

# When Life Gives You BC, Make Q-functions: Extracting Q-values from Behavior Cloning for On-Robot Reinforcement Learning

Author Names Omitted for Anonymous Review. Paper-ID 779

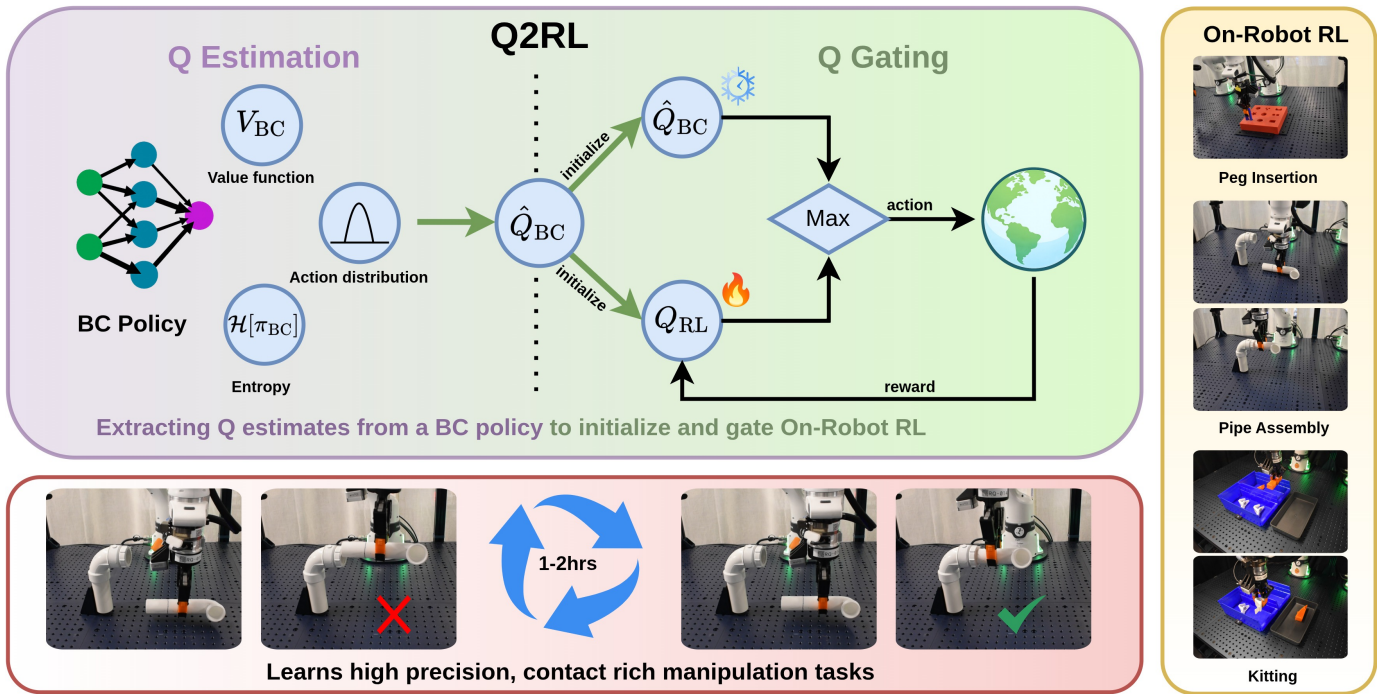


Fig. 1: *Q2RL* consists of Q-Estimation and Q-Gating. Q-Estimation extracts a Q function from a BC Policy. Then, during RL policy training, Q-Gating selects and executes the BC or RL action with highest respective Q value, updating the RL policy on the collected interactions. *Q2RL* learns contact-rich manipulation tasks in 1-2 hours of online interaction.

**Abstract**—Behavior Cloning (BC) has emerged as a highly effective paradigm for robot learning. However, BC lacks a self-guided mechanism for online improvement after demonstrations have been collected. Existing offline-to-online learning methods often cause policies to replace previously learned good actions due to a distribution mismatch between offline data and online learning. In this work, we propose *Q2RL*, Q-Estimation and Q-Gating from BC for Reinforcement Learning, an algorithm for efficient offline-to-online learning. Our method consists of two parts: (1) *Q-Estimation* extracts a Q function from a BC policy using a few interaction steps with the environment, followed by online RL with (2) *Q-Gating*, which switches between BC and RL policy actions based on their respective Q values to collect samples for RL policy training. Across manipulation tasks from D4RL and robomimic benchmarks, *Q2RL* outperforms SOTA offline-to-online learning baselines on success rate and time to convergence. *Q2RL* is efficient enough to be applied in an on-robot RL setting, learning robust policies for contact-rich and high precision manipulation tasks such as pipe assembly and kitting, in 1-2 hours of online interaction, achieving success rates of up to 100% and up to 3.75x improvement against the original BC policy.

Behavior Cloning (BC) has been highly successful in robot learning due to its simple supervised training objective, which directly models an action distribution conditioned on observed states. However, BC does not have a mechanism for self-guided online improvement. When the training data for a BC policy is sufficient for the scope of the task, and/or changes in the environment induce even subtle distribution shifts, BC policy performance can collapse due to compounding errors from covariate shift [27].

Several research directions have been proposed to address the shortcomings of behavior cloning. Interactive imitation learning approaches finetune BC policies online, but require expert annotators [13]. Offline reinforcement learning (RL) learns a value function offline [1, 14] to guide online learning, but require pretraining on large datasets of positive and negative offline samples [16] for good performance. Offline-to-online learning approaches have also been proposed, but they often result in unlearning good BC actions [24, 35]. Another challenge faced by offline-to-online methods is the tendency

to produce unsafe behaviors due to training an RL policy from scratch. In the on-robot RL setting in which online RL is trained directly on the robot, executing random exploratory actions is unsafe, causing significant wear and tear on the robot, workspace, and manipulated objects.

