final 33k global-step checkpoint	21 / 165 = 12.73%
10k IC-hinge-active contexts resolved by final actor	71 / 89 = 79.8%
Mean margin on 10k IC-hinge-active contexts	-0.207
Mean margin on the same contexts at final 33k checkpoint	+0.213
Mean margin improvement on the same contexts	+0.421

733 **Online intervention behavior.** Table 9 reports rollout-level process diagnostics during late online
 734 training. These statistics are not used as final autonomous evaluation metrics. Human-Controlled
 735 Step Ratio is computed at the low-level control-step level as the fraction of steps involving human
 736 control, and Avg. Human Interventions counts human takeover events per episode. Compared with
 737 the variant without the IC critic under the same late-training window, the full model shows lower
 738 human-intervention reliance and more intervention-free episodes during late online training. This
 739 is consistent with the intended role of Q_{ic} : reducing policy behaviors that tend to lead into human
 740 correction.

Table 9: **Latch-Lock online intervention diagnostics.** Statistics are computed over late online training rollouts and are not used as final autonomous evaluation metrics. Human-Controlled Step Ratio is computed at the low-level control-step level, and Avg. Human Interventions counts human takeover events per episode.

Variant / Window	Human-Controlled Step Ratio \downarrow	Avg. Human Interventions \downarrow	Intervention-Free Episodes \uparrow
w/o IC Critic, last 100 ep.	27.78%	13.40	11.0%
Full model, last 100 ep.	25.32%	10.68	29.0%
w/o IC Critic, last 50 ep.	25.42%	13.10	12.0%
Full model, last 50 ep.	18.12%	7.48	40.0%

741 **A.5 Deployment Efficiency**

742 **A.5.1 Deployment Inference Latency**

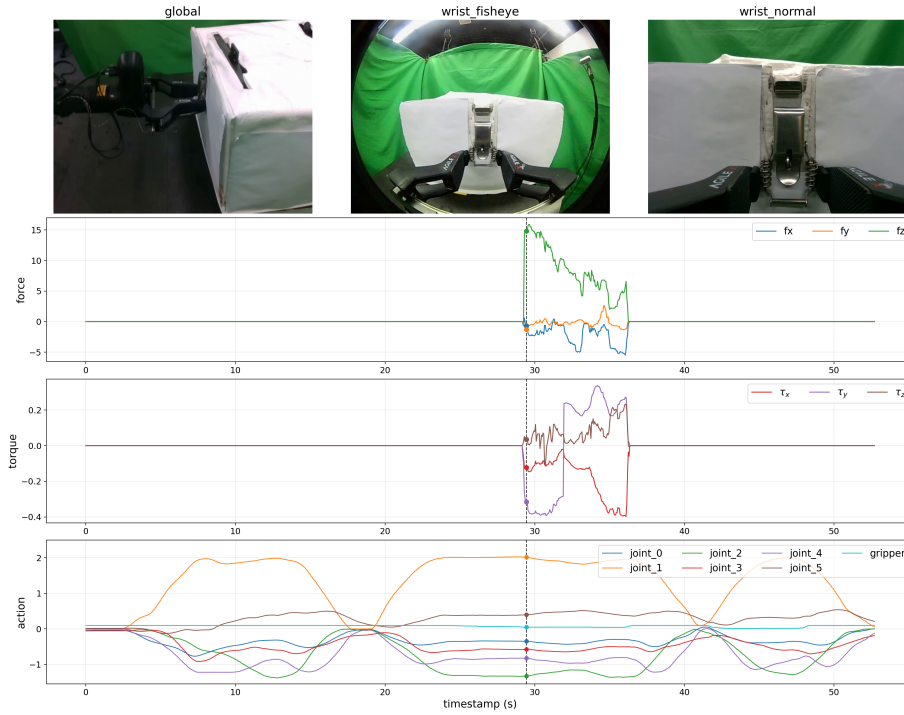
743 Table 10 reports deployment-time inference latency measured on a single NVIDIA RTX 4090
744 GPU with batch size one. We measure online policy computation only, excluding robot com-
745 munication, observation acquisition, low-level controller execution, and environment preparation.
746 All policy-query latencies are averaged over 30 timed warm-call queries after excluding the first
747 JIT/initialization call and using five warmup queries.

748 During a TORL-VLA policy query, the frozen wrench-aware VLA generates the reference action
749 sequence and predicted future wrench sequence. The stage estimator performs chunk-level routing
750 using policy-query features and does not require an additional VLA forward pass. When the cur-
751 rent chunk is routed to a contact-critical stage, the corresponding stage-specific online actor refines
752 the reference action chunk. The reported TORL-VLA total is measured end-to-end for this actor-
753 refinement route, rather than obtained by summing separately measured module times. The base
754 $\pi_{0.5}$ row is included only as a reference-policy comparison and is not part of the TORL-VLA total.

