

ALOHA Lightning: Learning Fast and Precise Manipulation

John Hua Yao Qi Wu Yihuai Gao Chelsea Finn Zipeng Fu
Stanford University

<https://aloha-lightning.github.io>

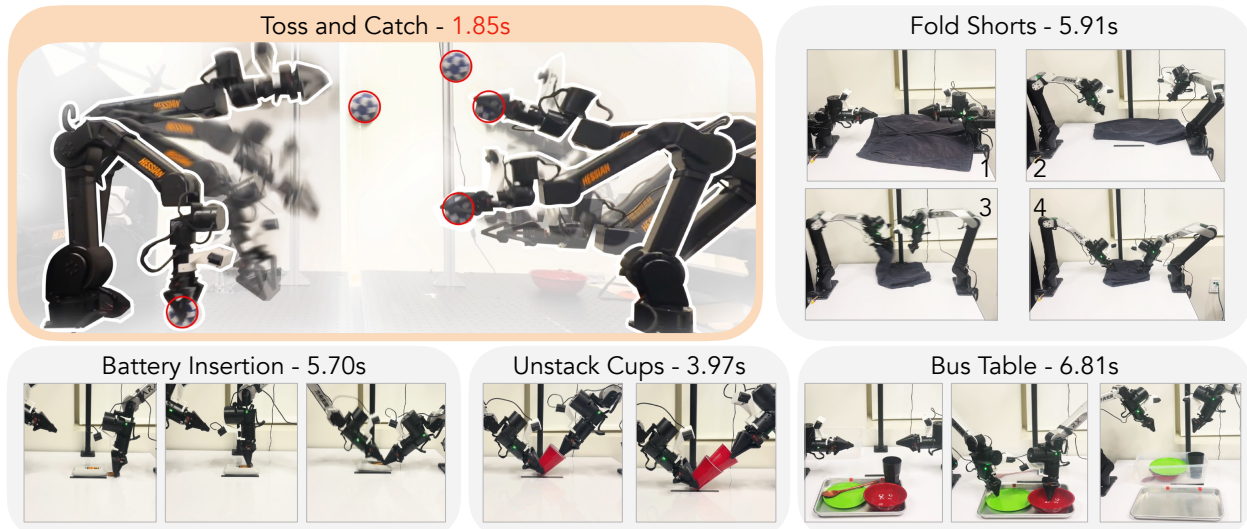


Fig. 1: *ALOHA Lightning*. We present a full-stack system for learning precise and high-speed manipulation policies, which can autonomously do bimanual toss and catch in under 2 seconds, bus a full dining set under 7, and fold a pair of shorts under 6.

Abstract—Learning from human demonstrations has enabled robots to acquire a wide range of manipulation skills, but learned policies typically execute far slower than ordinary humans. This speed gap is mainly due to lack of an interface for collecting demonstration data at high speed, and the difficulty in training policies that can robustly execute high-speed motions. In this paper, we present *ALOHA Lightning*, a system for learning fast and precise robotic manipulation. Our system uses kinesthetic teaching to intuitively collect near-human-speed demonstrations on a backdrivable bimanual platform, yielding natural and fast trajectories. We also present a learning pipeline that enables smooth high-speed execution through test-time action smoothing and aligns the visual data distribution between data collection and deployment with masking. Given 50 demonstrations for each task, *ALOHA Lightning* autonomously completes bimanual tasks such as folding shorts, and bussing tables for over 80% success rates, and ball tossing and catching with 50% success rate while close to human speed.

I. INTRODUCTION

Learning from human demonstrations has proven to be a powerful paradigm for teaching robots a wide array of complex skills. We have seen robots learn everything from fine-grained manipulation tasks such as threading a Velcro tape to mobile manipulation tasks such as sauteing and serving shrimp [1], [2]. However, despite these successes, a significant gap remains between robotic and

human performance: *speed*. While a person can perform many everyday chores with fluid and rapid motions, policies learned from demonstration typically execute with a slow and deliberate pace that feels distinctly non-human. This disparity is more than just an aesthetic issue; it is a fundamental barrier to deploying robots in practical human-centric environments where efficiency and responsiveness are important. For time-sensitive tasks like cooking and serving food, or collaborative tasks requiring synchronization with humans, the value of a robotic assistant diminishes greatly if it performs these tasks at a fraction of human speed.

The root of this speed gap lies in two challenges. First, there is a lack of accessible systems for collecting demonstration data that capture the dynamics of human-speed motion. Many teleoperation interfaces are not designed for high-speed interaction, due to latency, unintuitive control, and limited sensory feedback for human operators, causing heavy mental load for human operators to interpret feedback and control the robots. This forces human operators to move slowly and carefully, resulting in datasets of sluggish trajectories. The policy, in turn, learns to imitate this slow behavior. Second, even with a dataset of high-speed demonstrations, learning a policy that can reliably execute these motions is a

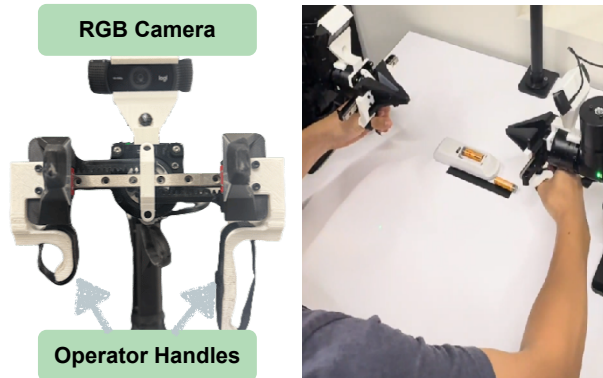


Fig. 2: **Hardware Setup.** Our printed handles allow for easy operation and data collection.

significant hurdle. Fast movements are inherently less forgiving, where small errors can quickly compound, and non-smooth actions can lead to instability. A policy must not only replicate the trajectory but also the underlying smoothness and stability that makes high-speed human motion possible.