In this paper, our objective is to combine the benefits of pre-training a BC policy with the iterative improvement capability of online RL. We aim to keep good BC policy behavior and use RL to explore when BC performs poorly. We propose Q-Estimation and Q-Gating from BC for Reinforcement Learning (*Q2RL*), an algorithm for efficient offline-to-online learning. In the *Q-Estimation* phase, we estimate the Q function of the pretrained BC policy from initial online rollouts, then use it to initialize off-policy reinforcement learning. We are able to derive these estimated Q values with just the action selection probability and entropy of the BC policy, without any additional information of the underlying policy class or training distribution. To ensure a stable transition from offline pretraining to online RL, we propose *Q-Gating*, which takes the Q function from the Q-Estimation phase and uses it to guide online reinforcement learning. Two copies of the Q function are used, one fixed as the initially extracted Q function to preserve good BC actions, and a learnable RL Q function to enable exploration of potentially better actions. At each state, *Q2RL* compares these Q estimates to selectively execute BC or RL actions, ensuring that the RL policy is able to initialize learning from good BC actions and explore online to improve beyond the BC policy.

We provide extensive experiments in simulation and real-world on-robot RL to validate the method. *Q2RL* outperforms SOTA offline RL and offline-to-online methods across D4RL and robomimic benchmark tasks, improving the success rates from 50-60% for the original BC policy up to 80-100% for most environments. Ablation experiments show that both Q-Estimation and Q-Gating are valuable components of the algorithm. Our real-world experiments show that *Q2RL* learns contact-rich, high-precision manipulation tasks in 1-2 hours of online interaction, outperforming baselines especially as tasks increase in difficulty, and improving the success rate compared to the original BC policies by up 100% and up to 3.75x improvement. Furthermore, our real-world experiments demonstrate that *Q2RL* can recover from distribution shifts in the environment.

Our contributions are as follows:

- “Q-Estimation”, an algorithm for estimating a Q function from a behavior cloning policy and a few online rollouts. Q-Estimation can be applied to any policy class that provides action likelihoods and entropy.
- “Q-Gating”, an algorithm that guides online reinforcement learning via a gating mechanism to select between BC and RL actions.
- Experiments demonstrating that *Q2RL*, which combines Q-Estimation and Q-Gating, outperforms baselines on D4RL and robomimic benchmarks.
- Real-world, on-robot reinforcement learning experiments with *Q2RL* that achieve performance over the original BC

policy on contact-rich manipulation tasks in 1-2 hours of environment interaction, and adapt to task distribution shifts unseen by the BC policy.

## I. RELATED WORK

### A. Offline Learning

Offline learning in robotics is broadly divided into Imitation Learning based approaches [34, 22] and Offline Reinforcement Learning approaches [14, 16]. Imitation learning approaches that train policies via supervised learning (i.e., Behavior Cloning) learn a direct mapping from states to an action distribution using demonstration data, while Offline RL methods learn value functions that estimate expected returns from data generated by a behavior policy. Offline RL techniques like CQL [14] penalize out-of-distribution actions to address dataset shift, but this can introduce pessimistic bias in value estimates. Furthermore, offline RL methods are more suited to learn from a large and potentially suboptimal datasets [8], and not necessarily from robotics datasets where successful demonstrations are easier to curate [25, 4]. BC approaches can further be divided into policy classes that explicitly parameterize the action distribution, such as Gaussian or Gaussian mixture heads paired with any network backbone [23, 29], and policy classes that implicitly represent the action distribution using generative denoising processes such as diffusion [5, 31] or flow-based models [17, 11]. *Q2RL* only requires access to the action selection probabilities and the entropy, which is easier to obtain for policy classes that explicitly parameterize the action distribution.

### B. Offline-to-Online Learning

Offline learning can provide strong initial performance from readily available datasets, but policies trained purely offline cannot improve further without online interaction or human intervention [13, 28]. Calibrated Q-Learning (CalQL) [24] aims to enable a smooth transition from offline to online RL by calibrating offline Q-values against a reference policy, but can struggle with limited offline data. WSRL [35] mitigates this by adding a warm-up phase that seeds the online replay buffer with rollouts from the offline policy. In contrast, other methods bypass offline RL by first pretraining an IL policy, then updating the behavior of the agent by adding corrective residual terms to the output action [30, 12, 33, 6]. These approaches require task-specific tuning of the residual action scale to ensure that the residual policy does not deviate from the BC policy and fail to optimize in a large action parameterization. *Q2RL* does not use a residual formulation and allows the RL policy to learn completely different actions from the original BC policy, while still relying on the BC policy for guidance. Imitation Bootstrapped RL (IBRL) [10], most closely related to our approach, also combines a pretrained IL/BC policy with an RL policy and selects between their actions using a randomly initialized critic. *Q2RL* estimates the Q-function of the BC policy from early online interaction and uses it in conjunction with an RL Q function to evaluate actions during learning. In our evaluations, we show that *Q2RL* outperforms

IBRL in simulated and real-world experiments, and that *Q2RL* does not require seeding the online replay buffer with high quality demonstrations, unlike IBRL.

### C. On-Robot Reinforcement Learning

Online reinforcement learning on a real robot is a challenging but necessary process to achieve robust robotic policies. SERL [19] introduces a system designed specifically for on-robot RL. It provides a software-suite to run popular online RL algorithms like SAC [9] and DrQ [32] using an async actor and learner. Since random exploration is costly on hardware, practical methods must begin with reasonable initial performance and adapt sample-efficiently. Extensions like HIL-SERL [21] incorporate human interventions to further improve efficiency, while other approaches like GCR [3] use reward shaping to guide learning. In contrast, *Q2RL* leverages readily available BC policies to provide strong initial behavior without additional reward engineering or human intervention. Although IBRL [10] also bootstraps RL from a BC policy, it starts with low initial performance that can induce aggressive exploration. *Q2RL* instead maintains a Q-estimate of the BC policy during online learning, recovering baseline performance and gates the exploration of the RL policy, making it more sample-efficient and practical for real-robot deployment.

## II. PROBLEM FORMULATION

Assume we are given a behavior cloning policy  $\pi_{\text{BC}}$ , pre-trained on a dataset  $\mathcal{D}$  of successful human demonstrations for a task  $\mathcal{M}_{\text{train}} \in \mathfrak{M}$ . Here  $\mathfrak{M}$  denotes a class of tasks that share the same success condition; in the case of task distribution shift,  $\mathcal{M}_{\text{test}}$  and  $\mathcal{M}_{\text{train}}$  are drawn from task class  $\mathfrak{M}$ . The dataset consists of  $M$  trajectories of observation-action tuples:

$$\mathcal{D} = \{\tau_j\}_{j=1}^M, \quad \tau_j = \{(o_{j,t}, a_{j,t})\}_{t=1}^{T_j}, \quad (1)$$

where  $T_j$  is the length of trajectory  $\tau_j$ , and  $o_{j,t}, a_{j,t}$  are the observation and action, respectively.

