# Bidirectional Decoding: Improving Action Chunking via Closed-Loop Resampling

**Yuejiang Liu,**[*] **Jubayer Ibn Hamid,**[*] **Annie Xie, Yoonho Lee, Maximilian Du, Chelsea Finn**
Department of Computer Science, Stanford University

## Abstract

Predicting and executing a sequence of actions without intermediate replanning, known as action chunking, is increasingly used in robot learning from human demonstrations. Yet, its reported effects on the learned policy are inconsistent: some studies find it crucial for achieving strong results, while others observe decreased performance. In this paper, we first dissect how action chunking impacts the divergence between a learner and a demonstrator. We find that action chunking allows the learner to better capture the temporal dependencies in demonstrations but at the cost of reduced reactivity in stochastic environments. To address this tradeoff, we propose *Bidirectional Decoding* (BID), a test-time inference algorithm that bridges action chunking with closed-loop operations. BID samples multiple predictions at each time step and searches for the optimal one based on two criteria: (i) backward coherence, which favors samples that align with previous decisions; (ii) forward contrast, which seeks samples of high likelihood for future plans. By coupling decisions within and across action chunks, BID promotes consistency over time while maintaining reactivity to unexpected changes. Experimental results show that BID boosts the performance of two state-of-the-art generative policies across seven simulation benchmarks and two real-world tasks. Code and videos are available at `https://bid-robot.github.io`.

## 1 Introduction

The increasing availability of human demonstrations has spurred renewed interest in behavioral cloning (Atkeson & Schaal, 1997; Argall et al., 2009). In particular, recent studies have highlighted the potential of learning from large-scale demonstrations to acquire a variety of complex skills (Zhao et al., 2023; Chi et al., 2023; Fu et al., 2024b; Lee et al., 2024; Khazatsky et al., 2024). Yet, existing methods still struggle with two common properties of human demonstrations: (i) strong temporal dependencies across multiple steps, such as idle pauses (Chi et al., 2023) and latent strategies (Xie et al., 2021; Ma et al., 2024), (ii) large style variability across different demonstrations, such as differences in proficiency (Belkhale et al., 2024) and preference (Kuefler & Kochenderfer, 2017). Often, both properties are prevalent yet unlabeled in collected data, posing significant challenges to the traditional behavioral cloning that maps an input state to an action.

In response to these challenges, recent works have pursued a generative approach equipped with action chunking: (i) predicting a sequence of actions over multiple time steps and executing all or part of the sequence (Zhao et al., 2023; Chi et al., 2023); (ii) modeling the distribution of action chunks and sampling from the learned model in an independent (Chi et al., 2023; Prasad et al., 2024) or weakly dependent (Janner et al., 2022; Zhao et al., 2023) manner for sequential decisions. Some studies find this approach crucial for learning a performant policy in laboratory scenarios (Zhao et al., 2023; Chi et al., 2023), while other recent work reports opposite outcomes under practical conditions (Lee et al., 2024). The reasons behind these conflicting observations remain unclear.

In this paper, we first dissect the influence of action chunking by examining the divergence between learned policies and human demonstrations. We find that, when a policy is built with limited context length – little or no history is used as input for robustness or efficiency (Mandlekar et al., 2020; Bharadhwaj et al., 2023b; Brohan et al., 2023a;b; Shi et al., 2023; Collaboration et al., 2023) – increasing the length of action chunks allows for implicit conditioning on more past actions, thereby improving its ability to capture the temporal dependencies inherent in demonstrations. However, this advantage comes at the cost of reduced access to recent state observations, which can be crucial for
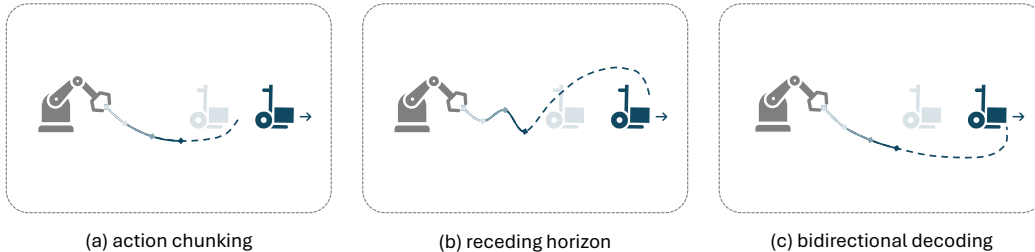
---

[*]Equal contribution.

Figure 1: Illustration of action chunking in a stochastic environment, where a robot is challenged to catch a moving trolley. (a) Vanilla action chunking (Zhao et al., 2023) executes actions based on previous predictions, resulting in delayed reactions to the latest box location. (b) Receding horizon (Chi et al., 2023) allows for faster reactions, but leads to a jittery trajectory in the presence of multimodal demonstrations (*e.g.*, both left- and right-handers). (c) Our Bidirectional Decoding explicitly searches for the optimal action from multiple predictions sampled at each time step, achieving both long-range consistency and short-term reactivity.

reacting to unexpected changes in stochastic environments, such as action noise or object motions. This tradeoff raises a crucial question: How can we preserve the strengths of action chunking in long-term consistency without suffering from its limitations in short-term reactivity?

To this end, we introduce *Bidirectional Decoding* (BID), an inference algorithm that bridges action chunking with closed-loop operations. Our main idea is to sample multiple predictions at each time step and search for the most desirable one. Specifically, BID operates on two decoding criteria: (i) backward coherence, which favors samples that are close to the sequence selected at the previous step; (ii) forward contrast, which favors samples that are close to the output of a stronger policy and distant from those of a weaker one. As illustrated in Fig. 1, BID updates the chunk of future actions based on the previous strategy, promoting temporal consistency over extended periods while remaining reactive to unexpected changes.

The main contributions of this paper are twofold: (i) a thorough analysis of action chunking (§3), and (ii) a decoding algorithm to improve it (§4). Empirically, we validate our theoretical analysis through a one-dimensional diagnostic simulation and evaluate our decoding method on two state-of-the-art generative policies across seven simulations and two real-world tasks (§5). Our experiment results show that the proposed BID boosts the performance of recent policies by more than 32% in relative performance. BID is model-agnostic, computationally efficient, and easy to implement, serving as a plug-and-play component to enhance generative behavior cloning at test time.

## 2 RELATED WORK

**Behavioral Cloning.** Learning from human demonstrations is becoming increasingly popular in robot learning due to recent advances in robotic teleoperation interfaces (Sivakumar et al., 2022; Zhao et al., 2023; Wu et al., 2023; Chi et al., 2024). Generative Behavior cloning, which models the distribution of demonstrations, is particularly appealing due to its algorithmic simplicity and empirical efficacy (Jang et al., 2022; Florence et al., 2022; Brohan et al., 2022; Shafiullah et al., 2022; Zhao et al., 2023; Chi et al., 2024; Brohan et al., 2023a). However, a significant limitation is compounding errors, where deviations from the training distribution accumulate over time (Ross et al., 2011; Ke et al., 2021). These errors can be mitigated by gathering expert correction data (Ross et al., 2011; Kelly et al., 2019; Menda et al., 2019; Hoque et al., 2021a;b) or injecting noise during data collection (Laskey et al., 2017; Brandfonbrener et al., 2023), but such strategies require additional time and effort from human operators. To address this, recent work proposes predicting a sequence of multiple actions into the future, known as action chunking, which reduces the effective control horizon (Lai et al., 2022; Zhao et al., 2023; George & Farimani, 2023; Bharadhwaj et al., 2023a). By handling sequences of actions, action chunking is also better at handling temporal dependencies in the data, such as idle pauses (Swamy et al., 2022; Chi et al., 2023) or multiple styles (Li et al., 2017; Kuefler & Kochenderfer, 2017; Gandhi et al., 2023; Belkhale et al., 2024). However, independently drawn action sequence samples may not preserve the necessary temporal dependencies for smooth and consistent execution. Our work provides a thorough analysis of action chunking and proposes a decoding algorithm to improve it.

**Sequential Decoding.** Decoding algorithms have been studied in generative sequence modeling for decades, with renewed attention driven by recent advances in large language modeling (LLM). One prominent approach focuses on leveraging internal metrics, *e.g.*, likelihood scores, to improve

the quality of generated sequences. Notable examples include beam search (Freitag & Al-Onaizan, 2017; Vijayakumar et al., 2018), truncated sampling (Fan et al., 2018; Hewitt et al., 2022), minimum Bayes risk decoding (Kumar & Byrne, 2004; Müller & Sennrich, 2021), and others (Welleck et al., 2019; Meister et al., 2023; Fu et al., 2024a). Another line of research explores the distinctions between multiple models to jointly optimize for the desired properties such as quality or efficiency (Li et al., 2023; Leviathan et al., 2023). More recently, several studies have highlighted the potential of guiding the decoding or sampling process through the use of an external model, such as a classifier (Dhariwal & Nichol, 2021) or reward model (Khanov et al., 2023). In the context of robot learning, recent works have explored guided decoding for long-horizon robotic planning (Huang et al., 2023) and manipulator geometry designs (Xu et al., 2024). Nevertheless, effective decoding strategies for low-level robotic actions remain lacking. Concurrent to our work, Nakamoto et al. (2024) propose to select action samples by querying a value function learned from reward-annotated demonstrations (Hansen-Estruch et al., 2023). Our work does not rely on a separate value function; instead, we propose a decoding strategy that addresses the consistency-reactivity tradeoff inherent in action chunking through sample comparison.

## 3 ANALYSIS: TRADEOFFS IN ACTION CHUNKING

### 3.1 PRELIMINARIES

Consider a dataset of demonstrations $\mathcal{D} = \{\tau_i\}_{i=1}^N$, where each demonstration $\tau_i$ consists of a sequence of state-action pairs $\tau_i = \{(s_1, a_1), (s_2, a_2), \cdots, (s_T, a_T)\}$ provided by a human expert. At each time step $t$, the demonstrated action $a_t$ is influenced not only by the observed state $s_t$, but also by latent variables $z_t$, such as planning strategies (*e.g.*, subgoals) and personal preferences (*e.g.*, handedness). These latent variables can persist across multiple time steps and vary significantly between different demonstrations. Fig. 2 illustrates the decision process of a human expert, highlighting the inherent temporal dependencies.

To model these temporal dependencies, some recent work (Zhao et al., 2023; Chi et al., 2023; Lee et al., 2024) utilize action chunking, *i.e.*, learning the joint distribution of future actions conditioned on past states $\pi(a_t, a_{t+1}, \cdots, a_{t+l}|s_{t-c}, \cdots, s_t)$, or in short $\pi(a_{t:t+l}|s_{t-c:t})$. Here, $c$ denotes the number of past steps for state inputs, and $l$ denotes the number of future steps for action outputs. Training such a policy typically involves minimizing the divergence of action distributions between the model $\pi$ and the expert $\pi^*$,

$$\pi = \arg\min_{\pi} \sum_{\tau \in \mathcal{D}} \sum_{\substack{s_{t-c:t} \\ a_{t:t+l}}} \mathcal{L}(\pi(a_{t:t+l}|s_{t-c:t}), \pi^*(a_{t:t+l}|s_{t-c:t})). \tag{1}$$

During deployment, the policy is operated by sampling a sequence of actions and executing a part of or the entire sequence for $h \in [1, l]$ time steps without re-planning. This approach essentially takes in $c$ states as context and executes $h$ actions, which we refer to as a $(c, h)$-policy.

The choices of context length $c$ and action horizon $h$ often play a crucial role in the effectiveness of the learned policy. Recent policies often use a short context length $c$, as extending the context can lead to performance degradation in the presence of limited training data (refer to Appendix A.2 for more details). Conversely, extending the action horizon $h$ has produced mixed results. Some studies report benefits in laboratory settings (Zhao et al., 2023; Chi et al., 2023), while others find it detrimental in real-world scenarios (Lee et al., 2024). The reasons behind these conflicting outcomes are not well understood yet.

Notably, Zhao et al. (2023) hypothesize that action chunking mitigates compounding errors, but it is unclear why this would hold when action chunking cannot correct deviations during open-loop execution. Similarly, Chi et al. (2023) argue that action chunking facilitates consistency and planning, yet the choice of action horizon is rather short and remains an empirical heuristic. The lack of a theoretical understanding of action chunking limits our ability to use it effectively across different methods and tasks. In the following sections, we seek to explicitly identify the strengths and weaknesses of action chunking, aiming to provide practical guidance on its suitable settings as well as to inform the design of better algorithms.