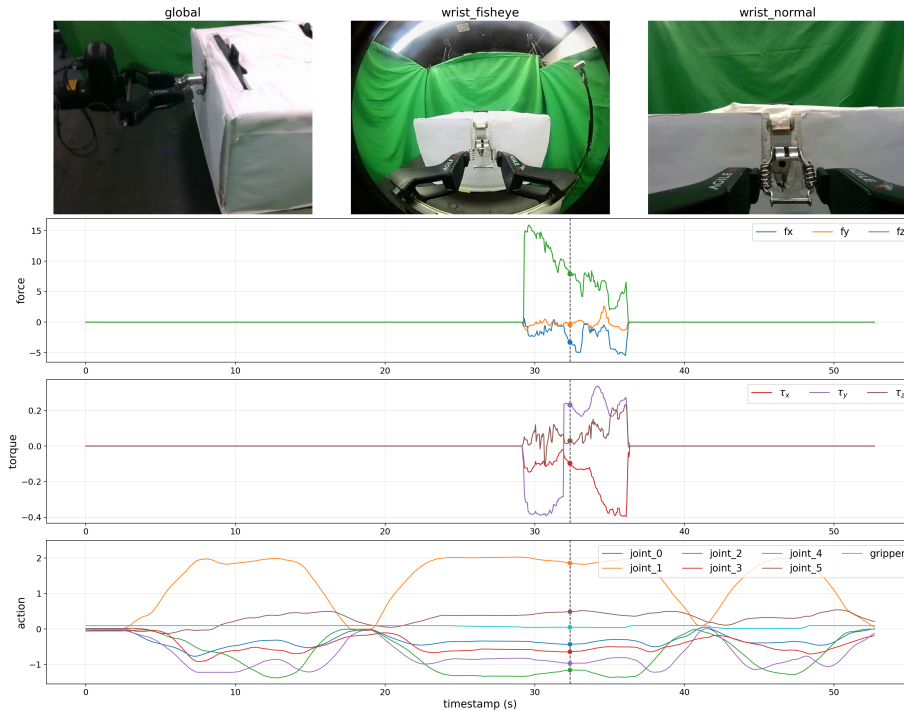
Table 10: **Deployment inference latency.** Mean warm-call latency is measured on a single NVIDIA RTX 4090 GPU with batch size one.

Module / Query	Timed Calls	Mean Latency
Base $\pi_{0.5}$ policy query (comparison)	30	94.7 ms
TORL-VLA wrench-aware reference query	30	100.2 ms
Stage estimator routing	30	1.7 ms
Stage-specific online actor	30	1.1 ms
End-to-end TORL-VLA query with actor refinement	30	103.0 ms

755 Compared with the base $\pi_{0.5}$ policy query, the full TORL-VLA actor-refinement route increases
756 mean policy-query latency by 8.3 ms, corresponding to an 8.8% relative overhead. The stage esti-
757 mator and online actor add only a few milliseconds of latency, indicating that TORL-VLA introduces
758 lightweight deployment overhead while retaining the same real-time policy-query regime.