Assume the BC policy  $\pi_{\text{BC}}$  perform sub-optimally on the task (i.e. failing to achieve the success condition), either because the amount of training data or number of training epochs were insufficient, or due to distribution shifts induced by changes in environment conditions such as lighting variations, the presence of visual distractors, modified robot hardware or controllers, or other changes to environment dynamics.

Our objective is to improve performance on task  $\mathcal{M}_{\text{test}}$  beyond the original BC policy’s performance. To achieve this, we assume access to the environment for online interactions and a sparse reward signal. For this online phase, we formulate the task as a Markov Decision Process (MDP)  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \rho_0, \gamma\}$ , where  $\mathcal{S} \in \mathbb{R}^n$  and  $\mathcal{A} \in \mathbb{R}^m$  represent states and actions,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function,  $\rho_0$  is the probability distribution over initial states, and  $\gamma \in [0, 1)$  is a discount factor.

### A. Background

Our proposed approach will involve estimating a Q function, or the expected return from taking an action  $a \in \mathcal{A}$  from state  $s \in \mathcal{S}$ , then following a policy  $\pi$  thereafter:

$$Q_\pi(s, a) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right]. \quad (2)$$

The value function for policies is defined as the expectation over Q values:

$$V(s) = \mathbb{E}_{a \sim \pi} [Q(s, a)]. \quad (3)$$

Although  $Q_\pi(s, a)$  is a standard action-value function, it can also be interpreted as defining an energy-based model over actions. In particular, treating  $E(s, a) = -Q(s, a)$  as an energy yields a Boltzmann policy

$$\pi(a \mid s) \propto \exp \left( \frac{1}{\alpha} Q(s, a) \right), \quad (4)$$

where  $\alpha$  is a temperature parameter.

## III. Q2RL: Q-ESTIMATION AND Q-GATING FROM BEHAVIOR CLONING FOR REINFORCEMENT LEARNING

In this section, we present *Q2RL*, an algorithm for estimating Q values from a BC policy to guide online reinforcement learning. *Q2RL* consists of two phases: Q-Estimation (Sec. III-A) and Q-Gating (Sec. III-B). See Alg. 1 for a high-level description and Fig. 1 for a visual overview.

---

### Algorithm 1 Q2RL

---

- 1: **Given:** behavior cloning policy  $\pi_{\text{BC}}$
  - 2: Roll out  $\pi_{\text{BC}}$  for  $N$  environment steps
  - 3: Estimate  $\hat{Q}_{\text{BC}}$  using Eq. 6 ▷ Q-Estimation
  - 4: Initialize  $Q_{\text{RL}}$  with  $\hat{Q}_{\text{BC}}$ , Initialize  $\pi_{\text{RL}}$
  - 5: **for**  $t = 0 \dots T$  environment steps **do**
  - 6:    $a_{\text{BC}} \sim \pi_{\text{BC}}(s_t)$
  - 7:    $a_{\text{RL}} \sim \pi_{\text{RL}}(s_t)$
  - 8:   Select  $a_t$  using Eq. 10 ▷ Q Gating
  - 9:   Observe next state  $s_{t+1}$  and reward  $r_t$
  - 10:   Store  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $\mathcal{B}$
  - 11:   Sample a minibatch  $\{(s, a, r, s')\} \in \mathcal{B}$
  - 12:   Update  $Q_{\text{RL}}, \pi_{\text{RL}}$  using off-policy RL
  - 13: **end for**
- 

### A. Q-Estimation from a Behavior Cloning Policy

In the first phase of our method, we estimate the Q-function of a pretrained BC policy  $\pi_{\text{BC}}$ . We first derive an analytic formula for the Q-function in terms of the value function, policy likelihoods, and entropy. We assume the action distribution of the human demonstrations  $\mathcal{D}$  used to train  $\pi_{\text{BC}}$  is approximately a Boltzmann distribution, which have been used to model human actions in prior work [15, 18]. Assuming the BC policy’s action distribution  $\pi_{\text{BC}}(a \mid s)$  can



Fig. 2: Tasks from D4RL (Kitchen-Complete, Adroit-Pen, Adroit-Door) and robomimic (Lift, Can, Square).

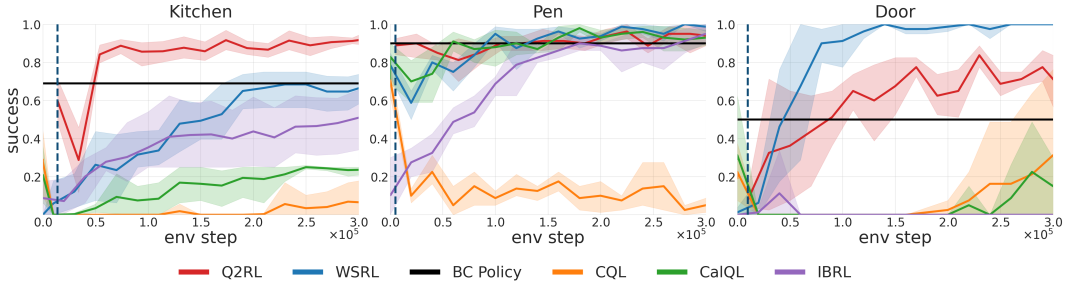


Fig. 3: Results on D4RL. Plots begin when a method starts online interaction. The dotted blue line signifies the end of the Q-Estimation phase and the start of online RL with Q-Gating for our method.

be expressed as a Boltzmann distribution with respect to  $Q_{BC}$  gives:

$$\pi_{BC}(a | s) \approx \frac{\exp(Q_{BC}(s, a)/\alpha)}{\int \exp(Q_{BC}(s, a')/\alpha) da'}. \quad (5)$$

Using Eq. 3 and Eq. 5, we derive the equation for the Q-function  $Q_{BC}$  using the value function, log probability of the action, and entropy:

$$\hat{Q}_{BC} = V_{BC}(s) + \log \pi_{BC}(a | s) + \mathcal{H}[\pi_{BC}(\cdot | s)]. \quad (6)$$

Appendix A provides a complete derivation for Eq. 6.