### 3.2 ANALYSIS

To understand the influence of action chunking, we focus on the last time step of an action chunk, where the discrepancy between the expert policy and the learned policy is most pronounced. At this
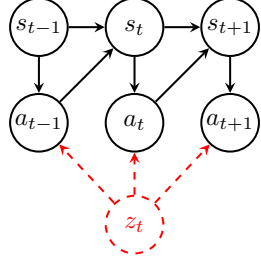
Figure 2: Illustration of the expert decision process, where a latent variable introduces temporal dependencies in actions.
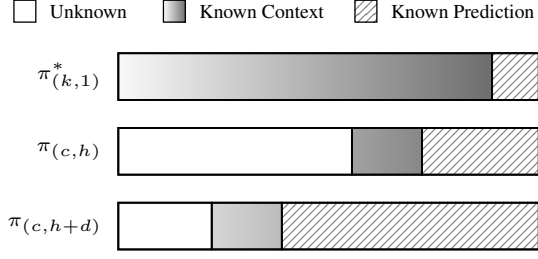
Figure 3: Illustration of $(k, 1)$-expert, $(c, h)$-learner and $(c, h + d)$-learner. Gray shades are observed *contexts*; darker indicates higher importance. Hatched areas denote executed *predictions*.

time step $t$, a $(k, 1)$-expert written as $\pi^* := \pi^*(a_t | s_{t-k:t}, z_{t-k:t})$ predicts $a_t$ by conditioning on $k$ steps of the past states and the corresponding latent variables. In contrast, a $(c, h)$-learner written as $\pi_{(c,h)} := \pi_{(c,h)}(a_t \mid s_{t-h-c:t-h}, a_{t-h:t-1})$ is constrained to observe $c$ steps of the past states and $h - 1$ steps of the past actions within the predicted chunk.

Considering that recent policies often use a short context length $c$, we assume the range of temporal dependency modeled by a $(c, h)$-policy is limited:

**Assumption 1.** The sum of context length and action horizon is less than the length of temporal dependency in expert demonstrations, $c + h < k$.

Additionally, since a $(c, h)$-policy observes only a subset of the states that the expert is conditioned on, we assume that an optimal policy must reconstruct all missing information correctly:

**Assumption 2.** An optimal $\pi_{c,h}$ must infer the unobserved states based on the observed states and actions by modeling the transition dynamics $P(s_{t'} \mid s_{t'-1}, a_{t'-1})$ accurately for all time step $t'$.

Under these assumptions, the divergence between a learned policy and an expert policy is attributed to two factors: (i) the importance of unobserved states in predicting the current action, and (ii) the difficulty of inferring unobserved states based on the available information.

To more clearly see the influence of action horizon on these factors, we next compare the performance of two policies that have the same context lengths but different action horizons, $\pi_h := \pi_{(c,h)}(a_t | s_{t-h-c:t-h}, a_{t-h:t-1})$ and $\pi_{h+d} := \pi_{(c,h+d)}(a_t | s_{t-h-d-c:t-h-d}, a_{t-h-d:t-1})$, where $d > 0$ is the extended action horizon. As illustrated in Fig. 3, each policy has access to unique information that is unavailable to the other. $\pi_h$ observes some recent states, where $\pi_{h+d}$ is only aware of the executed actions. On the other hand, $\pi_{h+d}$ has access to some earlier states and actions, which precede all information available to $\pi_h$. We characterize the *importance* of observations as follows (formal definitions in Appendix C.1):

**Definition (Expected Observation Advantage).** If a policy can observe a state $s_t$, we say that it has an *observation advantage* $\alpha_t$ over another policy that cannot observe it.

**Definition (Maximum Inference Disadvantage).** If a policy cannot observe a state $s_t$, the maximum divergence arising from inferring it incorrectly is $\epsilon_t$.

Hence, we denote the observation advantage that $\pi_h$ gains from the observed recent states by $\alpha_f$ and the inference disadvantage it incurs from the earlier unobserved states by $\epsilon_b$, whereas $\pi_{h+d}$ conversely gains $\alpha_b$ but incurs $\epsilon_f$.

The *difficulty* of inferring each unobserved state hinges on both the relevant observations as well as the environmental stochasticity. We quantify this difficulty as follows:

**Definition (Forward Inference).** Let $P_f := P(S_t = g_t | S_{t-1} = g_{t-1}, a_{t-1})$ where $g_t$ and $g_{t-1}$ are the ground truth states in the deterministic environment at time $t$ and $t - 1$, respectively. In deterministic environments, $P_f = 1$, whereas in stochastic settings, $P_f$ is smaller.

**Definition (Backward Inference).** Let $P_b := P(S_t = g_t | S_{t+1} = g_{t+1})$ where $g_t$ and $g_{t+1}$ are the ground truth states in the deterministic environment at time $t$ and $t + 1$, respectively. Since $P_b$ is not conditioned on any action, it has higher entropy in general. In stochastic environments, $P_b$ is small.

Given that the forward inference is generally easier than the backward inference, the performances of $\pi_h$ and $\pi_{h+d}$ differ as follows (proofs are deferred to Appendix C):

**Proposition 1** (Consistency-Reactivity Inequalities). *Let $\mathcal{L}$ be a non-linear convex function measuring the prediction error with respect to demonstrations. Let $\mathcal{S}^+ \subset \{s_{t-k:t}\}$ be the states both the*

$(c, h)$ *and the* $(c, h+d)$ *policies observe and let* $\mathcal{S}^- := \{s_{t-k:t}\} \setminus \mathcal{S}^+$. *Let* $C := \{a_{t-h-d:t-1}\} \cup \mathcal{S}^+$, $G := \{a_t, z_{t-k:t}\} \cup \mathcal{S}^-$. *Then, we can bound the expected loss of the* $(c, h+d)$-*policy and the* $(c, h)$-*policy as:*

$$\alpha_f - \epsilon_b(1 - P_b^{2d}) \leq \min_{\pi_{h+d}} \mathbb{E}_G\left[\mathcal{L}(\pi_{h+d}, \pi^*)|C\right] - \min_{\pi_h} \mathbb{E}_G\left[\mathcal{L}(\pi_h, \pi^*)|C\right] \leq -\alpha_b + \epsilon_f(1 - P_f^{2d}).$$

(2)

***Remark 1.1.*** Eq. (2) provides a general comparison of the performance of the two policies. Intuitively, the advantage of each policy stems from the additional information it has access to (*i.e.* $\alpha_f$ for $\pi_h$ and $\alpha_b$ for $\pi_{h+d}$) while the disadvantage is bounded by the maximum divergence arising from inferring missing information incorrectly (*i.e.* $\epsilon_b(1 - P_b^{2d})$ for $\pi_h$ and $\epsilon_f(1 - P_f^{2d})$ for $\pi_{h+d}$).

We next examine two specific environmental settings: deterministic and stochastic. In highly deterministic environments, while both policies need to infer the same number of unobserved states, $\pi_{h+d}$ benefits from conditioning on additional actions, which may significantly aid in inferring the corresponding states through its action chunk. If the maximum errors $\epsilon_f$ arising from inferring these states are bounded, $\pi_{h+d}$ becomes strictly advantageous:

**Corollary 2** (Consistency in Deterministic Environments). *In a highly deterministic environment, if* $a_t$ *is influenced by at least one state in* $\{s_{t-h-c-d:t-h-c-1}\}$ *and* $\epsilon_f$ *is finite, then*

$$\min_{\pi_{h+d}} \mathbb{E}_G\left[\mathcal{L}(\pi_{h+d}, \pi^*)|C\right] < \min_{\pi_h} \mathbb{E}_G\left[\mathcal{L}(\pi_h, \pi^*)|C\right].$$

(3)

Conversely, in highly stochastic environments, inferring the unobserved states is challenging, regardless of whether the corresponding actions are known. On the other hand, the recent states are likely more important than earlier states for predicting the current action $a_t$. As a result, $\pi_{h+d}$ becomes strictly more disadvantageous:

**Corollary 3** (Reactivity in Stochastic Environments). *In a highly stochastic environment, if temporal dependency decreases over time, i.e.,* $\alpha_f > \epsilon_b$, *then*

$$\min_{\pi_{h+d}} \mathbb{E}_G\left[\mathcal{L}(\pi_{h+d}, \pi^*)|C\right] > \min_{\pi_h} \mathbb{E}_G\left[\mathcal{L}(\pi_h, \pi^*)|C\right].$$

(4)

In summary, depending on the experimental conditions, action chunking may have varying effects on the learned policy. On the one hand, it benefits the modeling of temporal dependencies in the demonstrations, due to extended access to actions executed at past time steps. On the other hand, it hinders reactions to unexpected states in stochastic environments, due to reduced access to state observations at recent time steps. As a result, there is no universally optimal choice of action horizon across all conditions. When both temporal dependencies and transition stochasticities are significant, tuning the action horizon entails an inherent trade-off between the two opposing factors. Moreover, stochasticitic changes (*e.g.*, wind gusts or object drops) are often unpredictable in advance, raising a critical need to ensure the reactivity of a policy built with action chunking.

## 4 METHOD: BIDIRECTIONAL DECODING

As analyzed above, action chunking improves long-term consistency but sacrifices short-term reactivity. In this section, we propose to address this tradeoff by bridging action chunking with closed-loop operations. We will first outline the general framework in §4.1 and then describe two specific criteria in §4.2.

### 4.1 TEST-TIME SEARCH

Recall that for a policy with prediction horizon $l$, an action chunk $a := \{a_t^{(t)}, a_{t+1}^{(t)}, \cdots, a_{t+l}^{(t)}\}$ sampled at time $t$ is expected to follow a consistent strategy over the subsequent $l$ time steps. However, executing the action chunk in an open-loop manner leaves the policy vulnerable to unexpected state changes. Alternatively, one can maximize reactivity by resampling an action chunk and executing only the first immediate action at every time step. Yet, this simple closed-loop approach destroys the consistency preserved within each chunk, potentially leading to oscillations between different strategies. How can we combine the benefits of both approaches?

Our key idea is to bridge action chunking and closed-loop resampling by making use of additional computes at test time. In particular, we seek to restore temporal consistency in closed-loop operations by scaling up the number of candidate samples. Intuitively, while the probability of any single
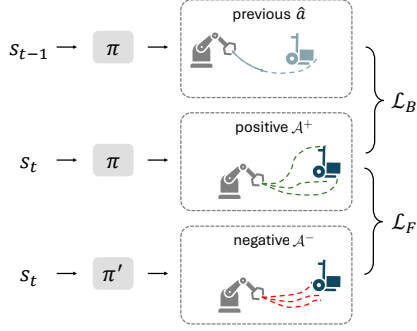
Figure 4: Illustration of bidirectional decoding.

**Algorithm 1** Bidirectional Decoding

**Require:** current state $s$, batch size $N$, mode size $K$, previous decision $\hat{a}$, strong policy $\pi$, weak policy $\pi'$
1: Generate $N$ samples from each policy $a \sim \pi(s)$, $a' \sim \pi'(s)$ to construct the initial sets $\mathcal{A}$ and $\mathcal{A}'$
2: Compute the backward loss $\mathcal{L}_B$ for each sample
3: Select K samples with minimal $\mathcal{L}_B$ from $\mathcal{A}$ and $\mathcal{A}'$ to construct $\mathcal{A}^+$ and $\mathcal{A}^-$, respectively
4: Compute the forward loss $\mathcal{L}_F$ for each sample
5: Select $a^* \in \mathcal{A}$ that minimizes the total loss
6: Update decision memory $\hat{a} \leftarrow a^*$

pair of samples sharing the same latent strategy is low, the likelihood of finding a consistent pair increases with the number of samples. We thus frame the problem of closed-loop action chunking as searching for the optimal action among a batch of samples drawn at each time step,

$$a^* = \arg\min_{a \in \mathcal{A}} \mathcal{L}_B(a) + \mathcal{L}_F(a), \tag{5}$$

where $\mathcal{A}$ is the set of sampled action chunks, $\mathcal{L}_B$ and $\mathcal{L}_F$ are two criteria approximating the optimality with respect to the backward decision and forward plan, which we will describe next.