To address these challenges, we introduce ALOHA Lightning, a system for learning bimanual manipulation near human speed. To enable data collection, we build a bimanual robot platform for intuitive kinesthetic teaching. During data collection, the end effector of a robot arm is attached with two handles for human operators to backdrive the robot’s arms and grippers for manipulation tasks. To make the robot arms feel less heavy, we implement real-time gravity compensation, which significantly reduces the physical effort required from operators. This method of direct physical guidance provides operators with rich real-time sensory feedback in vision, haptics, and sound, eliminating the delays inherent in teleoperation. Consequently, the resulting data faithfully captures the speed and fluidity of natural human motion.

To enable the robot to learn from high-speed manipulation data, we develop a learning pipeline specifically designed for speed and precision. Our approach incorporates two key techniques. First, to ensure smooth execution, we introduce test-time smoothing, a post-hoc real-time processing for smooth transitions between action chunks. This method helps mitigate jerkiness at high speed that can arise from action chunking, allowing the robot to perform fast motions without sacrificing stability. Second, to address the visual domain shift between data collection, where a human operator is present with hands holding the robot end-effectors, and autonomous execution, where the robot is alone, we employ a masking technique to align the visual domains between collection and deployment. By removing regions with human hands and arms appearing in cameras by using an analytical model,

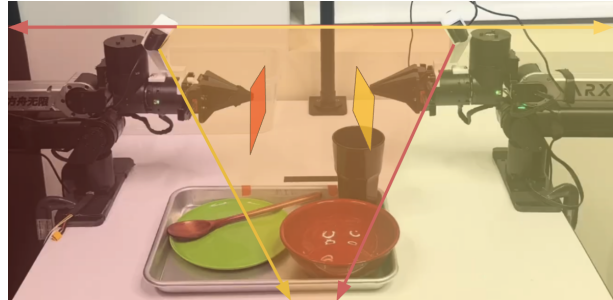


Fig. 3: **Masking.** Our masking method adjusts in real time based on the position of the end effectors of both arms.

we align the visual input distributions. This prevents the policy from being distracted by the absence of the human operator.

The synergy between our intuitive hardware for near-human-speed data collection and our robust learning pipeline allows ALOHA Lightning to learn fast, complex tasks from a small number of demonstrations. With just 50 demonstrations per task, our system can autonomously perform fine-grained tasks such as folding chino shorts, bussing a table, inserting a battery into a remote, and unstacking plastic cups, achieving over 80% success rates. Critically, these tasks are completed at speeds that are close to average human speeds, reducing the execution time by over 50% compared to that of teleoperation systems, and represent a significant step towards creating robots that can operate efficiently in the real world.

II. RELATED WORK

Learning from Demonstrations. Learning from demonstrations, or imitation learning, is a cornerstone of modern robotics, enabling robots to acquire complex skills by mimicking expert demonstrations [3]. The efficacy of this paradigm is heavily dependent on the interface used for data collection. A wide spectrum of interfaces has been explored, ranging from processing human videos to various forms of interactive teleoperation. While human videos offer immense scale and diversity, they introduce significant embodiment and observation gaps that complicate direct policy learning [4]–[6]. Teleop methods aim to reduce this gap. Common low-cost interfaces include keyboards [7], 3D space mice [8], game controllers [9] and phones [10], though they can be unintuitive for controlling 6-DoF motion. More immersive interfaces, such as Virtual Reality controllers, provide more direct 6-DoF end-effector control but can suffer from latency, a lack of haptic feedback, and a lack of feedback on kinematic constraints [11]–[14]. Other approaches shadow human motion in real time or offline using real-time pose estimation [15]–[18], exoskeletons [19], [20] and vision-based tracking [21]–[23], capturing natural kinematics but still facing the human-to-

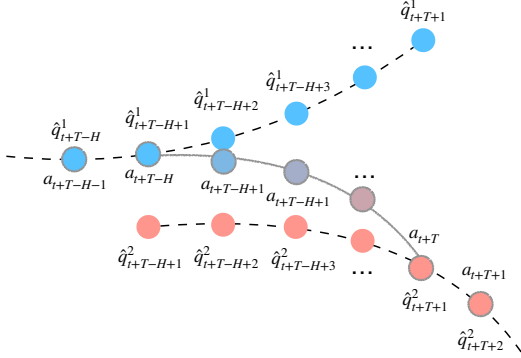


Fig. 4: **Test-Time Smoothing.** Given a previous chunk of target joint positions \hat{q}^1 and a later chunk \hat{q}^2 , we apply smoothing to obtain the final commanded actions a .

robot correspondence problem. To achieve a more direct mapping, leader-follower puppeteering systems provide secondary leader devices that the operator manipulates to control the follower robots providing intuitive joint-space control [1], [2], [24]–[28]. In contrast to these methods, ALOHA Lightning utilizes kinesthetic teaching where human operators directly hold and guide the robot motion, eliminating the cost and complexity of a duplicate leader setup for puppeteering and bypassing the latency or embodiment gap inherent in other interfaces.