Next, we detail how each term in Eq. 6 is computed or approximated. First, we must estimate the value function  $V_{BC}$ . We collect rollouts of the BC policy for a few initial online interaction episodes  $\{\tau_i\}_{i=1}^N$  and calculate the Monte Carlo returns under the online reward function:

$$\hat{V}_{BC}(s_t^{(i)}) = \frac{1}{N} \sum_{i=1}^N \sum_{k=t}^{\tau_i} \gamma^k r_{t+k}^{(i)}, \text{ for } t \in \{1, \dots, \tau_i\}. \quad (7)$$

These Monte Carlo returns are used to train a value estimator.

The remaining terms of Eq. 6 are derived from the BC policy  $\pi_{BC}$ . Our method is agnostic to the BC policy class, as long as it provides the log probabilities of actions and entropy. Nevertheless, we provide the remaining components of Eq. 6 using Gaussian policies as an illustrative case. The log probability of the actions  $\log \pi(a | s)$  for Gaussian policies has a closed-form expression:

$$\log \pi(a | s) = \frac{1}{2} \sum_{i=1}^d \left[ \frac{(a_i - \mu_i(s))^2}{\sigma_i^2(s)} + \log(2\pi\sigma_i^2(s)) \right], \quad (8)$$

where  $(\mu, \sigma)$  is Gaussian distribution output by the policy. Finally, the entropy  $\mathcal{H}[\pi(\cdot | s)]$  for Gaussian policies also has

a closed-form expression:

$$\mathcal{H}[\pi(\cdot | s)] = \frac{1}{2} \log(2\pi e \sigma^2(s)). \quad (9)$$

See Appendix A for a similar treatment for Gaussian Mixture Models (GMMs).

### B. Q-Gating for Online Reinforcement Learning

We now describe Q-Gating, our approach to using  $\hat{Q}_{BC}$  estimated from Eq. 6 to guide online reinforcement learning. We first initialize two separate Q-functions, both from  $\hat{Q}_{BC}$ :  $\hat{Q}_{BC}$  and  $Q_{RL}$ .  $\hat{Q}_{BC}$  is kept frozen to serve as a stable reference for the BC policy.  $Q_{RL}$  serves as the critic for the RL policy, and is updated through online interaction. During online training, we sample both the BC policy and RL policy for actions  $a_{BC}$  and  $a_{RL}$ , respectively. We evaluate  $a_{BC}$  using  $\hat{Q}_{BC}$  and  $a_{RL}$  using  $Q_{RL}$  to get their estimated Q values. The action corresponding to the greater of these two values is executed:

$$a = \begin{cases} a_{BC} \sim \pi_{BC}(s), & \hat{Q}_{BC}(s, a_{BC}) > Q_{RL}(s, a_{RL}) \\ a_{RL} \sim \pi_{RL}(s_t), & \text{otherwise.} \end{cases} \quad (10)$$

We propose this Q-Gating mechanism to preserve good BC action proposals while still enabling policy improvement through RL.  $\hat{Q}_{BC}$  represents the (estimated) value of taking an action in the current state, then following the BC policy thereafter. We evaluate BC actions through a frozen  $\hat{Q}_{BC}$ , so that BC actions with high value under  $\hat{Q}_{BC}$  are likely to be executed. On the other hand,  $Q_{RL}$  is only initialized as  $\hat{Q}_{BC}$ , and trains on samples with either BC or RL actions.  $Q_{RL}$  therefore improves its estimate of both BC and RL actions, enabling the reinforcement learning policy to surpass BC performance through online exploration and interaction. Note that our approach does not require access to the training

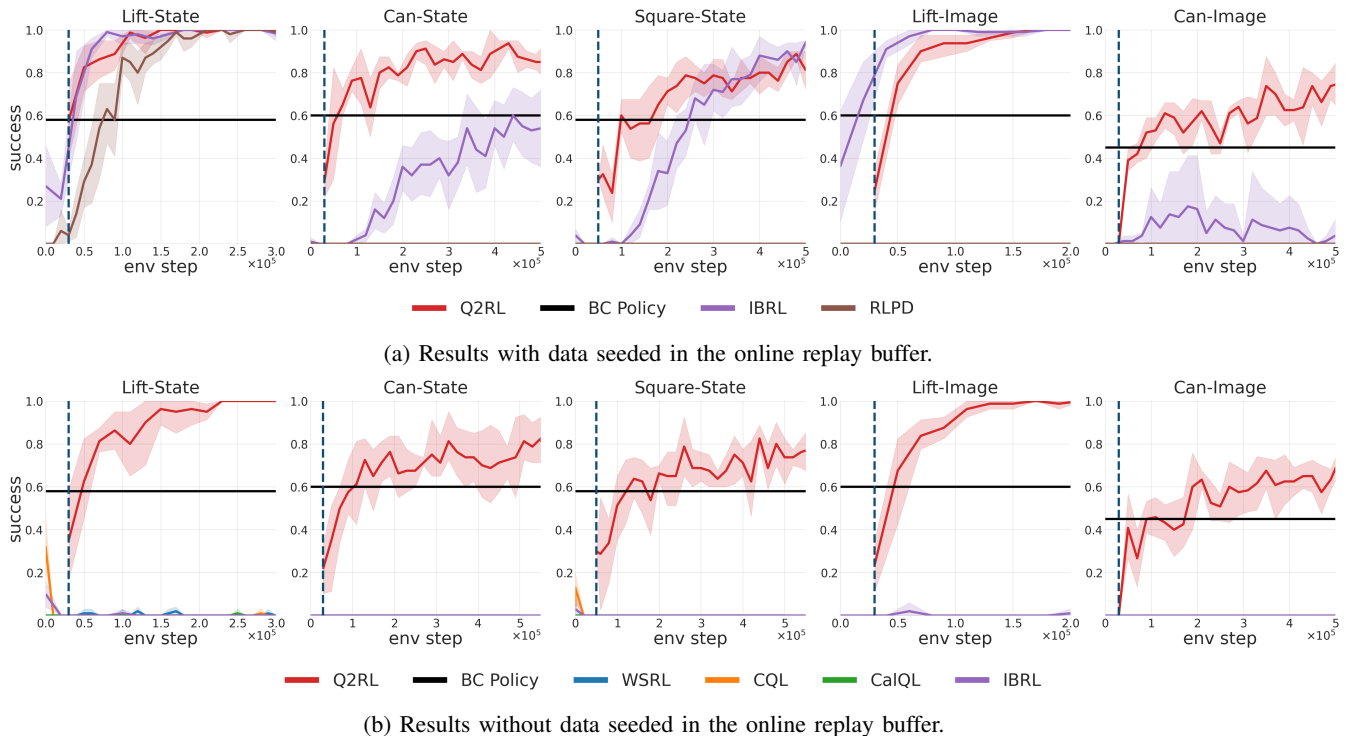


Fig. 4: Robomimic Results. The plot for each method begins at the start of online interaction. The dotted blue line signifies the end of the Q-Estimation phase and the start of online RL with Q-Gating for our method.

data for  $\pi_{BC}$ ; we rely on our Q-estimate  $\hat{Q}_{BC}$  to bootstrap online reinforcement learning.