## 4.2 BIDIRECTIONAL CRITERIA

**Backward coherence.** To preserve temporal dependencies in closed-loop operations, a sequence of actions should (i) commit to a consistent latent strategy over time, and (ii) react smoothly to unexpected changes. These desired properties motivate us to keep the action chunk selected at the previous time $\hat{a} := \{a_{t-1}^{(t-1)}, \cdots, a_{t+l-1}^{(t-1)}\}$ as a prior, and minimize the weighted Euclidean distance between the new action chunk and the prior across $l-1$ overlapping steps:

$$\mathcal{L}_B = \sum_{\tau=0}^{l-1} \rho^\tau \left\| a_{t+\tau}^{(t)} - a_{t+\tau}^{(t-1)} \right\|_2. \tag{6}$$

Here, $\rho$ is a decay hyperparameter to account for growing uncertainty over time. This backward objective encourages similar latent strategies between consecutive steps while allowing for gradual adaptation to unforeseen transition dynamics.

Nevertheless, the backward criterion alone presents a potential caveat: the prior chunk could be suboptimal due to the lack of information at the previous time step (*e.g.*, unexpected object motions). In such cases, selecting the next action chunk based solely on the prior may perpetuate suboptimality. Ideally, the sequential decision-making process should effectively correct suboptimal plans based on the latest observations. We next address this need through another forward criterion.

**Forward contrast.** Our design of the forward criterion is motivated by the need to identify the most optimal plan from a set of candidates. Within the same latent strategy, suboptimal samples may arise from (i) low likelihood under the learned model and (ii) divergence between the learned model and expert policy. To address this, we draw inspirations from LLM decoding techniques (Wang et al., 2022; Li et al., 2023) and introduce a forward contrast criterion. Specifically, we compare each candidate sample with two sets of reference samples: one positive set from a stronger policy and a negative set from a weaker one. The stronger policy is obtained from a well-trained checkpoint, whereas the weaker policy is taken from an early underfitting checkpoint and is expected to be further from the expert policy. Our forward objective is thus framed as minimizing the average distance between a candidate plan and the positive samples while maximizing its average distance from the negative ones,

$$\mathcal{L}_F = \frac{1}{N} \left( \sum_{a^+ \in \mathcal{A}^+} \sum_{\tau=0}^{l} \left\| a_{t+\tau}^{(t)} - a_{t+\tau}^+ \right\|_2 - \sum_{a^- \in \mathcal{A}^-} \sum_{\tau=0}^{l} \left\| a_{t+\tau}^{(t)} - a_{t+\tau}^- \right\|_2 \right), \tag{7}$$

where $\mathcal{A}^+ = \mathcal{A} \setminus \{a\}$ is the positive set predicted by the strong policy $\pi$, $\mathcal{A}^-$ is the negative set predicted by the weaker one $\pi'$, and $N$ is the sample size.
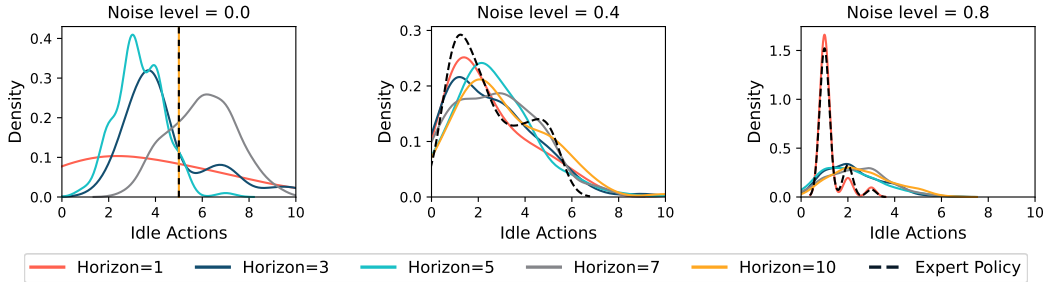
6

Figure 5: Probability distributions of idle actions taken by learners with varying action horizons in environments with varying stochasticity. The noise level in the environment grows from left to right.

Fig. 4 illustrates the combined effects of the backward coherence and forward contrast criteria on sample selection. Since samples in $\mathcal{A}^+$ and $\mathcal{A}^-$ are not all subject to the same strategy, we trim each set by removing samples that deviate significantly from the previous decision. This is achieved by summing over the $K$ smallest distance values for in the positive and negative sets in Eq. (7). The full process of our decoding method is outlined in Algorithm 1. Since all steps in BID can be computed in parallel, the overall computational overhead remains modest on modern GPU devices.

## 4.3 DISCUSSIONS

**Interpretation of our method.** Our method makes no changes to the learned policy; instead, it intervenes in the model distribution through sample selection. As illustrated in Fig. 10, randomly sampled sequences may be misaligned with both the previous decisions and the target demonstrations. Given a set of candidates, the backward step first identifies the behavioral mode from the past decision stored in memory; the forward step then removes the samples with low likelihood under the target distribution using prior knowledge of positive and negative samples. By comparing samples across time steps and model checkpoints, our method bridges the gap between the proposal and target distributions during inference.

**Relation to recent methods.** Our method builds upon the receding horizon (Chi et al., 2023) and temporal ensembling (Zhao et al., 2023) used in previous works, but with crucial distinctions. Receding horizon seeks a compromise between long-term consistency and short-term reactivity by using a moderate action horizon (*e.g.*, half of the prediction horizon), which is inevitably sup-optimal when both factors are prominent. Temporal ensembling strengthens dependency across chunks by averaging multiple decisions over time; however, weighted-averaging operations can be detrimental when consecutive decisions fall into distinct modes. Our method more effectively addresses cross-chunk dependency through dedicated behavioral search and is not mutually exclusive to previous methods. We will demonstrate in the next section that combining our method with the moving average can further improve closed-loop action chunking.

## 5 EXPERIMENTS

In this section, we present a series of experiments to answer the following questions:

1. How does our theoretical analysis on action chunking manifest under different conditions?
2. Can our decoding method improve closed-loop operations of policies built with action chunking?
3. Does our decoding method perform well across different policies, tasks, and environments?
4. Is our decoding method scalable to larger sample sizes and compatible with existing methods?

To this end, we will first validate our theoretical analysis through one-dimensional diagnostic simulations. We will then evaluate BID on seven tasks across three simulation benchmarks, including Push-T (Chi et al., 2023), RoboMimic (Mandlekar et al., 2022), and Franka Kitchen (Gupta et al., 2020). We will subsequently examine the generality and scalability of our method under various base policies and sample sizes. We will finally assess the effectiveness of BID in two challenging real-world tasks that require interactions with dynamic objects.

## 5.1 ONE-DIMENSIONAL DIAGNOSTIC EXPERIMENTS

**Setup.** We start with a diagnostic experiment in a one-dimensional state space $\{s_0, s_1, \cdots, s_{10}\}$, where $s_0$ is the starting state and $s_{10}$ is the goal state. The demonstrator plans to move forward
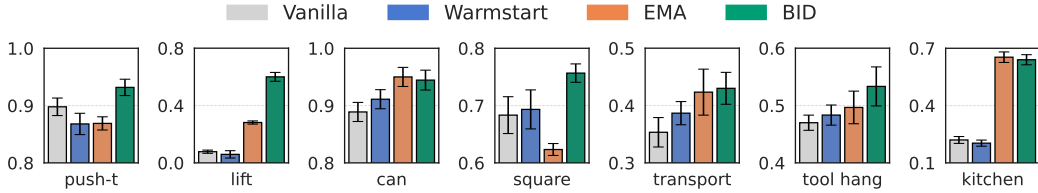
Figure 6: Comparison of different inference methods for closed-loop operations of diffusion policies. Each method is evaluated for 100 episodes on seven manipulation tasks in simulation benchmarks. Results are averaged across three seeds. BID significantly outperforms existing inference methods in most tasks.
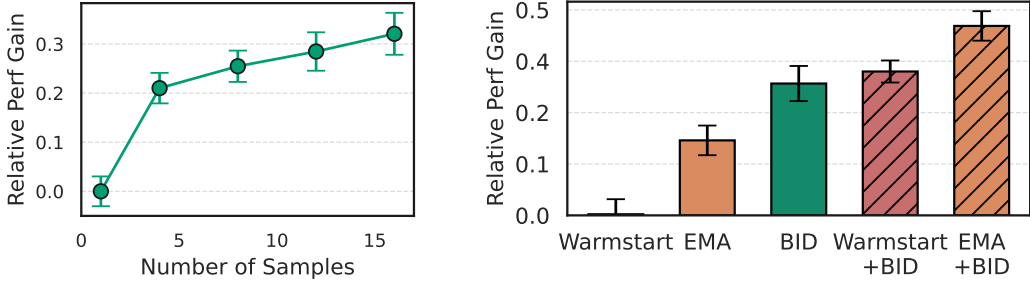


Figure 7: BID benefits from large sample sizes (left) and complements existing inference methods (right). Each method is evaluated on seven simulation tasks across three seeds. Relative performance gain is measured with respect to the vanilla baseline. When combined with EMA, BID results in 46% relative improvements.

by one step in each state, except in $s_5$ where it pauses unless the last five states visited were $s_5$. Each forward move has a probability of $1 - \delta$, where $\delta$ denotes the level of stochastic noise in the environment (as described in §3.2). Given these demonstrations, we train a collection of policies with different action horizons $h \in \{1, 2, 3, 5, 7, 10\}$. We investigate under what action horizon our learner can better imitate the distribution of idle actions taken by the expert over multiple rollouts.

**Result.** As shown in Fig. 5, when the environment is deterministic ($\delta = 0$), larger action horizons capture the expert distribution better, consistent with Corollary 2. With an action horizon of 10, the learner achieves zero total variation distance with the expert distribution. Conversely, when the environment is highly stochastic $\delta = 0.8$, an action horizon of 1 outperforms all other learners, corroborating with Corollary 3. With moderate noise $\delta = 0.4$, there is no discernible monotonic pattern due to the tradeoff revealed in Proposition 1. Refer to Appendix A for more detailed results.

## 5.2 SIMULATION EXPERIMENTS WITH STOCHASTIC NOISE

Next, we evaluate our decoding algorithm on seven simulation tasks of robot manipulation. We will first compare BID with existing inference methods in closed-loop operations. We will then assess the effectiveness of our method under different conditions, including policy classes, sample sizes, and levels of stochastic noise.

### 5.2.1 COMPARISON WITH EXISTING INFERENCE METHODS

**Setup.** In each manipulation task, we use Diffusion Policy (Chi et al., 2023) trained on human demonstrations as the base policy. We evaluate BID with a batch size of $N = 16$ and a mode size of $K = 3$. We consider three existing inference methods as baselines:

- *Vanilla (Chi et al., 2023)*: Execute the first action of a sampled chunk in a closed-loop manner.
- *Warmstart (Janner et al., 2022)*: Similar to *Vanilla*, but warm-start the initial noise for the diffusion process from the previous decision.
- *Exponential Moving Average (EMA) (Zhao et al., 2023)*: Smooth action chunking by averaging a new prediction $a$ with the previous one $\hat{a}$ for each overlapping step $a_t = \lambda a_t + (1 - \lambda)\hat{a}_t$. This method is also known as temporal ensembling. By default, we set $\lambda = 0.5$.

We evaluate each method for 100 episodes and average the results across three random seeds. Please refer to Appendix B for implementation details.

**Result.** Our main observation is that while existing inference methods offer some benefits for closed-loop operations, they lack robustness. As shown in Fig. 6, *Warmstart* yields mild perfor-

| Stochastic Noise | 0.0 | 1.0 | 1.5 |
|---|---|---|---|
| Vanilla Open-Loop | 61.0 | 39.0 | 19.4 |
| BID Open-Loop | **65.2** | 39.8 | 21.4 |
| Vanilla Closed-Loop | 52.0 | 50.4 | 44.2 |
| BID Closed-Loop | 56.6 | **54.8** | **54.4** |

Table 1: Success rates of VQ-BeT on the Push-T task under various conditions. BID consistently outperforms the vanilla counterpart. Closed-loop BID is particularly advantageous in stochastic settings.

| Sample Size | | Success (%) | Time (ms) |
|---|---|---|---|
| 1 | (vanilla) | 52.0 | 12.6 |
| 8 | (ours) | 53.8 | 25.6 |
| 16 | (ours) | 56.6 | 26.4 |
| 32 | (ours) | 56.6 | 27.3 |

Table 2: Success rates and inference times of VQ-BeT across varying sample sizes. BID benefits from a larger sample size at the cost of a doubled computational overhead, measured on an A5000 GPU.

mance gains on average, but degrades performance on 3 out of 7 tasks. Similarly, EMA leads to competitive results on several tasks, yet exhibits performance drops in 2 tasks. We conjecture that this robustness issue stems from independent sampling across chunks; when successive chunks follow distinct latent strategies, averaging them may not yield a plausible strategy, as further discussed in Appendix A.4. In comparison, BID offers substantial gains across all tasks. Notably, BID provides 32% relative improvements over the vanilla baseline, significantly outperforming EMA on Pust-T, Lift, Square, and Tool Hang, while achieving competitive performance on the other tasks.