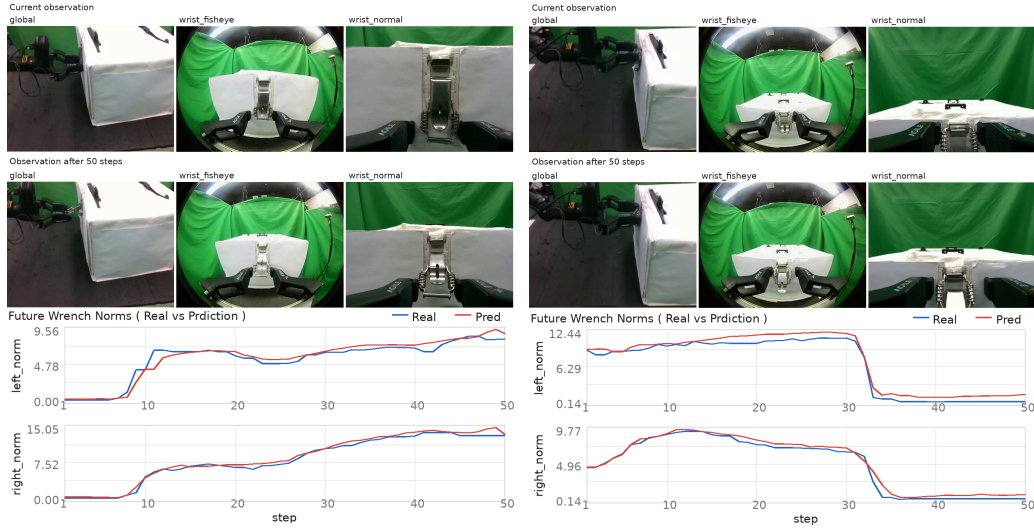


(a) Initial latch-grasping phase

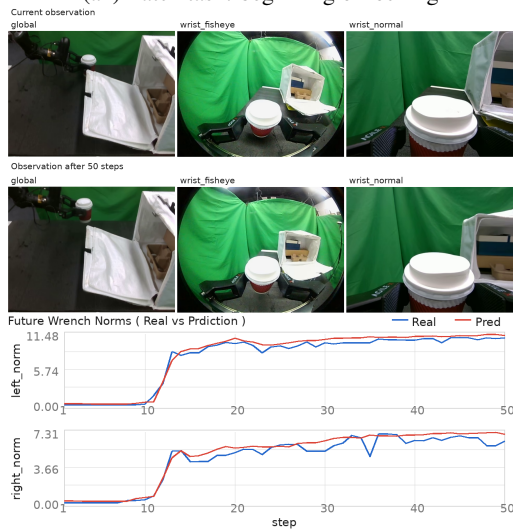
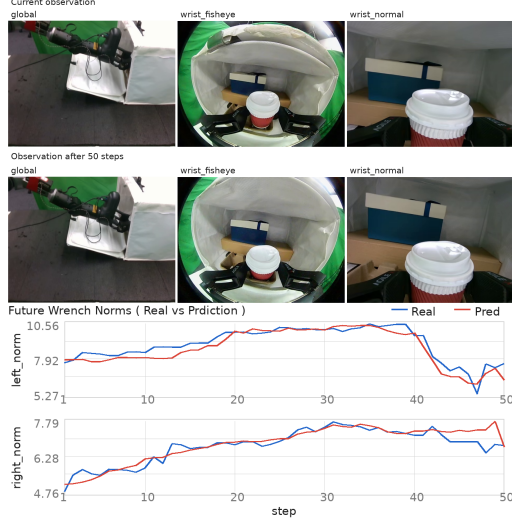


(b) Latch-flipping phase

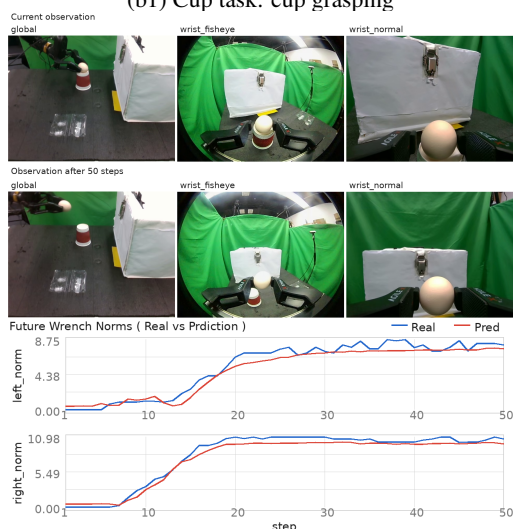
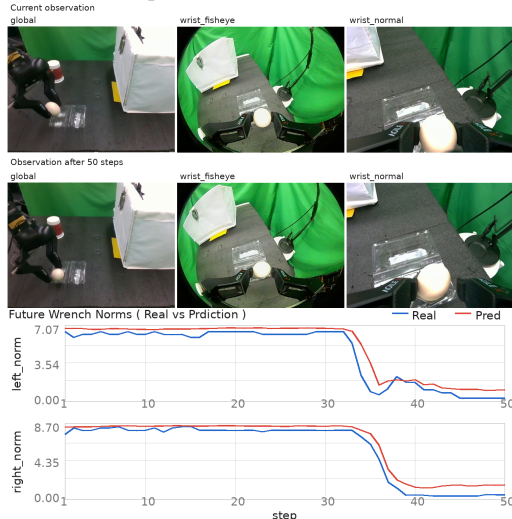
Figure 12: Measured contact evolution in the latch-locking task. Each panel shows synchronized camera views, the measured wrench $(f_x, f_y, f_z, \tau_x, \tau_y, \tau_z)$ from the right-fingertip tactile sensor, and action traces. The sequence covers the process from initially grasping the latch to flipping it.



(a2) Latch task: end of locking



(b2) Cup task: insertion into the holder



(c2) Egg task: end of grasping

Figure 13: Future-wrench prediction examples across the latch, cup, and egg tasks. Each panel shows the current and post-50-step observations above, and compares the predicted future wrench norms with the measured wrench norms from the left and right fingertip tactile sensors below.