### C. Implementation Details

$Q2RL$  can be applied to any behavior cloning policy that provides log probabilities of actions and entropy; in this work, we demonstrate the approach using Gaussian policies and GMMs. Similarly, though  $Q2RL$  is agnostic to the choice of off-policy RL algorithm, our implementation uses Soft Actor Critic [9].

To prevent the RL policy from deviating too far from the BC state-action distribution, we also weight the RL actor loss with an auxiliary BC loss [26]. While this regularization term can occasionally prevent the policy from reaching perfect success rates, our experiments in the following section show that training with BC regularization still outperforms baselines and plays a crucial role in producing safe, smooth behaviors on real robotic systems.

## IV. EXPERIMENTS

Our experiments are designed to answer two key questions: (i) does  $Q2RL$  improve performance beyond the BC policy faster and better than baseline approaches? (ii) Does  $Q2RL$  use high-value actions from the BC policy during online RL? To provide a comprehensive evaluation, we conduct experiments across multiple simulation environments using BC policies parameterized by Gaussian and GMM-RNN architectures, and compare against offline-to-online RL and BC-to-RL baselines.

Our evaluation spans both state-based and image-based observation settings. Finally, we validate the practicality of our approach through real-world RL experiments, assessing whether  $Q2RL$  can support online learning on physical hardware.

### A. Simulation Environments

We test  $Q2RL$  on tasks from the D4RL [7] and robomimic [23] benchmarks. All tasks are visualized in Fig 2.

1) **D4RL**: We run experiments on two Adroit manipulation tasks: *Pen* and *Door*. The Pen task involves dexterous in-hand manipulation, requiring the robotic hand to reorient a pen to a target configuration, while Door consists of a standard door-opening task. In addition, we evaluate performance on the *Kitchen* environment, which features a Franka robot interacting with multiple objects to achieve a multi-stage goal configuration. For these experiments, we use the offline Kitchen-Complete dataset, which contains 19 successful demonstrations which are representative of standard robot learning setups where we only have access to a few successful trajectories. We conduct all D4RL experiments without access to offline training data during the online RL phase.

2) **robomimic**: We run experiments on three widely used robomimic tasks, *Lift*, *Can*, and *Square*, all featuring a Franka robotic arm. In the Lift task, the robot must lift a cube from a randomly sampled position on the table. The Can task requires the robot to grasp a can from one side of the table and place it into the correct bin among four target bins located on the opposite side. The Square task involves picking up a square nut from the table and inserting it onto a square peg. For

robomimic tasks, we test both with and without access to offline training data during the online RL phase.

### B. Baselines

We compare  $Q2RL$  with the following baselines:

1) **CQL**: Conservative Q-Learning [14] is an offline RL method that trains exclusively on an offline dataset by prioritizing in-distribution actions while penalizing out-of-distribution actions. This induces a pessimistic bias in the learned Q-function, which often requires a recalibration phase during online interaction.

2) **CalQL**: Calibrated Q-Learning [24] addresses the pessimism of Q-value estimates in CQL by calibrating the critic loss with respect to the value of a reference policy. However, despite this calibration, transitioning to online training can still lead to unlearning of its offline pretraining.

3) **WSRL**: Warm-Start RL [35] performs for offline-to-online RL without explicit data retention. To mitigate unlearning during the transition from offline to online training, WSRL introduces a warm-up phase in which the online replay buffer is seeded with rollouts generated by the offline pretrained policy. However, a key limitation of WSRL is its dependence on the quality of the offline pretraining data.

4) **RLPD**: RL with Prior Demonstrations [2] incorporates offline data directly into online learning by uniformly sampling from both offline and online replay buffers during actor and critic updates.

5) **IBRL**: Imitation Bootstrapped RL [10] uses a pre-trained BC policy to bootstrap Q value estimates and action proposals for RL. IBRL randomly initializes the Q function and uses demonstration data to seed the online replay buffer.

### C. Results on D4RL

We present the results of D4RL tasks in Fig 3. For the **Kitchen** task, since the offline Kitchen-Complete dataset only contains successful demonstrations, representative of the current robot learning paradigm, the offline RL methods do not have enough data for a good pretraining. As offline RL pretraining fails, offline to online RL methods take longer time to achieve comparable performance. Whereas  $Q2RL$  is able to start with some initial performance provided by the BC pretraining, recover it further using our gating mechanism and continuously improve the performance with more environment interactions. In the **Pen** task, the offline data is able to provide sufficient pretraining for both BC and offline RL methods. We observe that  $Q2RL$  starts with near-optimal performance and maintains it throughout online learning. WSRL and CalQL are also able to reach optimal performance with online interaction, whereas CQL experiences a drop in performance, which we attribute to its conservative Q value estimates. IBRL is also eventually able to reach optimal performance. For the **Door** task, we observe that WSRL is able to converge to a more optimal policy than the baselines while  $Q2RL$  is not able to converge to optimal performance. From the rollout videos we observe that this is due to a difference in the behavior of an optimal RL policy and BC policy for the Door task. Though

the  $Q4RL$  policy does not converge to optimal performance, from the rollout videos we can see that they are more true to the human demonstrations and therefore more realistic to be executed in the real world whereas the RL policy would be infeasible to be carried out in the real world.