### 5.2.2 SCALABILITY AND COMPATIBILITY OF BID

**Setup.** We further assess two key properties of BID: scalability with growing batch sizes and compatibility with existing inference methods. For scalability, we experiment with batch sizes of $\{1, 4, 8, 12, 16\}$. For compatibility, we apply BID on top of two other baselines, Warmstart and EMA. The results are averaged across seven simulation tasks and three random seeds.

**Result.** As shown in Fig. 7, our method benefits from the large batch size, with performance gains not yet saturated at the default batch size used in §5.2.1. Moreover, the strength from BID is complementary to that of existing inference methods. Notably, combining BID with EMA further boosts the relative performance gain from 32% to 46%. These two properties highlight the potential of our method in practice.

### 5.2.3 GENERALITY AND EFFICIENCY OF BID

**Setup.** We next extend our experiment to VQ-BET (Lee et al., 2024), another state-of-the-art robot policy built with autoregressive transformers. We use the public checkpoint on the Push-T task provided by LeRobot (Cadene et al., 2024) as the base policy. We use a checkpoint terminated at 100 epochs as the weak policy in forward contrast. To simulate stochastic conditions, we add temporally correlated Gaussian noise to the executed action at each step, scaled by the action magnitude. We measure the computational time on a desktop equipped with an NVIDIA A5000 GPU.

**Result.** Table 1 summarizes the results of the baseline and our method. The vanilla random sampling performs significantly worse than BID in both closed and open-loop operations. Notably, the vanilla open-loop approach exhibits a rapid performance decline as the environment becomes increasingly stochastic. Even in closed-loop operations, the vanilla baseline still experiences a significant performance drop. In comparison, the closed-loop BID demonstrates much higher robustness to stochastic noise. Table 2 details the computational overhead associated with BID at varying batch sizes. The result shows that the performance gains of our method come with a doubled computational overhead. We expect that this overhead will be less of a constraint with higher-end GPUs.

### 5.3 REAL-WORLD EXPERIMENTS WITH DYNAMIC OBJECTS

We finally evaluate BID through two real-world experiments that require rapid reactions to dynamically moving objects.

**Setup.** We consider two pick-and-place tasks, where the target object undergoes unexpected movement during evaluation. In the dynamic placing task, a Franka Panda robot is required to deliver an object into a moving cup held by a human subject. In the dynamic picking task, a UR5 robot is required to grasp a moving cup pulled by a string and place it onto a nearby static saucer. In both tasks, we evaluate the performance of BID applied to pre-trained diffusion policies. For further experimental details, please refer to Appendix B.2.

**Result.** Figs. 8 and 9 compares the results of different inference methods on the two real-world tasks. Vanilla random sampling struggles to handle the diverse demonstrations and dynamic movements,
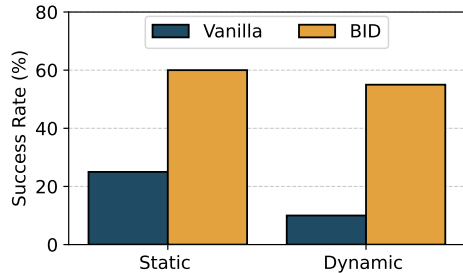
Figure 8: Success rate of object delivery. Each method-setting is evaluated across 20 episodes. BID achieves much higher success rate than the vanilla baseline, effectively handling the diverse demonstrations and dynamic target.
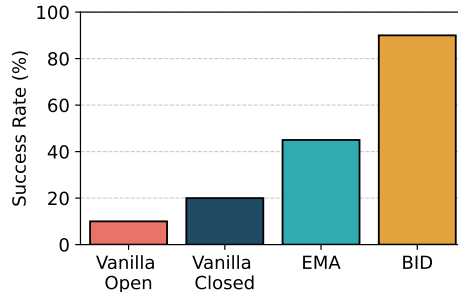
Figure 9: Success rate of cup replacement in the dynamic setting. Each method is evaluated across 20 episodes. Existing methods degrade substantially under slow cup movements, whereas BID retains a strong performance.

resulting in significantly lower success rates. In contrast, BID achieves high success rates in both static and dynamic conditions. Notably, in the dynamic picking task, BID achieves a 2x higher success rate than all other baselines, highlighting its potential for dynamic object interactions.

**Other experiments.** Please refer to Appendix A for additional analyses and ablations.

## 6 CONCLUSION

**Summary.** We have analyzed the strengths and limitations of action chunking for robot learning from human demonstrations. Based on our analysis, we proposed Bidirectional Decoding (BID), an inference algorithm that takes into account both past decisions and future plans for sample selection. Our experimental results show that BID can consistently improve closed-loop operations, scale well with computational resources, and complement existing methods. We hope these findings provide a new perspective on addressing the challenges of generative behavioral cloning at test time.

**Limitations.** One major limitation of BID lies in its computational complexity. While the decoding process can be parallelized on modern GPUs, it may remain prohibitive for high-frequency operations on low-cost robots. Designing algorithms that can generate quality yet diverse action chunks under batch size constraints can be an interesting avenue for future research. Additionally, our analysis and method have been limited to policies with short context lengths, driven by their empirical effectiveness with limited human demonstrations. Developing techniques capable of learning robust long-context policies can be another compelling direction for future research.

## REPRODUCIBILITY STATEMENT

To facilitate the reproducibility of our results, we release our code for the experiments with Diffusion Policy at `https://github.com/YuejiangLIU/bid_diffusion` and our code for the experiments with VQ-BET at `https://github.com/Jubayer-Hamid/bid_lerobot`. Additionally, we provide our detailed experimental setups in Appendix B, and complete proofs of our theoretical analysis in Appendix C.

## ACKNOWLEDGMENTS

## REFERENCES

Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 57(5):469–483, May 2009.

Christopher G. Atkeson and Stefan Schaal. Robot Learning From Demonstration. In *Proceedings of the Fourteenth International Conference on Machine Learning*, ICML '97, pp. 12–20, July 1997.

Suneel Belkhale, Yuchen Cui, and Dorsa Sadigh. Data quality in imitation learning. *Advances in Neural Information Processing Systems*, 36, 2024.

Homanga Bharadhwaj, Jay Vakil, Mohit Sharma, Abhinav Gupta, Shubham Tulsiani, and Vikash Kumar. Roboagent: Generalization and efficiency in robot manipulation via semantic augmentations and action chunking. *2024 IEEE International Conference on Robotics and Automation*, 2023a.

Homanga Bharadhwaj, Jay Vakil, Mohit Sharma, Abhinav Gupta, Shubham Tulsiani, and Vikash Kumar. RoboAgent: Generalization and Efficiency in Robot Manipulation via Semantic Augmentations and Action Chunking, September 2023b.

David Brandfonbrener, Stephen Tu, Avi Singh, Stefan Welker, Chad Boodoo, Nikolai Matni, and Jake Varley. Visual backtracking teleoperation: A data collection protocol for offline image-based reinforcement learning. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11336–11342. IEEE, 2023.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control, July 2023a.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: Robotics Transformer for Real-World Control at Scale. In *Robotics: Science and Systems XIX*. Robotics: Science and Systems Foundation, July 2023b. ISBN 978-0-9923747-9-2.

Remi Cadene, Simon Alibert, Alexander Soare, Quentin Gallouedec, Adil Zouitine, and Thomas Wolf. Lerobot: State-of-the-art machine learning for real-world robotics in pytorch. `https://github.com/huggingface/lerobot`, 2024.

Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion Policy: Visuomotor Policy Learning via Action Diffusion. In *Robotics: Science and Systems XIX*. Robotics: Science and Systems Foundation, July 2023. ISBN 978-0-9923747-9-2.

Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. *arXiv preprint arXiv:2402.10329*, 2024.

Open X.-Embodiment Collaboration, Abhishek Padalkar, Acorn Pooley, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Animesh Garg, Anthony Brohan, Antonin Raffin, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Brian Ichter, Cewu Lu, Charles Xu, Chelsea Finn, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Chuer Pan, Chuyuan Fu, Coline Devin, Danny Driess, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Federico Ceola, Fei Xia, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Giulio Schiavi, Gregory Kahn, Hao Su, Hao-Shu Fang, Haochen Shi, Heni Ben Amor, Henrik I. Christensen, Hiroki Furuta, Homer Walke, Hongjie Fang, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jan Peters, Jan Schneider, Jasmine Hsu, Jeannette Bohg, Jeffrey Bingham, Jiajun Wu, Jialin Wu, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jitendra Malik, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Joseph J. Lim, João Silvério, Junhyek Han, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Zhang, Krishan Rana, Krishnan Srinivasan, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Lionel Ott, Lisa Lee, Masayoshi Tomizuka, Max Spero, Maximilian Du, Michael Ahn, Mingtong Zhang, Mingyu Ding, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J. Joshi, Niko Suenderhauf, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Pannag R. Sanketi, Paul Wohlhart, Peng Xu, Pierre Sermanet, Priya Sundaresan, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Martín-Martín, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shuran Song, Sichun Xu, Siddhant Haldar, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Sudeep Dasari, Suneel Belkhale, Takayuki Osa, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Vidhi Jain, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiaolong Wang, Xinghao Zhu, Xuanlin Li, Yao Lu, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zhuo Xu, and Zichen Jeff Cui. Open X-Embodiment: Robotic Learning Datasets and RT-X Models, December 2023.

Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.

Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical Neural Story Generation. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 889–898, Melbourne, Australia, July 2018. Association for Computational Linguistics.

Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In *Conference on Robot Learning*, pp. 158–168. PMLR, 2022.

Markus Freitag and Yaser Al-Onaizan. Beam Search Strategies for Neural Machine Translation. In Thang Luong, Alexandra Birch, Graham Neubig, and Andrew Finch (eds.), *Proceedings of the First Workshop on Neural Machine Translation*, pp. 56–60, Vancouver, August 2017. Association for Computational Linguistics.

Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. Break the Sequential Dependency of LLM Inference Using Lookahead Decoding, February 2024a.

Zipeng Fu, Tony Z. Zhao, and Chelsea Finn. Mobile ALOHA: Learning Bimanual Mobile Manipulation with Low-Cost Whole-Body Teleoperation, January 2024b.

Kanishk Gandhi, Siddharth Karamcheti, Madeline Liao, and Dorsa Sadigh. Eliciting compatible demonstrations for multi-human imitation learning. In *Conference on Robot Learning*, pp. 1981–1991. PMLR, 2023.

Abraham George and Amir Barati Farimani. One act play: Single demonstration behavior cloning with action chunking transformers. *arXiv preprint arXiv:2309.10175*, 2023.

Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. In *Conference on Robot Learning*, pp. 1025–1037. PMLR, 2020.

Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine. IDQL: Implicit Q-Learning as an Actor-Critic Method with Diffusion Policies, May 2023.

John Hewitt, Christopher Manning, and Percy Liang. Truncation Sampling as Language Model Desmoothing. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 3414–3427, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.

Ryan Hoque, Ashwin Balakrishna, Ellen Novoseller, Albert Wilcox, Daniel S Brown, and Ken Goldberg. Thriftydagger: Budget-aware novelty and risk gating for interactive imitation learning. *arXiv preprint arXiv:2109.08273*, 2021a.

Ryan Hoque, Ashwin Balakrishna, Carl Putterman, Michael Luo, Daniel S Brown, Daniel Seita, Brijen Thananjeyan, Ellen Novoseller, and Ken Goldberg. Lazydagger: Reducing context switching in interactive imitation learning. In *2021 IEEE 17th international conference on automation science and engineering (case)*, pp. 502–509. IEEE, 2021b.