Kinesthetic Teaching. Kinesthetic teaching, where an operator physically guides a robot’s end-effector through a task, is one of the most direct and intuitive methods for programming robot behaviors [3]. By moving the robot by hand, the demonstrator can implicitly convey not only the desired trajectory but also information about contact points and forces. This method requires the robot to be either passively compliant or actively backdrivable, often necessitating gravity compensation to make the robot feel “weightless” and reduce operator fatigue. Early work demonstrated the effectiveness of this approach for learning simple trajectories and dynamic movement primitives [29]–[35]. However, a primary challenge has been the difficulty of recording fast, dynamic motions due to robot inertia and motor friction, which can make the robot feel sluggish and heavy to the operator. Furthermore, the physical presence of the human operator in the workspace during demonstration can create a visual domain shift, as the policy is later deployed in a scene without the operator’s hands and arms. ALOHA Lightning addresses the challenge of collecting high-speed data by using a lightweight backdrivable bimanual platform with low-gear-ratio actuators with real-time gravity compensation, significantly reducing operator effort. Together with our masking and real-time smoothing, ALOHA Lightning can capture and learn end-to-end policies for precise and fast skills from fluid near-human-speed motions that are difficult to obtain and execute with kinesthetic teaching

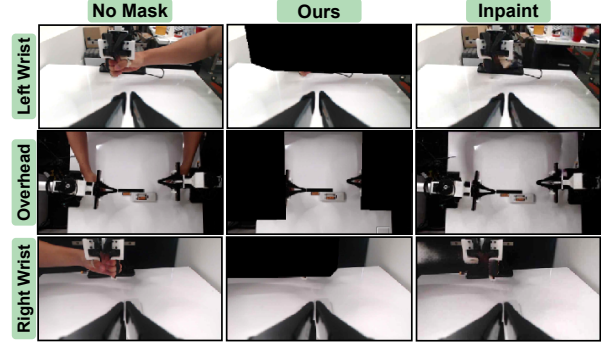


Fig. 5: **Masking Comparisons.** We show the same initial starting frame of *Battery Insertion*. Our masking method completely aligns the visual distributions of training and test times. Inpainting via diffusion creates unwanted artifacts.

systems.

High-Speed Manipulation. Achieving human-level speed in manipulation is a long-standing goal in robotics, crucial for applications where efficiency is paramount. Traditional approaches to high-speed robotics often rely on precise analytical models of robot and object dynamics, using techniques like trajectory optimization to generate feasible, time-optimal paths. These methods have enabled impressive feats, such as industrial robots performing rapid pick-and-place operations or research platforms batting [36], playing table tennis [37]–[41], moving objects and avoiding collisions at speed [42], tossing [43], dynamically manipulate deformable objects [44]–[46], and catching objects in mid-air [47]. However, these approaches typically require extensive, task-specific engineering and can be brittle to unmodeled contact dynamics or perturbations. Learning-based methods offer a promising alternative by acquiring dynamic skills directly from data. However, a primary bottleneck has been the lack of systems capable of collecting demonstration data that captures near-human-speed motion [48], [49].

Furthermore, even with a suitable dataset, learning a policy to robustly execute fast motions is a significant hurdle [21], [50], [51]. Fast movements are less forgiving, where latency discrepancies can cause out-of-sync actions, and discontinuous commands from a policy can lead to instability.

III. HARDWARE SETUP & DATA COLLECTION

We set up our hardware with the consideration of using kinesthetic teaching for fast and precise data collection. We prioritize comfort for human operators, pursuing maximal achievable speed using our hardware setup.

Shown in Fig. 2, ALOHA Lightning consists of a bimanual setup of two different types of 6-DoF robot arms, HESSIAN and ARX, each with one parallel jaw gripper. We leverage off-the-shelf robot arms using

Method	Masking																	
	Unstack Cups				Battery Insertion				Bus Table				Fold Shorts					
	Grab	Pull	Whole	Time(s)	Pick	Place	Prepare	Push	Whole	Time(s)	C&S	B&P	Whole	Time(s)	Up	Right	Whole	Time(s)
No Mask	90	90	80	4.12	95	95	85	75	75	6.44	100	100	100	6.97	100	90	90	6.19
Inpainting	100	85	85	3.92	95	95	70	65	65	5.65	100	100	100	6.57	100	80	80	6.20
Ours	100	90	90	3.97	100	100	90	85	85	5.70	100	100	100	6.81	100	100	100	5.94

TABLE I: *Our masking improves performance.* Our masking strategy almost uniformly beats non-masked and diffusion inpainted in policies, enabling tasks such as battery insertion.

planetary actuators, and human operators can back-drive the robot arms during kinesthetic teaching with relatively low internal friction in the joints, compared to others like harmonic actuators. We only use the HESSIAN setup for the highly dynamic Toss & Catch task since this task requires high joint velocity limits for high-speed trajectory following and very low joint friction for high-speed data collection.

We implement real-time gravity compensation based on the joint positions and masses of each robot link through torque control at 100Hz [52]. We retrofit the end effector of a robot arm with two handles for human operators to hold and backdrive the robots arms for completing manipulation tasks. One operator handle is attached on each jaw of the gripper. The smaller handle is for the thumb, and the other handle is for the other four fingers. Fingers are held in place with an adjustable strap, providing human operators with firm and intuitive grasps of the end effectors of the robot arms. The long finger handle, allows operators to pull downwards and sideways, while direct contact with the bottom of the arm assists with upwards acceleration. The whole setup consists of three RGB cameras capturing images at 50Hz, with one wrist camera mounted on each arm facing towards the region between the gripper jaws, and one top-down camera capturing the workspace from above.