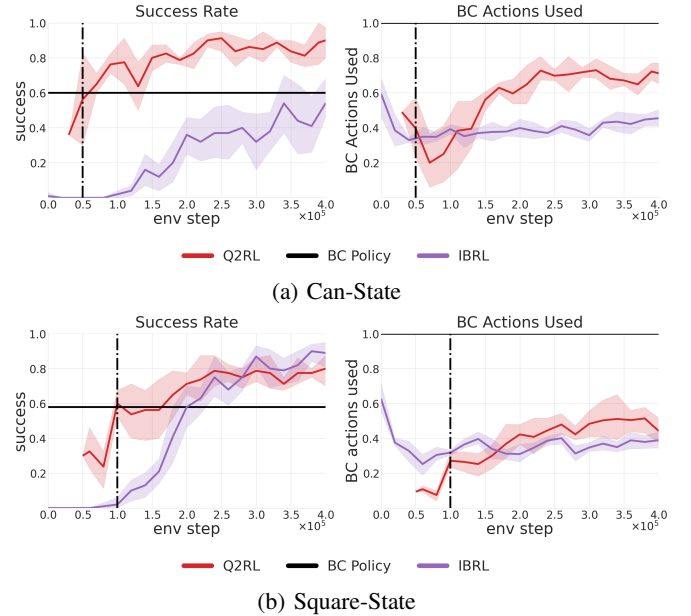
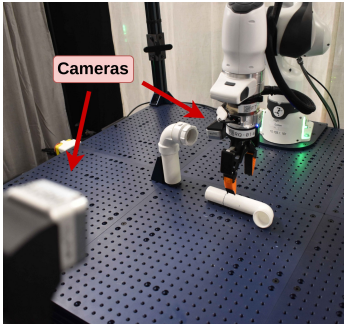


Fig. 5: BC vs. RL Actions Used During Online Training.  $Q2RL$  optimizes BC vs. RL action selection via Q-Gating, and achieves higher success rate earlier than IBRL.

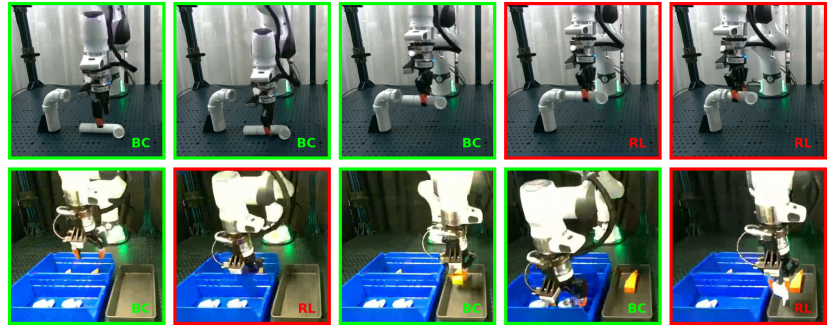
### D. Robomimic Results

Next, we evaluate  $Q2RL$  on both state-based and image-based robomimic tasks. We use two cameras for image-based tasks, the agentview of the workspace and wrist camera. We implement  $Q2RL$  and all baselines with a learnable CNN encoder for the images to keep the implementation standard and comparisons fair. As these tasks are more challenging than the D4RL benchmarks, we consider two experimental settings: one in which training data is available in the online replay buffer, and another in which the policy does not have access to any offline training data during online learning.

1) **Access to BC Training Data**: Results for the robomimic tasks with access to training data during online learning are shown in Fig. 4a. When training data is available in the online replay buffer, IBRL performs competitively and is able to converge to optimal performance. While IBRL is simple and effective, it does not explicitly preserve the initial performance of the BC policy, causing online interaction to begin with limited or no meaningful performance. In contrast,  $Q2RL$  begins online interaction with non-trivial initial performance rather than starting from near-zero success, which is particularly important for real-world deployment where unsafe or unproductive exploration is costly. Moreover, by restricting learning to targeted policy improvements,  $Q2RL$  converges to high-performing policies more rapidly. We also observe that for harder high-dimensional task, Can-Image, IBRL is unable



(a) On-Robot RL Setup.



(b) Trained  $Q2RL$  policies for pipe assembly and kitting tasks. Green: BC actions, red: RL actions.

Fig. 6: Real World Experiments. (a) We use a Franka Panda arm with a Robotiq 2F-85 gripper, and RGB images from both wrist and workspace Intel RealSense D405s. (b) We analyze usage of BC vs. RL actions in representative trajectories executed by trained  $Q2RL$  policies. For pipe assembly,  $Q2RL$  relies on BC for grasping and initial alignment, then switches to RL for high-precision, contact-rich insertion. For kitting, the BC policy was trained in a task setting with only one object present in each bin;  $Q2RL$  learns to complete the task with multiple objects in each bin. In the modified setting,  $Q2RL$  uses BC actions for behaviors common to both task settings (e.g. moving between bins), and uses RL actions for part grasping and placement.

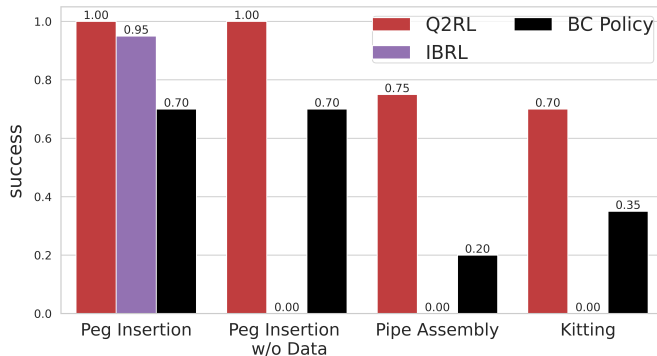


Fig. 7: On-Robot RL Results. Task success is reported over 20 trials per method.  $Q2RL$  significantly improves BC policy performance on every task. IBRL fails when (1) demonstrations are not seeded in the replay buffer (w/o Data), and (2) when tasks involve large action spaces and long horizon manipulation (even with seeding).

to reach the base policy performance. Though, we note that in the IBRL paper, it is able to converge for Can-Image task, we hypothesize that this is due to the other optimization that IBRL uses for their RL policy like specialized ViT encoders.  $Q2RL$  on the other hand is able to outperform the base policy with only the standard set of encoders. RLPD succeeds only on simpler tasks such as Lift and struggles to converge on the more challenging Can and Square tasks.

2) *No Access to BC Training Data*: Results for the robomimic tasks without access to training data during online learning are shown in Fig. 4b. Offline RL and offline to online RL methods are unable to learn for harder robotic tasks with limited offline datasets. Without access to data, IBRL is also unable to learn any meaningful performance. Early in training, the randomly initialized RL policy can propose arbitrary actions that may receive spuriously high Q-values from an untrained critic. To mitigate this issue, IBRL relies

on seeding the replay buffer with training data so that the critic learns to assign higher value to BC-generated actions than to random RL actions. In the absence of a seeded replay buffer, IBRL has no mechanism to calibrate Q-values, leading to unreliable action selection. Our method is able to continually improve the performance even in the absence of training data. We also conduct a thorough evaluation for seeding the replay buffer with different amounts of data in Appendix B.