Wenlong Huang, Fei Xia, Dhruv Shah, Danny Driess, Andy Zeng, Yao Lu, Pete Florence, Igor Mordatch, Sergey Levine, Karol Hausman, and Brian Ichter. Grounded Decoding: Guiding Text Generation with Grounded Models for Embodied Agents. In *Thirty-Seventh Conference on Neural Information Processing Systems*, November 2023.

Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*, pp. 991–1002. PMLR, 2022.

Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with Diffusion for Flexible Behavior Synthesis. In *International Conference on Machine Learning (ICML)*, May 2022.

Liyiming Ke, Jingqiang Wang, Tapomayukh Bhattacharjee, Byron Boots, and Siddhartha Srinivasa. Grasping with chopsticks: Combating covariate shift in model-free imitation learning for fine manipulation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 6185–6191. IEEE, 2021.

Michael Kelly, Chelsea Sidrane, Katherine Driggs-Campbell, and Mykel J Kochenderfer. Hgdagger: Interactive imitation learning with human experts. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 8077–8083. IEEE, 2019.

Maxim Khanov, Jirayu Burapacheep, and Yixuan Li. ARGS: Alignment as Reward-Guided Search. In *The Twelfth International Conference on Learning Representations*, October 2023.

Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, Peter David Fagan, Joey Hejna, Masha Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, Arhan Jain, Abraham Lee, Youngwoon Lee, Marius Memmel, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Albert Zhan, Kevin Black, Cheng Chi, Kyle Beltran Hatch, Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pannag R. Sanketi, Archit Sharma, Cody Simpson, Quan Vuong, Homer Rich Walke, Blake Wulfe, Ted Xiao, Jonathan Heewon Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia, Rohan Baijal, Mateo Guaman Castro, Daphne Chen, Qiuyu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul Foster, Jensen Gao, David Antonio Herrera, Minho Heo, Kyle Hsu, Jiaheng Hu, Donovon Jackson, Charlotte Le, Yunshuang Li, Kevin Lin, Roy Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mirchandani, Daniel Morton, Tony Nguyen, Abigail O'Neill, Rosario Scalise, Derick Seale, Victor Son, Stephen Tian, Emi Tran, Andrew E. Wang, Yilin Wu, Annie Xie, Jingyun Yang, Patrick Yin, Yunchu Zhang, Osbert Bastani, Glen Berseth, Jeannette Bohg, Ken Goldberg, Abhinav Gupta, Abhishek Gupta, Dinesh Jayaraman, Joseph J. Lim, Jitendra Malik, Roberto Martín-Martín, Subramanian Ramamoorthy, Dorsa Sadigh, Shuran Song, Jiajun Wu, Michael C. Yip, Yuke Zhu, Thomas Kollar, Sergey Levine, and Chelsea Finn. DROID: A Large-Scale In-The-Wild Robot Manipulation Dataset, March 2024.

Alex Kuefler and Mykel J Kochenderfer. Burn-in demonstrations for multi-modal imitation learning. *arXiv preprint arXiv:1710.05090*, 2017.

Shankar Kumar and William Byrne. Minimum Bayes-Risk Decoding for Statistical Machine Translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pp. 169–176, Boston, Massachusetts, USA, May 2004. Association for Computational Linguistics.

Lucy Lai, Ann Zixiang Huang, and Samuel J Gershman. Action chunking as policy compression. 2022.

Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection for robust imitation learning. In *Conference on robot learning*, pp. 143–156. PMLR, 2017.

Seungjae Lee, Yibin Wang, Haritheja Etukuru, H. Jin Kim, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Behavior Generation with Latent Actions, March 2024.

Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast Inference from Transformers via Speculative Decoding. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 19274–19286. PMLR, July 2023.

Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive Decoding: Open-ended Text Generation as Optimization. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12286–12312, Toronto, Canada, July 2023. Association for Computational Linguistics.

Yunzhu Li, Jiaming Song, and Stefano Ermon. Infogail: Interpretable imitation learning from visual demonstrations. *Advances in neural information processing systems*, 30, 2017.

Xiao Ma, Sumit Patidar, Iain Haughton, and Stephen James. Hierarchical Diffusion Policy for Kinematics-Aware Multi-Task Robotic Manipulation, March 2024.

Ajay Mandlekar, Fabio Ramos, Byron Boots, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Dieter Fox. IRIS: Implicit Reinforcement without Interaction at Scale for Learning Control from Offline Robot Manipulation Data. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4414–4420, May 2020.

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. In *Proceedings of the 5th Conference on Robot Learning*, pp. 1678–1690. PMLR, January 2022.

Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. Locally Typical Sampling. *Transactions of the Association for Computational Linguistics*, 11:102–121, 2023.

Kunal Menda, Katherine Driggs-Campbell, and Mykel J Kochenderfer. Ensembledagger: A bayesian approach to safe imitation learning. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5041–5048. IEEE, 2019.

Mathias Müller and Rico Sennrich. Understanding the Properties of Minimum Bayes Risk Decoding in Neural Machine Translation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 259–272, Online, August 2021. Association for Computational Linguistics.

Mitsuhiko Nakamoto, Oier Mees, Aviral Kumar, and Sergey Levine. Steering Your Generalists: Improving Robotic Foundation Models via Value Guidance. In *8th Annual Conference on Robot Learning*, September 2024.

Aaditya Prasad, Kevin Lin, Jimmy Wu, Linqi Zhou, and Jeannette Bohg. Consistency Policy: Accelerated Visuomotor Policies via Consistency Distillation, May 2024.

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*, pp. 627–635. JMLR Workshop and Conference Proceedings, 2011.

Nur Muhammad Shafiullah, Zichen Cui, Ariuntuya Arty Altanzaya, and Lerrel Pinto. Behavior transformers: Cloning $k$ modes with one stone. *Advances in neural information processing systems*, 35:22955–22968, 2022.

Lucy Xiaoyang Shi, Archit Sharma, Tony Z. Zhao, and Chelsea Finn. Waypoint-Based Imitation Learning for Robotic Manipulation. In *Proceedings of The 7th Conference on Robot Learning*, pp. 2195–2209. PMLR, December 2023.

Aravind Sivakumar, Kenneth Shaw, and Deepak Pathak. Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube. *Robotics: Science and Systems (RSS)*, 2022.

Gokul Swamy, Sanjiban Choudhury, Drew Bagnell, and Steven Wu. Causal imitation learning under temporally correlated noise. In *International Conference on Machine Learning*, pp. 20877–20890. PMLR, 2022.

Ashwin Vijayakumar, Michael Cogswell, Ramprasaath Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. Diverse Beam Search for Improved Description of Complex Scenes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), April 2018.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-Consistency Improves Chain of Thought Reasoning in Language Models. *arXiv:2203.11171 [cs]*, April 2022.

Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural Text Generation With Unlikelihood Training. In *International Conference on Learning Representations*, September 2019.

Philipp Wu, Yide Shentu, Zhongke Yi, Xingyu Lin, and Pieter Abbeel. Gello: A general, low-cost, and intuitive teleoperation framework for robot manipulators. *arXiv preprint arXiv:2309.13037*, 2023.

Annie Xie, Dylan Losey, Ryan Tolsma, Chelsea Finn, and Dorsa Sadigh. Learning Latent Representations to Influence Multi-Agent Interaction. In *Proceedings of the 2020 Conference on Robot Learning*, pp. 575–588. PMLR, October 2021.

Xiaomeng Xu, Huy Ha, and Shuran Song. Dynamics-Guided Diffusion Model for Robot Manipulator Design, February 2024.

Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning Fine-Grained Bimanual Manipulation with Low-Cost Hardware. In *Robotics: Science and Systems (RSS) 2023*, April 2023.
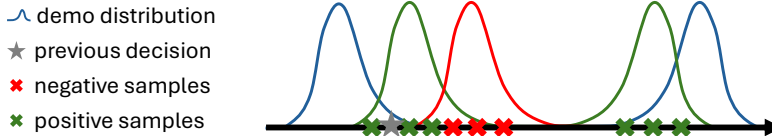
Figure 10: Distributional interpretation of BID. The backward criterion (Equation 6) favors samples close to the past decision; the forward criterion (Equation 7) promotes samples with a high likelihood under the target distribution.

| Noise Level | Action Horizon | | | | |
|---|---|---|---|---|---|
| | 1 | 3 | 5 | 7 | 10 |
| 0.0 | 4.21 | 1.75 | 1.55 | 1.28 | **0.00** |
| 0.4 | 0.55 | **0.30** | 0.95 | 0.53 | 0.93 |
| 0.8 | **0.04** | 0.98 | 1.23 | 1.26 | 1.44 |

Table 3: Total variation distance between the action distributions of each model and the expert in environments with varying noise levels. Lower values indicate better performance.

## A  ADDITIONAL EXPERIMENTS

### A.1  ONE-DIMENSIONAL SIMULATIONS

In addition to Fig. 5, we summarize the total variation distance between each learned policy and the demonstration in the one-dimensional simulation. Our results indicate that a shorter action horizon is more effective in noisier environments, whereas a longer action horizon yields better performance in static environments.

### A.2  ACTION HORIZON VS. CONTEXT LENGTH

**Setup.** Our work builds on the premise that the action horizon is longer than the context length, as commonly designed for recent policies. While BID mitigates the inherent limitations of this design choice through test-time decoding, an important question remains: could extending the context length yield stronger policies? To understand this, we trained diffusion policies with varying combinations of action horizons and context lengths on the Push-T task. Specifically, we use a short context length ($c = 2$) and a short action horizon ($h = 2$) as our baseline, and incrementally increase these parameters to larger values $6, 10, 14$ to assess their impact.

**Result.** Fig. 11 compares the performance of the policy learned with different $\Delta h = h - c$. As expected, the policy with both a short action horizon and a short context length does not perform well, due to its limited capability to model long-range temporal dependencies. Interestingly, extending the context length initially boosts performance ($\Delta h = -4$), but this trend reverses as the context length becomes too long ($\Delta h \leq -8$), likely due to overfitting to an increased number of spurious features. In contrast, expanding the action horizon results in more robust performance improvements, validating its pivotal role in imitation learning from human demonstrations.

### A.3  ABLATION STUDY OF FORWARD CONTRAST

**Setup.** To understand the effect of forward contrast (Equation 7), we evaluate the full version of our method against three reduced variants in open-loop operations: *Vanilla* (without forward contrast), *Positive* (without negative samples), and *Negative* (without positive samples). Similar to §5.2.2, our ablation study is conducted in seven simulation tasks, each with three random seeds.

**Result.** Fig. 12 summarizes the result of this ablation study. Notably, both positive and negative samples are essential for effective sample selection, and omitting either leads to significant performance declines. Without negative samples, our decoding method reduces to an approximate maximum a posteriori estimation, which can result in suboptimal decisions due to modeling errors. Conversely, without positive samples, the sampling process may be biased towards rare instances. This result highlights the importance of both components and suggests the potential for extending this paradigm in future work.
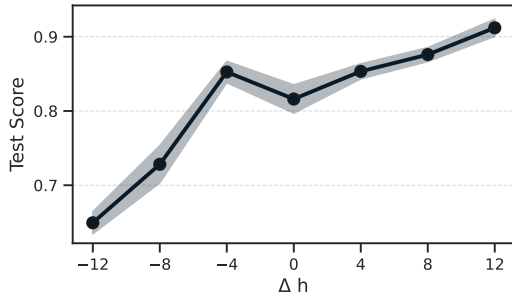
Figure 11: Effect of prediction horizon ($h$) and context length ($c$) on diffusion policies in the Push-T task. The baseline is set at $h = 2$ and $c = 2$, with $\Delta h = h - c = 0$. Extending the prediction horizon ($h > 2$) consistently improves performance, whereas extending the context length ($c > 2$) can cause substantial performance declines. Each model is trained for $5k$ epochs. Results are averaged over the last five checkpoints.
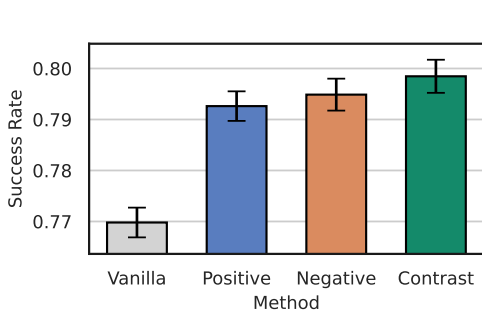


Figure 12: Effect of positive and negative samples on forward contrast. Performance of ablated variants of forward contrast is evaluated across seven simulation tasks. The absence of either positive or negative samples prevents achieving the full performance gains observed with the contrast objective.
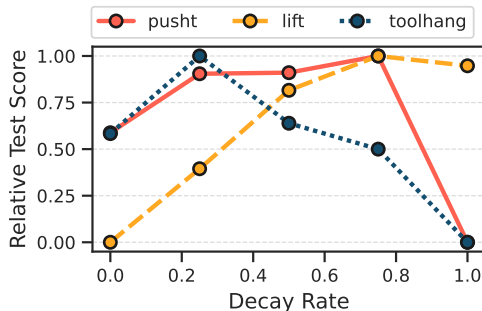
Figure 13: Effect of the decay rate for the exponential moving average. In each task, we measure the relative performance among different decay rates. The optimal decay rate varies by task, leading to a practical challenge of identifying a universal temporal ensembling strategy (Zhao et al., 2023).