During data collection, we record three 360 x 640 RGB images I_1, I_2, I_3 from the three cameras, joint positions q , and joint velocities \dot{q} at 50Hz. Through kinesthetic teaching where the human operators directly hold and guide the robots to do different tasks, human operators obtain various real-time sensory feedbacks, including 3D vision, haptics through the rigid bodies of the robots, and contact sound without any delays that exist in common teleoperation platforms.

IV. DATA PROCESSING, POLICY LEARNING & DEPLOYMENT

Given a dataset of kinesthetic teaching data for fast and precise manipulation, directly applying imitation learning algorithms fail for two reasons. First, a human operator is present with hands holding the robot end-effectors during data collection, while the robots are alone during the autonomous execution phase after deployment. This creates a visual distribution shift between data

collection and deployment. Second, fast movements are less forgiving, as small errors in execution can compound into unrecoverable failures and non-smooth actions can cause instability. The high-speed nature of the manipulation data collected presents a hurdle for learning a reliable policy.

To address the distribution shift between data collection and deployment, we employ a masking technique to align the visual domains for training and testing, illustrated in Fig. 3. By removing regions with human hands and arms appearing in cameras using an analytical model, we align the visual input distributions, preventing the policy from being distracted by the presence or absence of the human operator. And to ensure smooth execution of fast manipulation, we introduce test-time smoothing, a post-hoc real-time processing for smooth transitions between action chunks. Illustrated in Fig. 4, this method helps mitigate the jerkiness at high speed that can arise from action chunking [53], allowing the robot to perform fast motions without sacrificing stability.

Masking in Image Space. During training, we pre-process the collected data, aiming to mask out the appearances of hands and arms of the human operators. Given that the joint positions of both arms and robot description files are known, we can compute the poses of end-effectors of both robot arms in their respective base frames using forward kinematics (FK). We also know the transformation between the base frames of both arms through camera calibration. Thus, we can transform the end-effector pose of one arm into the base frame of the other arm. Since the relative pose of the handles for the human hands is fixed with respect to the end-effector pose, we deterministically calculate the center pose of the left hand in the right camera, and similarly for the right hand in the left camera, and mask out a small square region perpendicular to the end-effector pointing direction near the handles by setting the image RGB values for that region to be zero. For the top camera, we similarly mask out two rectangular regions for both hands and arms based on the end-effector poses. We provide an illustration for masking in Fig. 3. During deployment, we perform masking in real time for every frame. Formally, at each time step t , we perform the masking operation M and get three processed images

Method	Smoothing																	
	Unstack Cups				Battery Insertion					Bus Table				Fold Shorts				
Grab	Pull	Whole	Time(s)	Pick	Place	Prepare	Push	Whole	Time(s)	C&S	B&P	Whole	Time(s)	Up	Right	Whole	Time(s)	
None	100	65	65	4.23	100	85	75	75	75	5.29	100	80	80	6.70	95	80	80	6.16
TE	95	35	35	7.75	100	100	85	55	55	10.18	100	100	100	12.23	100	90	90	11.60
Ours	100	90	90	3.97	100	100	90	85	85	5.70	100	100	100	6.81	100	100	100	5.94

TABLE II: *Smoothing improves performance.* Test time smoothing outperforms or matches other methods in success rate in every category

$I'_{1,t}, I'_{2,t}, I'_{1,t}$, where

$$I'_{1,t}, I'_{2,t}, I'_{1,t} = M(I_{1,t}, I_{2,t}, I_{1,t}, q_t, \text{FK}).$$

Policy Learning. Unlike teleoperation where the action commands for the robots are explicitly sent through devices like leader arms or VR controllers, we only capture the visual and proprioceptive readings from the robots. We aim to train policies that can predict future joint positions which the robots should achieve given current images and proprioception. We follow the recent progress in imitation learning by using action chunking [1], which mitigates compounding error problems that is common in behavior cloning methods and improves overall performance. Given a policy π and an action chunk size of T , we train the policy by maximizing the probability

$$\pi([q_{t+1}, \dots, q_{t+T+1}] | I'_{1,t}, I'_{2,t}, I'_{1,t}, q_t),$$

where the predicted future joint positions will serve as the target joint positions for a low-level PD controller to execute. In this paper, we use ACT [1] as the policy representation, but we expect other policy representations for imitation learning with action chunks would also work.

Test-Time Smoothing. When moving at increased speeds, we observe that discontinuities between action chunks [53] become larger and more visible, affecting policy performance. We introduce a simple post-hoc test-time smoothing technique. Illustrated in Fig. 4, we tackle the discontinuity between two adjacent chunks by blending the last few actions of a chunk with the first few actions of the later chunk where these two action chunks have a small overlap in corresponding timesteps. For a smoothing window size of H smaller than the action chunk size T , an action chunk $[\hat{q}_{t+1}^1, \dots, \hat{q}_{t+T+1}^1]$ is first predicted at timestep t , and another one $[\hat{q}_{t+T-H+1}^2, \dots, \hat{q}_{t+2T-H+1}^2]$ is later predicted at timestep $t+T-H$. At timestep t' , the action command $a_{t'}$ is the raw predicted target joint position $\hat{q}_{t'+1}$ at regions where there is no chunk overlap. Within the overlapping regions, we weigh the two corresponding target joint positions by using a sigmoid coefficient $w(t') = \text{sigmoid}((t' - (t + T - H)) - H/2)/\sigma)$ where σ is a time scale, and the final blended action command

is

$$a_{t'} = (1 - w(t')) * \hat{q}_{t'+1}^1 + w(t') * \hat{q}_{t'+1}^2.$$

By adjusting the smoothing horizon and time scale of the sigmoid coefficient, we can adjust how sharply and reactively the robot will change to follow a later action chunk from a previous obsolete action chunk.