#### E. Analysis of BC Actions Used During Online Training

Next, we investigate the relationship between success rate and the ratio of BC vs. RL actions used during online training. Fig. 5 shows the success rate and the fraction of BC actions selected by  $Q2RL$  vs. IBRL on Can-State and Square-State tasks. The vertical dotted line serves as a visual reference indicating when online performance approximately matches the original BC policy performance.

We observe that  $Q2RL$  rapidly recovers the original BC performance within only a few online interaction steps, demonstrating that it identified and exploited high-value actions of the BC policy. This performance is achieved while using BC actions less than half of the time, indicating that the RL policy also learned to produce effective actions. With continued training, the ratio of BC actions for our method increases, indicating that our Q-Gating algorithm enables higher success rates as  $Q_{RL}$  improves.

In contrast, IBRL selects a similar fraction of BC actions, but requires substantially more interaction to reach the same performance. Further, the ratio of BC actions used by IBRL during training stays relatively flat, indicating that it trades off between BC and RL actions less efficiently than  $Q2RL$ .

#### F. Real World RL Experiments

1) *Experiment Setup*: We run all experiments on a 7-DOF Franka Panda Robot with a Robotiq 2F-85 gripper. For all tasks, the robot receives the end effector pose, gripper width,

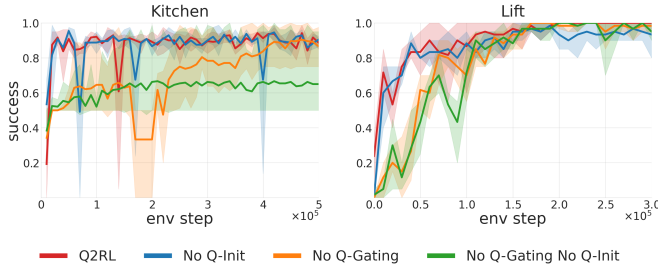


Fig. 8: Q-Initialization and Q-Gating Ablations.

and  $84 \times 84$  RGB images from a workspace camera and wrist camera (see Fig. 6a). The robot outputs delta pose actions executed by a Cartesian impedance controller. We evaluate on the following tasks (see Fig. 1):

**Peg Insertion.** This insertion task from FMB [20] has been used in recent on-robot RL methods [19, 35]. Success occurs when the peg is fully inserted into the board.

**Plumbing Pipe Assembly.** This task involves both grasping and low tolerance insertion. Insertion requires rotating the gripper to align the pipe with the fixture. Success occurs when the pipe is fully inserted into the fixture.

**Kitting with Task Distribution Shift.** In this task, two parts must be retrieved from their respective bins and placed in a kitting tray. As the longest horizon task, this task requires multiple picks and places. Success occurs when the robot has placed both parts in the kitting tray. The BC policy for this task was only trained on demonstrations with a single object per bin. To evaluate robustness to task distribution shift, the test task setting has two objects per bin, placed in new locations.

GMM-RNN BC policies were trained for each task on 50 to 100 demos collected by a human operator using a 3D Connexion Spacemouse. We compare  $Q2RL$  with IBRL for offline-to-online improvement with the pre-trained BC policies. During online learning, a termination signal and sparse reward signal were provided by a human operator when the task’s success conditions were met. Environment resets were also handled by a human operator. Additional details on the real experiment setup can be found in Appendix D.

2) *Real World Results:* For the Kitting task, all results are for the modified setting (two parts in each bin instead of one). The best checkpoint after up to 2.5 hours of online training is used for each method, representing no more than 170k gradient steps and 80k environment steps (much fewer than in simulation). From Fig. 7,  **$Q2RL$  significantly improves performance compared to the original BC policy, achieving a 1.4x improvement and 100% success on the Peg Insertion task, 3.75x improvement on Pipe Assembly, and 2x improvement on Modified Kitting.** While IBRL performed on par with  $Q2RL$  for the Peg Insertion task, it achieved zero success when demonstration samples were not seeded in the online replay buffer. Even with data seeding, when tasks involved large action spaces or were long horizon, IBRL had zero success, failing to recover BC performance. We attribute IBRL’s poor performance to its use of a single, randomly initialized Q function to select BC vs. RL actions,

which hinders it from accurately rating BC actions, especially for long horizon, sparse reward tasks. In contrast, we initialize  $Q_{RL}$  from  $Q_{BC}$ , and our proposed Q-Gating procedure uses both a frozen  $Q_{BC}$  and trainable  $Q_{RL}$  to score respective actions more accurately.

**Safety.** During the course of training IBRL for Peg Insertion, we noted two safety violations, in which the Franka exerted excessive force against the board and caused the robot to fault. These safety violations occurred despite careful tuning of controller gains during preliminary testing that was subsequently kept fixed for all reported runs.  $Q2RL$  experienced no such safety violations, and we qualitatively observed that it produced relatively smoother and safer actions early in training: we attribute this behavior to a combination of the auxiliary BC loss, Q-Gating, and  $Q_{RL}$  initialization with  $Q_{BC}$ .

**Qualitative Results.** Fig. 6b shows frames from representative rollouts of trained  $Q2RL$  policies, annotated to indicate BC vs. RL action selection. We observe that  $Q2RL$  selects BC actions to perform behaviors supported by the BC policy’s training data, e.g. initial pipe grasping and alignment for the Pipe Assembly task, before switching to RL actions for high-precision, contact-rich insertion. With Modified Kitting, BC actions were selected for moving between bins, with RL actions handling parts of the task that had shifted from their original setting, e.g. grasping parts from new locations. Other cases of switching between BC and RL actions corresponded to recovery behaviors, such as insertion re-alignment and re-grasping. Not all rollouts or trained  $Q2RL$  policies exhibit these behaviors exactly as described here, but we generally find meaningful switching between BC and RL actions. Please refer to the supplementary material for videos of annotated rollouts, and Appendix D for additional details on results.

## G. Ablations

We ablate Q-Initialization and Q-Gating in Fig. 8, and show that both are critical for achieving high performance. When Q-Gating is enabled, initializing  $Q_{RL}$  with  $\hat{Q}_{BC}$  yields performance comparable to a randomly initialized critic; however, we retain this initialization because it provides a safer starting distribution that remains closer to the BC policy, which is preferable for on-robot RL.