## A.4 CHALLENGES FOR TEMPORAL ENSEBMLING

EMA exhibits competitive performance in Fig. 6. However, tuning its decay rate can be difficult in practice. Fig. 13 shows the sensitivity of EMA to the decay rate across three different tasks, where the optimal choices differ significantly. We conjecture that this high sensitivity stems from the variability in the latent strategies between consecutive predictions. When consecutive predictions follow similar strategies, a lower decay rate (*i.e.*, stronger moving average) can enhance smoothness and improve performance. Conversely, when consecutive predictions diverge in their underlying strategies, averaging them can introduce adverse effects. Our method promotes coherence in latent strategies and thus effectively complements temporal ensembling, as evidenced in Fig. 7.

## B ADDITIONAL DETAILS

### B.1 SIMULATION EXPERIMENT DETAILS

#### B.1.1 ENVIRONMENT DETAILS

Our simulation experiments are conducted on three robot manipulation benchmarks. We use the training data collected from human demonstrations in each benchmark.

*Push-T*: We adopt the Push-T environment introduced in (Chi et al., 2023), where the goal is to push a T-shaped block on a table to a target position. The action space is two-dimensional end-effector velocity control. The training dataset contains 206 demonstrations collected by humans.

*Robomimic*: We use five tasks in the Robomimic suite (Mandlekar et al., 2022), namely Lift, Can, Square, Transport, and Tool Hang. The training dataset for each task contains 300 episodes collected from multi-human (MH) demonstrations.

| name | value |
|---|---|
| batch size $N$ | 16 |
| mode size $K$ | 3 |
| prediction horizon $l$ | 16 |
| temporal coherence decay $\rho$ | 0.5 |
| moving average decay $\lambda$ | 0.5 |

Table 4: Default hyper-parameters in our experiments.

*Franka Kitchen*: We use the Franka Kitchen environment from (Gupta et al., 2020), featuring a Franka Panda arm with a seven-dimensional action space and 566 human-collected demonstrations. The learned policy is evaluated on test cases involving four or more objects (p4), a challenging yet practical task for robotic manipulation in household contexts.

### B.1.2 IMPLEMENTATION DETAILS.

Our implementation of BID for Diffusion Policy is built upon the official code of Chi et al. (2023), with modifications made solely to the inference process. The policy takes in state inputs and predicts a chunk of 16 actions as outputs. For each simulation task, we train the model for 100-1000 epochs to reach near-optimal performance. We evaluate it in closed-loop operations, *i.e.*, action horizon is set to 1. For forward contrast, we train the weak policy for 10-100 epochs, resulting in a suboptimal policy for each task. The core hyperparameters are summarized in Table 4.

Our implementation of BID for VQ-BeT (Lee et al., 2024) is built upon the code of LeRobot (Cadene et al., 2024). We use the best public checkpoint as the strong policy and a checkpoint trained for 100k iterations as the weak policy. Since BID requires sample diversity, we set the temperature as 0.5 for our methods, while vanilla sampling uses the default temperature of 0.1.

## B.2 REAL-WORLD EXPERIMENT DETAILS

### B.2.1 DYNAMIC PLACING

**Task.** We consider a task where the robot is to deliver an object held in its gripper into a cup held by a human. As shown in Fig. 14, this task comprises four main stages and presents two core challenges. First, due to the similar size of the object and the cup, the robot must achieve high precision to place the object accurately into the cup. Second, the position of the cup is not fixed, requiring the robot to adjust its plans based on the latest position continuously. This task mirrors real-world scenarios where robots interact with a dynamic environment, accommodating moving objects and agents.

**Demonstration.** In light of temporal dependencies and style variations in human behaviors, we intentionally collect a diverse set of demonstration data, differing in factors such as average speed, idling pause, and overall trajectory. We gather a total of 150 demonstration episodes: 50 clean and consistent demonstrations, and 100 noisy and diverse demonstrations. All demonstrations successfully accomplish the task. Additional, the location of the cup is fixed and static within each episode.

**Robot.** Following previous works (Chi et al., 2023; Prasad et al., 2024), we use a Franka Panda as the robot hardware and the vision-based diffusion policy for its operation. The robot is equipped with two cameras: one ego-centric camera mounted at the wrist of the robot, one third-person camera mounted at a static bracket. Both cameras provide visual observations at a resolution of $256 \times 256$ pixels. The robot operates at a frequency of 10 Hz, with a prediction horizon of 16 time steps.

**Evaluation.** We evaluate our method in comparison to vanilla random sampling under two conditions: *static target*, where the target cup remains fixed throughout the evaluation, and *dynamic target*, where the target cup is gradually moved. In the dynamic setting, the location of the cup stays within the range of training locations, but the movement is not encountered during training. This evaluation protocol is designed to explicitly assess the ability of the policy to react to unexpected dynamics in the environment. Each method-setting pair is tested over 20 episodes, with both the initial and target locations randomized across different episodes.

**Result.** We summarize the result of the real-world experiments in Fig. 8. The success rate of vanilla random sampling is generally limited due to oscillations between different latent strategies, which quickly diverge from the distribution of demonstrations. This issue is particularly pronounced in the dynamic setting, where the vanilla baseline struggles to account for the target movements

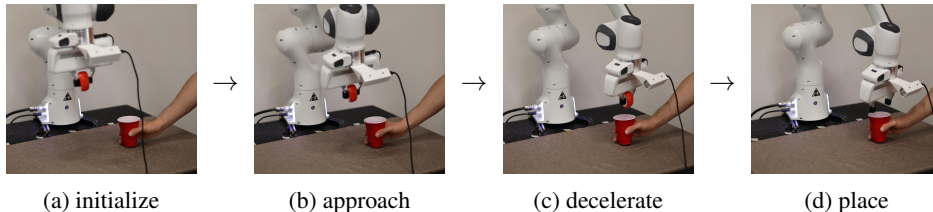(a) initialize      (b) approach      (c) decelerate      (d) place

Figure 14: Human demonstrations on a Franka Panda robot for a real-world object delivery task. The robot is tasked with delivering an object held in its gripper into a cup held by a human. Each demonstration consists of four main stages: (a) initialize the robot position randomly, (b) approach the target cup, (c) slow down near the target cup, and (d) release the object. The position of the target cup may change during an episode.



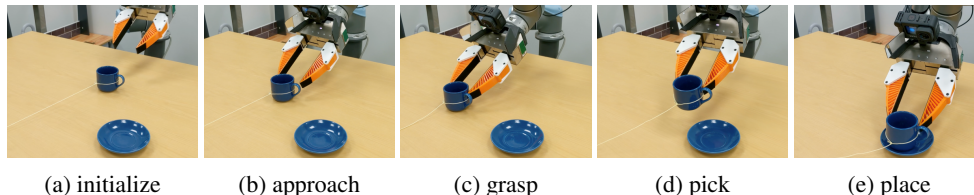(a) initialize     (b) approach     (c) grasp     (d) pick     (e) place

Figure 15: The robot is tasked with picking up a cup and placing it on a saucer nearby. The four main stages are (a) initializing the robot, (b) approaching the target cup, (c) grasping the target cup, (d) picking up the cup, and (e) placing the cup on the target saucer. The position of the target cup may change during an episode.

within an action chunk lasting for 1.6 seconds. In contrast, the proposed BID method significantly improves performance in both static and dynamic settings. Notably, BID maintains a similar success rate in the dynamic setting as in the static setting, suggesting its potential to extend action chunking into uncertain environments.

### B.2.2 DYNAMIC PICKING

**Task.** Next, we consider a task where the robot is required to pick up a cup and place it onto a nearby saucer. The cup was pulled with a string until the robot's gripper successfully grasped it. The task consists of five main stages, which are illustrated in Fig. 15. This setup also tests the robot's capability to interact with a dynamic environment, a critical challenge in real-world applications.

**Policy.** We utilized the publicly available diffusion policy checkpoint from UMI (Chi et al., 2024) without any additional fine-tuning. Notably, the policy was originally trained using demonstrations in a static setting, where the cup's position remained constant throughout the task. Our experimental setup mirrored the one described by UMI, using the same UR5 robot hardware. This allowed us to directly evaluate the policy's transferability to a dynamic environment, where the cup's position changes during the task. Due to the absence of an early checkpoint, we omitted negative samples in forward contrast, focusing solely on positive consistency discussed in Appendix A.3.

**Evaluation.** We evaluated BID against three baselines: vanilla random sampling in both open-loop and closed-loop configurations, and EMA (closed-loop). These methods were tested under two conditions: *static target*, where the cup remained in a fixed position, and *dynamic target*, where the cup was moved using the string. Each method-setting combination was tested across 20 episodes, with the initial positions of the cup and saucer kept consistent to ensure controlled comparisons.

**Results.** The results, summarized in Fig. 9, highlight the challenges of the dynamic setting. Open-loop vanilla sampling performed poorly due to its inability to adapt to the cup's movements, often failing to approach the cup as it was pulled. While closed-loop vanilla sampling showed improved reactivity, it suffered from inconsistent trajectories, resulting in jittery behavior when attempting to grasp and place the cup. Similarly, closed-loop EMA sampling demonstrated higher adaptability to environmental changes but often failed to firmly grasp the cup, likely due to the limitations of naive averaging, which compromises commitment to a specific strategy. In contrast, BID achieved at least a 2x improvement in success rate compared to all other methods in the dynamic setting, while maintaining its performance in the static setting, demonstrating both adaptability and precision in dynamic environments.

## C  PROOFS

First, we establish the following lemma which will help us compare different function classes:

**Lemma 4.** *Let $\mathcal{L}$ be a convex function and let $X$ and $Y$ be two random variables. Let $G$ be the class of functions $g(X)$ that accept $X$ as an input. Then*

$$\min_{g(X) \in G} \mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X,Y), g(X)) \right] = \mathbb{E}_X \left[ \min_{c \in \mathbb{R}} \mathbb{E}_Y \left[ \mathcal{L}(f(X,Y), c) | X \right] \right].$$

*Proof.* The left hand side is less than or equal to the right hand side by the following logic:

$$\mathbb{E}_X \left[ \min_{c \in \mathbb{R}} \mathbb{E}_Y \left[ \mathcal{L}(f(X,Y), c) | X \right] \right] = \mathbb{E}_X \left[ \mathbb{E}_Y \left[ \mathcal{L}(f(X,Y), c^*(X)) | X \right] \right]$$

$$\geq \min_{g(X) \in G} \mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X,Y), g(X)) \right]$$

where we used $c^*(X) := \arg\min_c \mathbb{E}_X[\mathcal{L}(f(X,Y), c)|X]$. We get the inequality by recognizing that $\mathbb{R} \subsetneq G$. On the other hand, the left hand side is greater than or equal to the right hand side. For any $g(X)$, we have:

$$\mathbb{E}[\mathcal{L}(f(X,Y), g(X))] = \mathbb{E}_X \left[ \mathbb{E}_Y \left[ \mathcal{L}(f(X,Y), g(X)) | X \right] \right]$$

$$\geq \mathbb{E}_X \left[ \min_g \mathbb{E}_Y \left[ \mathcal{L}(f(X,Y), g(X)) | X \right] \right]$$

$$= \mathbb{E}_X \left[ \min_c \mathbb{E}_Y \left[ \mathcal{L}(f(X,Y), c) | X \right] \right].$$

With these two inequalities, we conclude. $\square$

Next, we prove the following lemma. This straightforward, and almost trivial, result is provided as a separate lemma because we simplify terms in this manner quite often throughout our proofs.