V. EXPERIMENTS

In our experiments, we aim to determine the extent to which our masking strategy as well as our smoothing method benefit the overall success rate of policies. Additionally, we study the effects of several comparative approaches in these two domains. Finally, through user trials, we verify the intuitive understanding that our directly operated data collection platform enables faster training of human operators. All tasks are shown in Fig. 1.

A. Tasks

In *Unstack Cups*, the left arm reaches out and grasps the bottom of two stacked red solo cups. Then, the right arm subsequently grasps the lip of the top cup without contacting the bottom, and pulls. The grippers must carefully balance their grasps; if the left grips too hard, the right can't tilt and grab. If the right over-grips, it pulls the entire stack. The cups' initial starting positions are randomized along a fifteen centimeter piece of tape. This task is collected with a timer of four seconds.

Battery Insertion. During execution, one arm picks up a battery and places it at an oblique angle into a remote control, taking care not to touch the spring. It then pushes one end of the battery into the spring. The other gripper then raises and forcefully pushes down on the battery, slotting it into place. The battery and the remote are randomized along the tape. This task is collected in under seven seconds.

Bus Table involves removing common dining ware from a pan and placing them into a tray. An arm first grasps a nearby tray and brings it closer to the tray. Then, both arms engage, alternating picking up and placing four items from the pan into the tray: a cup, a bowl, a spoon, and a plate. The position of the pan is randomized along the tape. This task is collected with a timer of six seconds.

Fold Shorts involves folding a pair of shorts placed

Method	Speed Comparisons																	
	Unstack Cups				Battery Insertion				Bus Table				Toss & Catch					
	Grab	Pull	Whole	Time(s)	Pick	Place	Prep	Push	Whole	Time(s)	C&S	B&P	Whole	Time(s)	Pick & Toss	Catch	Whole	Time(s)
Puppeteering	85	75	75	8.84	100	90	60	55	55	17.87	100	90	90	16.05	-	-	-	-
Puppet Double Freq	80	65	65	5.47	80	65	35	15	15	12.90	95	85	85	9.30	-	-	-	-
Puppet Skip Actions	85	50	50	5.58	95	95	55	30	30	10.30	95	85	85	9.29	-	-	-	-
Ours	100	90	90	3.97	100	100	90	85	85	5.70	100	100	100	6.81	100	50	50	1.85

TABLE III: *Speed Comparisons*. We note a sharp drop in performance on precise tasks when applying either speed-up approach.

flat. The arms must reach and grasp the edge of the bottom trunk of the shorts and toss the trunk over from bottom to top flush against the top cover. Then, the left gripper must grasp the edge of the shorts along the waist, and toss it from left to right while the right arm creates a crease. The shorts are randomly rotated thirty degrees in its initial position, and this task is collected against a timer of five seconds.

Toss & Catch. The left arm first approaches and picks up a small soft ball on the table. Then it swings forward and up dynamically in high speed while releasing the ball. The right arm approaches the tossed ball-in-air swiftly with closing grippers to catch the ball. The task last for merely less than 2 seconds.

B. Comparisons

We compare our method to *Inpainting*, *No Mask*, *Temporal Ensembling* (TE), *No Smooth*, *Puppeteering*, *Puppet Double Freq*, and *Puppet Skip Actions*.

Inpainting [54]. A policy trained on diffusion inpainted video, where the operator hands are segmented and replaced with the background by a pretrained segmentation [55] and pretrained diffusion network [56].

No Mask. A policy trained with the raw video.

Temporal Ensembling [1]. A post-process smoothing method where an action chunk is predicted every timestep, and exponentially weighted and averaged with past action chunks.

No Smooth. No smoothing method applied.

Puppeteering [1]. The operator collects and deploys policies on a teleoperated setup.

Puppet Double Freq. On the puppeteering setup, the policy is deployed with double the controller frequency.

Puppet Skip Actions. On the puppeteering setup, the policy acts on every other predicted action per chunk.

C. Masking Aligns Data Collection and Deployment

We evaluate the efficacy of masking, which aligns the visual distribution of the policy between collection and deployment time. Table I details the results of our experimental trials comparing policies trained with masking, no masking, and on diffusion inpainted video. We note an increase in performance on most tasks with masking compared to simply training with hands included and inpainted video.

The policies trained on inpainted video suffer signifi-

cant performance losses compared to the base video. We hypothesize that this is due to two factors. First, as can be seen in Fig. 5, the shadows of the arms are still visible against the reflective surface of the table. Another factor is the noisy diffusing of frames due to lack of context. These visual artifacts not visible at test time may confuse the policy. Given a more controlled environment, this method’s performance may improve.

D. Test-Time Smoothing Improves Performance

In Fig. 4, we see performance comparisons between non-smoothed, test-time smoothed, and temporal aggregation-smoothed trajectories. In all tasks, test-time smoothing achieves the best performance, and by a wide margin in tasks *Unstack Cups* and *Battery Insertion*. The resulting instabilities in the high speed non-smoothed policy rollouts result in imprecise grasps and pushes, and while the temporal-aggregation smoothed policies have mixed results, their execution speeds drop in half dramatically from collection time. This is due to the increased processing time necessary to compute inference every timestep.