## V. CONCLUSION

We presented  $Q2RL$ , an algorithm that improves the performance of a behavior cloning policy through online reinforcement learning.  $Q2RL$  estimates a Q function from the BC policy, then uses it to initialize and guide online RL through Q-gating. Our approach does not require access to the BC policy’s original training data, and is compatible with any BC policy class that provides action likelihoods and entropy. In on-robot reinforcement learning experiments,  $Q2RL$  successfully improves performance on high precision, contact-rich manipulation tasks within a few hours of online interaction.

## REFERENCES

- [1] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In *International conference on machine learning*, pages 104–114. PMLR, 2020.
- [2] Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning with offline data. In *International Conference on Machine Learning*, pages 1577–1594. PMLR, 2023.
- [3] Ondrej Biza, Thomas Weng, Lingfeng Sun, Karl Schmeckpeper, Tarik Kelestemur, Yecheng Jason Ma, Robert Platt, Jan-Willem van de Meent, and Lawson L.S. Wong. On-robot reinforcement learning with goal-contrastive rewards. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4797–4805, 2025. doi: 10.1109/ICRA55743.2025.11128466.
- [4] Lawrence Yunliang Chen, Simeon Adebola, and Ken Goldberg. Berkeley UR5 demonstration dataset. <https://sites.google.com/view/berkeley-ur5/home>.
- [5] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025.
- [6] Lakshita Dodeja, Karl Schmeckpeper, Shivam Vats, Thomas Weng, Mingxi Jia, George Konidaris, and Stefanie Tellex. Accelerating residual reinforcement learning with uncertainty estimation. *arXiv preprint arXiv:2506.17564*, 2025.
- [7] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- [8] Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch deep reinforcement learning algorithms. *arXiv preprint arXiv:1910.01708*, 2019.
- [9] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. Pmlr, 2018.
- [10] Hengyuan Hu, Suvir Mirchandani, and Dorsa Sadigh. Imitation bootstrapped reinforcement learning, 2023.
- [11] Sunshine Jiang, Xiaolin Fang, Nicholas Roy, Tomás Lozano-Pérez, Leslie Pack Kaelbling, and Siddharth Ancha. Streaming flow policy: Simplifying diffusion / flow-matching policies by treating action trajectories as flow trajectories. *arXiv preprint arXiv:2505.21851*, 2025.
- [12] Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo, Avinash Kumar, Matthias Loskyll, Juan Aparicio Ojea, Eugen Solowjow, and Sergey Levine. Residual reinforcement learning for robot control. In *2019 international conference on robotics and automation (ICRA)*, pages 6023–6029. IEEE, 2019.
- [13] Michael Kelly, Chelsea Sidrane, Katherine Rose Driggs-Campbell, and Mykel J. Kochenderfer. Hg-dagger: Interactive imitation learning with human experts. In *International Conference on Robotics and Automation, ICRA 2019, Montreal, QC, Canada, May 20-24, 2019*, pages 8077–8083. IEEE, 2019.
- [14] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *Advances in neural information processing systems*, 33:1179–1191, 2020.
- [15] Cassidy Laidlaw and Anca Dragan. The boltzmann policy distribution: Accounting for systematic suboptimality in human models. In *ICLR*, 2022.
- [16] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- [17] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- [18] R Duncan Luce et al. *Individual choice behavior*, volume 4. Wiley New York, 1959.
- [19] Jianlan Luo, Zheyuan Hu, Charles Xu, You Liang Tan, Jacob Berg, Archit Sharma, Stefan Schaal, Chelsea Finn, Abhishek Gupta, and Sergey Levine. Serl: A software suite for sample-efficient robotic reinforcement learning. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 16961–16969. IEEE, 2024.
- [20] Jianlan Luo, Charles Xu, Fangchen Liu, Liam Tan, Zipeng Lin, Jeffrey Wu, Pieter Abbeel, and Sergey Levine. Fmb: a functional manipulation benchmark for generalizable robotic learning. *The International Journal of Robotics Research*, 44(4):592–606, 2025.
- [21] Jianlan Luo, Charles Xu, Jeffrey Wu, and Sergey Levine. Precise and dexterous robotic manipulation via human-in-the-loop reinforcement learning. *Science Robotics*, 10(105):eads5033, 2025.
- [22] Ajay Mandlekar, Danfei Xu, Roberto Martín-Martín, Silvio Savarese, and Li Fei-Fei. Learning to generalize across long-horizon tasks from human demonstrations. *arXiv preprint arXiv:2003.06085*, 2020.
- [23] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.
- [24] Mitsuhiko Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. *Advances in Neural Information Processing Systems*, 36:62244–62269, 2023.
- [25] Soroush Nasiriany, Tian Gao, Ajay Mandlekar, and Yuke Zhu. Learning and retrieval from prior data for skill-based imitation learning. *arXiv preprint arXiv:2210.11435*, 2022.

- [26] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2018.
- [27] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 661–668. JMLR Workshop and Conference Proceedings, 2010.
- [28] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011.
- [29] Ziyad Sheebaelhamd, Michael Tschannen, Michael Muehlebach, and Claire Vernade. Quantization-free autoregressive action transformer. *arXiv preprint arXiv:2503.14259*, 2025.
- [30] Tom Silver, Kelsey Allen, Josh Tenenbaum, and Leslie Kaelbling. Residual policy learning. *arXiv preprint arXiv:1812.06298*, 2018.
- [31] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- [32] Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *International conference on learning representations*, 2021.
- [33] Xiu Yuan, Tongzhou Mu, Stone Tao, Yunhao Fang, Mengke Zhang, and Hao Su. Policy decorator: Model-agnostic online refinement for large policy model. *arXiv preprint arXiv:2412.13630*, 2024.
- [34] Tianhao Zhang, Zoe McCarthy, Owen Jow, Dennis Lee, Xi Chen, Ken Goldberg, and Pieter Abbeel. Deep imitation learning for complex manipulation tasks from virtual reality teleoperation. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 5628–5635. Ieee, 2018.
- [35] Zhiyuan Zhou, Andy Peng, Qiyang Li, Sergey Levine, and Aviral Kumar. Efficient online reinforcement learning fine-tuning need not retain offline data. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=HN0CYZbAPw>.