**Lemma 5.** *Let $\mathcal{L}$ be a convex function and let $X, Y$ be two random variables. Then,*

$$\min_f \mathbb{E}_{X,Y} \left[ P(X' = X)\mathcal{L}(f(X'), S(X,Y)) \right] + \mathbb{E}_{X,Y} \left[ \sum_{X' \neq X} P(X')\mathcal{L}\left(f(X'), S(X,Y)\right) \right]$$

$$\leq \min_f \{ \mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X), S(X,Y)) \right] \} + \epsilon$$

*where $\epsilon = \max_{X' \neq X, X, Y} \{\mathcal{L}(f^*(X'), S(X,Y)\}$ and $f^* = \arg\min_f \{\mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X), S(X,Y)) \right]\}$*

*Proof.*

$$\min_f \mathbb{E}_{X,Y} \left[ P(X' = X)\mathcal{L}(f(X'), S(X,Y)) \right] + \mathbb{E}_{X,Y} \left[ \sum_{X' \neq X} P(X')\mathcal{L}\left(f(X'), S(X,Y)\right) \right]$$

$$\leq \min_f \mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X), S(X,Y)) \right] + \mathbb{E}_{X,Y} \left[ \sum_{X' \neq X} P(X')\mathcal{L}\left(f(X'), S(X,Y)\right) \right]$$

$$\leq \min_f \{ \mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X), S(X,Y)) \right] \} + \mathbb{E}_{X,Y} \left[ \sum_{X' \neq X} P(X')\mathcal{L}\left(f^*(X'), S(X,Y)\right) \right]$$

$$\leq \min_f \{ \mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X), S(X,Y)) \right] \} + \mathbb{E}_{X,Y} \left[ \sum_{X' \neq X} P(X')\epsilon \right]$$

$$\leq \min_f \{ \mathbb{E}_{X,Y} \left[ \mathcal{L}(f(X), S(X,Y)) \right] \} + \epsilon$$

$\square$

## C.1 Definitions

Note that the informal definitions provided in §3.2 for expected observation advantage and maximum inference disadvantage have some deviations from the mathematical definitions. The informal definitions only attempt to provide intuition for what these terms might mean but they are not sufficient to describe the full construction.

To define the terms formally, we, first, analyze the effect of reducing context horizon. We show that, provided action horizon is constant, decreasing context horizon causes performance of the optimal policy to decrease.

Consider a $(c, h)$-policy whose probability of taking action $a_t$ at time $t$ in a chunk generated at $t$ is referred to as

$$\pi_{(c,h)} := \pi_{(c,h)}(a_t | s_{t-c:t}).$$

On the other hand, consider a $(c + 1, h)$-policy whose probability of taking action $a_t$ in a chunk generated at time $t$ is referred to as

$$\pi_{(c+1,h)} := \pi_{(c+1,h)}(a_t | s_{t-c-1:t}).$$

Lastly, consider a $(k, 1)$-expert whose probability of taking action $a_t$ at time $t$ is $\pi^*$.

**Proposition 6** (Backward Context is valuable). *Let $\mathcal{L}$ be a non-linear, convex function. Let $c < k$. Let $G := \{a_t, s_{t-k:t-c-1}, z_{t-k:t}\}$ and let $C := \{s_{t-c:t}\}$. Then,*

$$\min_{\pi_{(c+1,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+1,h)}, \pi^*) \Big| C \right] \leq \min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*) \Big| C \right]$$

*Proof.* We refer to the class of functions that accept $a_t$ and $s_{t-c-1:t}$ as inputs as $X_{+,+}$. Similarly, the class of functions that do not accept $a_t$ as inputs but accept $s_{t-c-1:t}$ as inputs is $X_+$. The function class that accepts only $s_{t-c:t}$ and not $s_{t-c-1}$ or $a_t$ as inputs are elements of $X_-$. Lastly, the function class that accepts $s_{t-c:t}$ and $a_t$ as inputs, but not $s_{t-c-1}$, are elements of $X_{-,-}$.

$$\min_{\pi_{(c+1,h)} \in X_{+,+}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+1,h)}, \pi^* \Big| C \right]$$

$$= \mathbb{E}_{a_t} \left[ \min_{\pi'_{(c+1,h)} \in X_+} \mathbb{E}_{s_{t-c-1}} \left[ \mathbb{E}_{s_{t-k:t-c-2}, z_{t-k:t}} \left[ \mathcal{L}(\pi'_{(c+1,h)}, \pi^*) \Big| a_t, s_{t-c-1}, C \right] \Big| a_t, C \right] \Big| C \right]$$
(Lemma 4)

$$= \mathbb{E}_{a_t} \left[ \mathbb{E}_{s_{t-c-1}} \left[ \min_{\pi'_{(c,h)} \in X_-} \mathbb{E}_{s_{t-k:t-c-2}, z_{t-k:t}} \left[ \mathcal{L}(\pi'_{(c,h)}, \pi^*) \Big| a_t, s_{t-c-1}, C \right] \Big| a_t, C \right] \Big| C \right] \quad \text{(Lemma 4)}$$

$$\leq \mathbb{E}_{a_t} \left[ \min_{\pi'_{(c,h)} \in X_-} \mathbb{E}_{s_{t-c-1}} \left[ \mathbb{E}_{s_{t-k:t-c-2}, z_{t-k:t}} \left[ \mathcal{L}(\pi'_{(c,h)}, \pi^*) \Big| a_t, s_{t-c-1}, C \right] \Big| a_t, C \right] \Big| C \right]$$
(Jensen's inequality)

$$= \min_{\pi_{(c,h)} \in X_{-,-}} \mathbb{E}_{a_t} \left[ \mathbb{E}_{s_{t-c-1}} \left[ \mathbb{E}_{s_{t-k:t-c-2}, z_{t-k:t}} \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*) \Big| a_t, s_{t-c-1}, C \right] \Big| a_t, C \right] \Big| C \right] \quad \text{(Lemma 4)}.$$

Use the law of total expectation to conclude. $\qquad \square$

Now, we formalize the definitions of *Expected Observation Advantage* and *Maximum Inference Disadvantage*.

Recall that, in §3.2, we have two policies: $\pi_{(c,h)}$ and $\pi_{(c,h+d)}$; the former sees more recent states while the latter remembers more past states. First, we define an agent that gets access to all the information that both learners, combined, have: a $(c + d, h)$-policy whose probability of taking action $a_t$ in a chunk generated at time $t - h$ is

$$\pi_{(c+d,h)} := \pi_{(c+d,h)}(a_t | s_{t-h-d-c:t-h}, a_{t-h:t-1}).$$

Observe that $\pi_{(c+d,h)}$ has access to more context than $\pi_{(c,h)}$, particularly the knowledge of states $s_{t-h-c-d:t-h-c-1}$.

**Definition (Expected Observation Advantage ($\alpha_b$)).** We know, using Proposition 6, $\pi_{(c+d,h)}$ has lower divergence with respect to $\pi^*$ than $\pi_{(c,h)}$. We say that the advantage $\pi_{(c+d,h)}$ gets from the extra information is $\alpha_b$. More formally, we say that

$$0 \leq \alpha_b = \min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*)) \Big| C \right] - \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] \qquad (8)$$

where $C = \{s_{t-h-c:t-h}, a_{t-h:t-1}\}$ and $G = \{a_t, s_{t-k:t-h-c-1}, s_{t-h+1:t}, z_{t-k:t}\}$. In particular, $\alpha_b = 0$ when $s_{t-h-d-c:t-h-c-1}$ can be deterministically inferred by $\pi_{(c,h)}$ or when the expert policy is independent of them. However, this is extremely unlikely since $\pi_{(c,h)}$ does not know the actions taken in those time steps (even more unlikely in a stochastic environment) and the expert's action depends on the last $k$ time steps.

**Definition (Maximum Inference Disadvantage ($\epsilon_f$)).** Consider the maximum divergence that can be accumulated by the $(c, h + d)$-policy from not knowing the recent states at time steps $s_{t-h-d+1:t-h}$ and let that be $\epsilon_f$. More formally, we say that, for fixed $C$ from Proposition 1, any state in $\mathcal{S}^-$ and any $z_{t-k:t}$ and any $\hat{s}_{t-h-d+1:t-h} \neq s_{t-h-d+1:t-h}$:

$$\mathcal{L}(\pi_{(c+d,h)}(a_t | s_{t-h-d-c:t-h-d}, \hat{s}_{t-h-d+1:t-h} \neq s_{t-h-d+1:t-h}, a_{t-h:t-1}), \pi^*) \leq \epsilon_f. \qquad (9)$$

Here, $\pi_{(c+d,h)} := \arg\min_{\pi_{(c+d,h)}} \mathbb{E}_G[\mathcal{L}(\pi_{(c+d,h)}, \pi^*)|C]$ is the optimal $(c + d, h)$-policy.

Intuitively, maximum inference disadvantage captures the maximum dependency on relative *time steps* whereas expected observation advantage captures the expected advantage from observing some given *states*.

To define $\alpha_f$ and $\epsilon_b$, we prove a second version of Proposition 6. Consider a $(c, h)$-policy whose probability of taking action $a_t$ at time $t$ in a chunk generated at $t$ is referred to as

$$\pi_{(c,h)} := \pi_{L(c,h)}(a_t | s_{t-c:t}).$$

On the other hand, consider a $(c - 1, h + 1)$-policy whose probability of taking action $a_t$ in a chunk generated at time $t - 1$ is referred to as

$$\pi_{(c-1,h+1)} := \pi_{(c-1,h+1)}(a_t | s_{t-c:t-1}).$$

Lastly, consider a $(k, 1)$-expert whose probability of taking action $a_t$ at time $t$ is $\pi^*$.

**Proposition 7** (Forward Context is valuable). *Let $\mathcal{L}$ be a non-linear, convex function. Let $c < k$. Let $G := \{a_t, s_{t-k:t-c-1}, s_t, z_{t-k:t}\}$ and let $C := \{s_{t-c:t-1}, a_{t-1}\}$. Then,*

$$\min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*) \Big| C \right] \leq \min_{\pi_{(c-1,h+1)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c-1,h+1)}, \pi^*) \Big| C \right]$$

*Proof.* The proof is similar to that of Proposition 6. We refer to the class of functions that accept $a_t$ and $s_{t-c:t}$ as inputs as $X_{+,+}$. Similarly, the class of functions that do not accept $a_t$ as inputs but accept $s_{t-c:t}$ as inputs is $X_+$. The function class that accepts only $s_{t-c:t-1}$ and not $s_t$ or $a_t$ as inputs are elements of $X_-$. Lastly, the function class that accepts $s_{t-c:t-1}$ and $a_t$ as inputs, but not $s_t$, are elements of $X_{-,-}$.

$$\min_{\pi_{(c,h)} \in X_{+,+}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^* \Big| C \right]$$

$$= \mathbb{E}_{a_t} \left[ \min_{\pi'_{(c,h)} \in X_+} \mathbb{E}_{s_t} \left[ \mathbb{E}_{s_{t-k:t-c-1}, z_{t-k:t}} \left[ \mathcal{L}(\pi'_{(c,h)}, \pi^*) \Big| a_t, s_t, C \right] \Big| a_t, C \right] \Big| C \right] \qquad \text{(Lemma 4)}$$

$$= \mathbb{E}_{a_t} \left[ \mathbb{E}_{s_t} \left[ \min_{\pi'_{(c-1,h+1)} \in X_-} \mathbb{E}_{s_{t-k:t-c-2}, z_{t-k:t}} \left[ \mathcal{L}(\pi'_{(c-1,h+1)}, \pi^*) \Big| a_t, s_t, C \right] \Big| a_t, C \right] \Big| C \right] \qquad \text{(Lemma 4)}$$

$$\leq \mathbb{E}_{a_t} \left[ \min_{\pi'_{(c-1,h+1)} \in X_-} \mathbb{E}_{s_t} \left[ \mathbb{E}_{s_{t-k:t-c-1}, z_{t-k:t}} \left[ \mathcal{L}(\pi'_{(c-1,h+1)}, \pi^*) \Big| a_t, s_t, C \right] \Big| a_t, C \right] \Big| C \right]$$
$$\text{(Jensen's inequality)}$$

$$= \min_{\pi_{(c-1,h+1)} \in X_{-,-}} \mathbb{E}_{a_t} \left[ \mathbb{E}_{s_t} \left[ \mathbb{E}_{s_{t-k:t-c-1}, z_{t-k:t}} \left[ \mathcal{L}(\pi_{(c-1,h+1)}, \pi^*) \Big| a_t, s_t, C \right] \Big| a_t, C \right] \Big| C \right] \qquad \text{(Lemma 4)}.$$