A failure mode of temporal aggregation is its tendency to choose an averaged path when forced to change course [53]. We hypothesize that, at higher speeds, the tracking of robot controllers often become less reliable which forces policies to react and frequently change course to adjust for error. Thus, even if temporal aggregation were run on faster hardware and had no inference delay, it would perform even more poorly at speed due to this phenomenon. This also partially explains the presence of discontinuities in the first place, even in relatively uni-modal data with little strategy-switching.

E. Speed Comparisons

In this experiment, we investigate the effectiveness of training a policy on slow data, collected with puppeteering. We artificially increase the policy execution speed in two different ways across three different tasks: *Bus Table*, *Unstack Cups*, and *Battery Insertion*. We are unable to complete the *Toss & Catch* task with puppeteering due to the difficulty of catching a dynamic object through teleoperation. The results are visible in Table III, where we see sharp drops in performance after speeding up data collected through puppeteering to match the speed

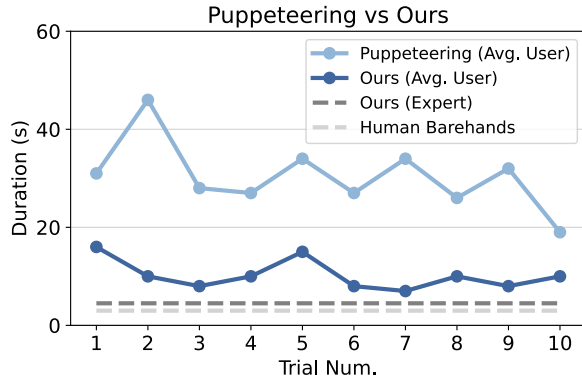


Fig. 6: *User Study*. Results on *Battery Insertion* for timing operators on both Ours and Puppeteering demonstrate a 50%+ gap in duration for both experienced and new operators.

of ours.

Puppet Double Freq simply doubles the controller frequency at execution time. This results in policies that appear out of control, with large jumps at chunk boundaries. With actions that jump larger distances, such as in *Bus Table*, this effect becomes more noticeable. Despite still moving slower than our method, double-timed policies from puppeteering have significantly larger discontinuities than in our method. We hypothesize that these discontinuities lead to deviations between collection and test time, further worsening performance.

Puppet Skip Actions is running the controller frequency at the same rate, but skipping every other action. This method does not introduce discontinuity, but creates new failure modes such as oscillating above the battery instead of pushing it in, and missing grasps. This is due to the policy oftentimes skipping critical steps, which also reduces overall performance.

F. User Study

We conduct a user study with seven volunteers with little to no teleoperation experience, having them operate our setup and puppeteering. For the study, we choose *Battery Insertion*, a task requiring precision.

In the study, in both setups, users spend three minutes familiarizing themselves with the objects in the scene and the robotic arms. Next, they are shown an expert demonstration and given basic tips on operation procedures. Finally, they are timed in ten trials, attempting to insert the battery into the remote as quickly as possible.

We observe in Fig. 6 an ability for users to quickly learn how to operate with our hardware and at speeds up to three times as fast as puppeteering. This primarily results from the operator sitting closer to the environment, giving more direct visual feedback, and haptic sensations from operating the actual robot arms. This results in more natural data collection, even for beginners.

VI. CONCLUSION AND ACKNOWLEDGEMENT

In this paper, we present ALOHA Lightning, a full-stack system for learning fast and precise manipulation. Our system presents a hardware setup for collecting high-speed demonstration, and a learning pipeline for smooth autonomous execution of fast and precise movements. Currently, the masking strategy that we use in the paper over-conservatively masks regions near human hands and arms, erasing potentially useful visual cues near human hands in the image space. Additionally, the whole system incorporates vision and proprioception as only two modalities, preventing the system from learning dexterous manipulation requiring contact-rich perception and control. In the future we would like to research adaptive masking strategies and including tactile sensing for better performance. This work was supported by the Robotics and AI Institute. Zipeng Fu was supported by Pierre and Christine Lamond Fellowship. We appreciate the help from HESSIAN and ARX.

REFERENCES

- [1] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn, "Learning fine-grained bimanual manipulation with low-cost hardware," *RSS*, 2023.
- [2] Z. Fu, T. Z. Zhao, and C. Finn, "Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation," in *CoRL*, 2024.
- [3] A. Billard, S. Calinon, R. Dillmann, and S. Schaal, "Robot programming by demonstration," in *Springer handbook of robotics*. Springer, 2008, pp. 1371–1394.
- [4] S. Kareer, D. Patel, R. Punamiya, P. Mathur, S. Cheng, C. Wang, J. Hoffman, and D. Xu, "Egomimic: Scaling imitation learning via egocentric video," in *2025 ICRA*, 2025.
- [5] M. Lepert, J. Fang, and J. Bohg, "Phantom: Training robots without robots using only human videos," *arXiv preprint arXiv:2503.00779*, 2025.
- [6] L. Y. Zhu, P. Kuppili, R. Punamiya, P. Aphiwetsa, D. Patel, S. Kareer, S. Ha, and D. Xu, "Emma: Scaling mobile manipulation via egocentric human data," *arXiv preprint arXiv:2509.04443*, 2025.
- [7] H. Xiong, R. Mendonca, K. Shaw, and D. Pathak, "Adaptive mobile manipulation for articulated objects in the open world," *arXiv preprint arXiv:2401.14403*, 2024.
- [8] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu, "Viola: Imitation learning for vision-based manipulation with object proposal priors," in *CoRL*. PMLR, 2023, pp. 1199–1210.
- [9] S. Ross, N. Melik-Barkhudarov, K. S. Shankar, A. Wendel, D. Dey, J. A. Bagnell, and M. Hebert, "Learning monocular reactive uav control in cluttered natural environments," in *2013 ICRA*, 2013.
- [10] A. Mandlekar, Y. Zhu, A. Garg, J. Booher, M. Spero, A. Tung, J. Gao, J. Emmons, A. Gupta, E. Orbay *et al.*, "Roboturk: A crowdsourcing platform for robotic skill learning through imitation," in *CoRL*. PMLR, 2018, pp. 879–893.
- [11] T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel, "Deep imitation learning for complex manipulation tasks from virtual reality teleoperation," in *2018 ICRA*, 2018.
- [12] X. Cheng, J. Li, S. Yang, G. Yang, and X. Wang, "Open-television: Teleoperation with immersive active visual feedback," *arXiv preprint arXiv:2407.01512*, 2024.
- [13] Y. Park, J. S. Bhatia, L. Ankile, and P. Agrawal, "Dexhub and dart: Towards internet scale robot data collection," *arXiv preprint arXiv:2411.02214*, 2024.
- [14] M. Seo, S. Han, K. Sim, S. H. Bang, C. Gonzalez, L. Sentis, and Y. Zhu, "Deep imitation learning for humanoid loco-manipulation through human teleoperation," in *2023 Humanoids*, 2023.

- [15] A. Handa, K. Van Wyk, W. Yang, J. Liang, Y.-W. Chao, Q. Wan, S. Birchfield, N. Ratliff, and D. Fox, "Dexpilot: Vision-based teleoperation of dexterous robotic hand-arm system," in *2020 ICRA*, 2020.
- [16] A. Sivakumar, K. Shaw, and D. Pathak, "Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube," *arXiv preprint arXiv:2202.10448*, 2022.
- [17] Z. Fu, Q. Zhao, Q. Wu, G. Wetzstein, and C. Finn, "Humanplus: Humanoid shadowing and imitation from humans," in *CoRL*, 2024.
- [18] T. He, Z. Luo, W. Xiao, C. Zhang, K. Kitani, C. Liu, and G. Shi, "Learning human-to-humanoid real-time whole-body teleoperation," *arXiv preprint arXiv:2403.04436*, 2024.
- [19] H. Fang, H.-S. Fang, Y. Wang, J. Ren, J. Chen, R. Zhang, W. Wang, and C. Lu, "Airexo: Low-cost exoskeletons for learning whole-arm manipulation in the wild," in *2024 ICRA*, 2024.
- [20] A. Purushottam, C. Xu, Y. Jung, and J. Ramos, "Dynamic mobile manipulation via whole-body bilateral teleoperation of a wheeled humanoid," *IEEE Robotics and Automation Letters*, 2023.
- [21] C. Chi, Z. Xu, C. Pan, E. Cousineau, B. Burchfiel, S. Feng, R. Tedrake, and S. Song, "Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots," in *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- [22] C. Wang, H. Shi, W. Wang, R. Zhang, L. Fei-Fei, and C. K. Liu, "Dexcap: Scalable and portable mocap data collection system for dexterous manipulation," *arXiv preprint arXiv:2403.07788*, 2024.
- [23] H. Etukuru, N. Naka, Z. Hu, S. Lee, J. Mehu, A. Edsinger, C. Paxton, S. Chintala, L. Pinto, and N. M. M. Shafiqullah, "Robot utility models: General policies for zero-shot deployment in new environments," in *2025 ICRA*, 2025.
- [24] R. Jenness and C. Wicker, "Master-slave manipulators and remote maintenance at the oak ridge national laboratory," Oak Ridge National Lab., Tenn.(USA), Tech. Rep., 1975.
- [25] P. Wu, Y. Shentu, Z. Yi, X. Lin, and P. Abbeel, "Gello: A general, low-cost, and intuitive teleoperation framework for robot manipulators," in *2024 IROS*, 2024.
- [26] T. Buamane, M. Kobayashi, Y. Uranishi, and H. Takemura, "Bi-act: Bilateral control-based imitation learning via action chunking with transformer," in *2024 AIM*, 2024.
- [27] H. Ito, K. Yamamoto, H. Mori, and T. Ogata, "Efficient multitask learning with an embodied predictive model for door opening and entry with whole-body control," *Science Robotics*, 2022.
- [28] Y. Jiang, R. Zhang, J. Wong, C. Wang, Y. Ze, H. Yin, C. Gokmen, S. Song, J. Wu, and L. Fei-Fei, "Behavior robot suite: Streamlining real-world whole-body manipulation for everyday household activities," *arXiv preprint arXiv:2503.05652*, 2025.
- [29] C. Finn, S. Levine, and P. Abbeel, "Guided cost learning: Deep inverse optimal control via policy optimization," in *International conference on machine learning*, 2016.
- [30] A. Xie, F. Ebert, S. Levine, and C. Finn, "Improvisation through physical understanding: Using novel objects as tools with visual foresight," *arXiv preprint arXiv:1904.05538*, 2019.
- [31] B. Akgun, M. Cakmak, J. W. Yoo, and A. L. Thomaz, "Trajectories and keyframes for kinesthetic teaching: A human-robot interaction perspective," in *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*, 2012.
- [32] P. Kormushev, D. N. Nenchev, S. Calinon, and D. G. Caldwell, "Upper-body kinesthetic teaching of a free-standing humanoid robot," in *2011 ICRA*, 2011.
- [33] P. Kormushev, S. Calinon, and D. G. Caldwell, "Imitation learning of positional and force skills demonstrated via kinesthetic teaching and haptic input," *Advanced Robotics*, 2011.
- [34] A. Weiss, J. Igelsbock, S. Calinon, A. Billard, and M. Tscheligi, "Teaching a humanoid: A user study on learning by demonstration with hoap-3," in *RO-MAN 2009-The 18th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 2009.
- [35] H. Ben Amor, E. Berger, D. Vogt, and B. Jung, "Kinesthetic bootstrapping: Teaching motor skills to humanoid robots through physical interaction," in *Annual conference on artificial intelligence*. Springer, 2009.
- [36] T. Senoo, A. Namiki, and M. Ishikawa, "Ball control in high-speed batting motion using hybrid trajectory generator," in *Proceedings 2006 ICRA.*, 2006.
- [37] D. B. D'Ambrosio, S. Abeyruwan, L. Graesser, A. Iscen, H. B. Amor, A. Bewley, B. J. Reed, K. Reymann, L. Takayama, Y. Tassa *et al.*, "Achieving human level competitive robot table tennis," in *2025 ICRA*, 2025.
- [38] D. Büchler, S. Guist, R. Calandra, V. Berenz, B. Schölkopf, and J. Peters, "Learning to play table tennis from scratch using muscular robots," *IEEE Transactions on Robotics*, 2022.
- [39] K. Mülling, J. Kober, O. Kroemer, and J. Peters, "Learning to select and generalize striking movements in robot table tennis," *The International Journal of Robotics Research*, 2013.
- [40] Y. Zhu, Y. Zhao, L. Jin, J. Wu, and R. Xiong, "Towards high level skill learning: Learn to return table tennis ball using monte-carlo based policy gradient method," in *2018 IEEE international conference on real-time computing and robotics (RCAR)*, 2018.
- [41] C. Liu, Y. Hayakawa, and A. Nakashima, "Racket control for a table tennis robot to return a ball," *SICE Journal of Control, Measurement, and System Integration*, 2013.
- [42] J. Ichnowski, Y. Avigal, Y. Liu, and K. Goldberg, "Gomp-fit: Grasp-optimized motion planning for fast inertial transport," in *2022 iCRA*. IEEE, 2022, pp. 5255–5261.
- [43] A. Zeng, S. Song, J. Lee, A. Rodriguez, and T. Funkhouser, "Tossingbot: Learning to throw arbitrary objects with residual physics," *IEEE Transactions on Robotics*, 2020.
- [44] H. Zhang, J. Ichnowski, D. Seita, J. Wang, H. Huang, and K. Goldberg, "Robots of the lost arc: Self-supervised learning to dynamically manipulate fixed-endpoint cables," in *2021 ICRA*. IEEE, 2021, pp. 4560–4567.
- [45] C. Chi, B. Burchfiel, E. Cousineau, S. Feng, and S. Song, "Iterative residual policy: for goal-conditioned dynamic manipulation of deformable objects," *The International Journal of Robotics Research*, 2024.
- [46] H. Ha and S. Song, "Flingbot: The unreasonable effectiveness of dynamic manipulation for cloth unfolding," in *CoRL*, 2022.
- [47] S. Kim, A. Shukla, and A. Billard, "Catching objects in flight," *IEEE Transactions on Robotics*, 2014.
- [48] N. Masuya, H. Sato, K. Yamane, T. Kusume, S. Sakaino, and T. Tsuji, "Variable-speed teaching-playback as real-world data augmentation for imitation learning," *Advanced Robotics*, 2025.
- [49] L. Guo, Z. Xue, Z. Xu, and H. Xu, "Demospeedup: Accelerating visuomotor policies via entropy-guided demonstration acceleration," *arXiv preprint arXiv:2506.05064*, 2025.
- [50] N. R. Arachchige, Z. Chen, W. Jung, W. C. Shin, R. Bansal, P. Barroso, Y. H. He, Y. C. Lin, B. Joffe, S. Kousik *et al.*, "Sail: Faster-than-demonstration execution of imitation learning policies," *arXiv preprint arXiv:2506.11948*, 2025.
- [51] D. D. Yuan, T. Z. Zhao, K. Burns, and C. Finn, "Speedtuning: Speeding up policy execution with lightweight reinforcement learning," in *2025 ICRA*, 2025.
- [52] K. M. Lynch and F. C. Park, *Modern robotics*. Cambridge University Press, 2017.
- [53] K. Black, M. Y. Galliker, and S. Levine, "Real-time execution of action chunking flow policies," *arXiv preprint arXiv:2506.07339*, 2025.
- [54] M. Xu, H. Zhang, Y. Hou, Z. Xu, L. Fan, M. Veloso, and S. Song, "Dexumi: Using human hand as the universal manipulation interface for dexterous manipulation," *arXiv preprint arXiv:2505.21864*, 2025.
- [55] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár, and C. Feichtenhofer, "Sam 2: Segment anything in images and videos," *arXiv preprint arXiv:2408.00714*, 2024.
- [56] S. Zhou, C. Li, K. C. Chan, and C. C. Loy, "ProPainter: Improving propagation and transformer for video inpainting," in *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, 2023.