Use the law of total expectation to conclude. $\qquad \square$

Using this, we can define $\epsilon_b$ and $\alpha_f$ in a similar manner:

**Definition (Expected Observation Advantage ($\alpha_f$)).** Recall that we have two models: $\pi_{(c,h)}$ and $\pi_{(c,h+d)}$ and a hypothetical $(c, h + d)$-policy that has access to all the information both our learners have (as in Eq. (8) and Eq. (9)). Observe that $\pi_{(c+d,h)}$ has access to more context than $\pi_{(c,h+d)}$, particularly the knowledge of states $s_{t-h-d+1:t-h}$. Therefore, we know, using Proposition 7, $\pi_{(c+d,h)}$ has lower divergence with respect to $\pi^*$ than $\pi_{(c,h+d)}$. We say that the advantage $\pi_{(c+d,h)}$ gets from the extra information is $\alpha_f$. More formally, we say that

$$0 \leq \alpha_f = \min_{\pi_{(c,h+d)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h+d)}, \pi^*)) \Big| C \right] - \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] \tag{10}$$

where $C = \{s_{t-h-d-c:t-h-d}, a_{t-h-d:t-1}\}$ and $G = \{a_t, s_{t-k:t-h-d-c-1}, s_{t-h-d+1:t}, z_{t-k:t}\}$. In particular, $\alpha_f = 0$ when $\pi_{(c,h+d)}$ can infer $s_{t-h-d+1:t-h}$ perfectly *i.e.* when the environment is completely static with $P_f = 1$. This makes sense–in the static environment, observing these states does not provide any advantage since the optimal $\pi_{(c,h+d)}$ can infer these states anyway using the actions taken at those time steps.

Note how this formal definition has some difference from the informal one. In particular, $\pi_h$ only observes $s_{t-h-c:t-h}$, not necessarily all of $s_{t-h-d+1:t-h}$. So, some of the value of $\alpha_f$ can be informally called the "advantage" $\pi_h$ gets over $\pi_{h+d}$, but not necessarily all of it. Nevertheless, our proofs will be using this formal definition.

**Definition (Maximum Inference Disadvantage ($\epsilon_b$)).** Consider the maximum divergence that can be accumulated by the $(c, h)$-model from not knowing the past states $s_{t-h-d-c:t-h-c-1}$ and let that be $\epsilon_b$. More formally, we say that, for fixed $C$ from Proposition 1, any state in $\mathcal{S}^-$ and any $z_{t-k:t}$ and any $\hat{s}_{t-h-d-c:t-h-c-1} \neq s_{t-h-d-c:t-h-c-1}$:

$$\mathcal{L}(\pi_{(c+d,h)}(a_t | \hat{s}_{t-h-d-c:t-h-c-1} \neq s_{t-h-d-c:t-h-c-1}, s_{t-h-c:t-h}, a_{t-h:t-1}), \pi^*) \leq \epsilon_b. \tag{11}$$

Here, $\pi_{(c+d,h)} \coloneqq \arg\min_{\pi_{(c+d,h)}} \mathbb{E}_G[\mathcal{L}(\pi_{(c+d,h)}, \pi^*)|C]$ is the optimal $(c + d, h)$-policy.

The intuitive relationship between $\alpha_f$ and $\epsilon_f$ (and the same for $\alpha_b$ and $\epsilon_b$) holds:

**Proposition 8.** $\alpha_f \leq \epsilon_f$ *and* $\alpha_b \leq \epsilon_b$.

*Proof.* We prove the first inequality; the second can be proven in the same manner. We use Assumption 2 to write $\pi_{(c,h+d)} = \mathbb{E}_{s_{t-h-d+1:t-h} \sim P} \left[ \pi_{(c+d,h)} \mid s_{t-h-d-c:t-h-d}, a_{t-h-d:t-1} \right]$ where $P$ is the environment's transition dynamics. Let

$$P_{\text{correct inference}} = P(\hat{s}_{t-h-d+1:t-h} = s_{t-h-d+1:t-h} | s_{t-h-d-c:t-h-d}, a_{t-h-d:t-1})$$

and

$$P_{\text{incorrect inference}} = P(\hat{s}_{t-h-d+1:t-h} \neq s_{t-h-d+1:t-h} | s_{t-h-d-c:t-h-d}, a_{t-h-d:t-1}).$$

Then,

$$\begin{aligned}
\alpha_f &= \min_{\pi_{(c,h+d)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h+d)}, \pi^*) \Big| C \right] - \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] \\
&= \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(P_{\text{correct inference}} \pi_{(c+d,h)} + P_{\text{incorrect inference}} \pi_{(c+d,h)}, \pi^*) \Big| C \right] \\
&\quad - \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] \\
&\leq \min_{\pi_{(c+d,h)}} \{ \mathbb{E}_G \left[ P_{\text{incorrect inference}} \mathcal{L}(\pi_{(c+d,h)}(\text{conditioning on incorrect inference}), \pi^*) \Big| C \right] \\
&\quad + \mathbb{E}_G \left[ P_{\text{correct inference}} \mathcal{L}(\pi_{(c+d,h)}(\text{conditioning on correct inference}), \pi^*) \Big| C \right] \} \\
&\quad - \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] \hspace{3cm} \text{(Convexity)} \\
&\leq \mathbb{E}_G \left[ P_{\text{incorrect inference}} \mathcal{L}(\hat{\pi}^*_{(c+d,h)}(\text{conditioning on incorrect inference}), \pi^*) \Big| C \right] \\
&\quad + \mathbb{E}_G \left[ \mathcal{L}(\pi^*_{(c+d,h)}, \pi^*) \Big| C \right] - \mathbb{E}_G \left[ \mathcal{L}(\pi^*_{(c+d,h)}, \pi^*) \Big| C \right] \\
&\hspace{3cm} \text{(Bounding probabilities by 1 and Lemma 5)} \\
&\leq \mathbb{E}_G \left[ P_{\text{incorrect inference}} \epsilon_f \mid C \right] \\
&\leq \epsilon_f
\end{aligned}$$

Here, $\pi^*_{(c+d,h)} := \arg\min_{\pi_{(c+d,h)}} \mathbb{E}_G\left[\mathcal{L}(\pi_{(c+d,h)}, \pi^*)\Big|C\right]$. $\qquad\square$

**Definition (Forward and Backward Inference).** For a fixed time step $t$ and $C$ (as in §3.2), consider the time steps $\{t - h - d + 1 : t - h\}$. Define

$$P_f(t') := P(S_{t'} = g_{t'}|S_{t'-1} = g_{t'-1}, A_{t'-1} = a_{t'-1})$$

for all $t' \in [t - h - d + 1 : t - h]$ with $g_{t'}, g_{t'-1}, a_{t'-1}$ being the ground truth states and action in the deterministic environment. We assume that $P_f(t') = 1$ in a deterministic environment. In a stochastic environment, $P_f(t') < 1$ for all $t'$ and as the stochasticity increases, these values decrease and approach 0. Then, we define

$$P_f := \sup\{P_f(t')|t' \in [t - h - d + 1 : t - h]\}.$$

Now, consider the time steps $\{t - h - d - c : t - h - c - 1\}$. Define

$$P_b(t') := P(S_{t'} = g_{t'}|S_{t'+1} = g_{t'+1})$$

for all $t' \in [t-h-d-c : t-h-c-1]$. Since this is not conditioned on any action and as conditioning on $a_{t'}$ reduces entropy, we assume that $P_b(t') < 1$ for all $t' \in [t - h - d - c : t - h - c - 1]$ in all environments. As stochasticity increases, $P_b(t')$ decreases and approaches 0. Then, define

$$P_b := \sup\{P_b(t')|t' \in [t - h - d - c : t - h - c - 1]\}.$$

## C.2 Discussion on Assumption 2

Before we prove the next result, we briefly discuss Assumption 2. Note that, using law of total probability, we can already write:

$$
\pi_{(c,h)}(a_t|s_{t-h-c:t-h}, a_{t-h:t-1})
$$
$$
= \sum_{\substack{s_{t-k:t-h-c-1}, \\ s_{t-h+1:t}}} \hat{P}(s_{t-k:t-h-c-1}, s_{t-h+1:t}|s_{t-h-c:t-h}, a_{t-h:t-1})\pi_{(k,1)}(a_t|s_{t-k:t}, a_{t-h:t-1})
$$

$$\text{(Law of Total Probability)}$$

$$
= \mathbb{E}_{s_{t-k:t-h-c-1}, s_{t-h+1:t} \sim \hat{P}}\left[\pi_{(k,1)}(a_t|s_{t-k:t})|s_{t-h-c:t-h}, a_{t-h:t-1}\right].
$$

Here $\hat{P}$ is the policy's *learned* environment dynamics. Assumption 2 allows us to replace $\hat{P}$ with the true environment dynamics $P$. In other words, we assume that an optimal policy has already learned these distributions as optimally as possible. Given we are talking about the optimal policy trained with infinite data, this is a reasonable assumption. Note that this does not make inference trivial - the policy learns the distribution but, given the distribution has non-zero entropy, the policy can still infer the wrong state.

Thus, using assumption 2, we can write the optimal policy as

$$
\pi_{(c,h)} = \mathbb{E}_{\hat{s}_{t-k:t-h-c-1}, \hat{s}_{t-h+1:t} \sim P}\left[\pi_{(k,1)}|s_{t-h-c:t-h}, a_{t-h:t-1}\right]
$$

where $P(S_{t+1} = s_{t+1}|S_t = s_t, a_t)$ is the environment's transition dynamics and $\pi_{(k,1)}$ is the distribution of $a_t$ under a $(k, 1)$-model.

## C.3 Consistency-Reactivity Inequalities

Now we prove the Consistency-Reactivity Inequalities. We prove the upper and lower bound separately:

**Proposition 1** (Consistency-Reactivity Inequalities - Upper Bound). Let $\mathcal{L}$ be a non-linear, convex loss function. Let $\mathcal{S}^+ \subset \{s_{t-k:t}\}$ be the states both models observe and let $\mathcal{S}^- := \{s_{t-k:t}\} \setminus \mathcal{S}^+$. Let $C := \{a_{t-h-d:t-1}\} \cup \mathcal{S}^+$, $G := \{a_t, z_{t-k:t}\} \cup \mathcal{S}^-$. Then, we can bound the expected loss of the $(c, h + d)$-policy and the $(c, h)$-policy as:

$$
\min_{\pi_{h+d}} \mathbb{E}_G\left[\mathcal{L}(\pi_{h+d}, \pi^*)|C\right] \leq \min_{\pi_h} \mathbb{E}_G\left[\mathcal{L}(\pi_h, \pi^*)|C\right] - \alpha_b + \epsilon_f(1 - P_f^{2d}).
$$

*Proof.* For ease of notation, we will write $x_{a:}^b$ to mean $x_{a:b}$. Additionally, for greater clarity, we will explicitly include the context length of each model, so $\pi_{(c,h)} = \pi_h$ and $\pi_{(c,h+d)} = \pi_{h+d}$. We start

by writing, using Assumption 1,

$$\pi_{(c,h+d)}(a_t|s_{t-h-d-c:t-h-d}, a_{t-h-d:t-1})$$

$$= \pi_{(c,h+d)}(a_t|s^{t-h-d}_{t-h-d-c:}, a^{t-1}_{t-h-d:})$$

$$= \mathbb{E}_{\hat{s}_{t-h-d+1:t-h}}\left[\pi_{(c+d,h)}(a_t|s^{t-h-d}_{t-h-d-c:}, \hat{s}^{t-h}_{t-h-d+1:}, a^{t-1}_{t-h-d:})\Big|s^{t-h-d}_{t-h-d-c:}, a^{t-1}_{t-h-d:}\right].$$

Using this, we expand the left hand side of our inequality:

$$\min_{\pi_{(c,h+d)}} \mathbb{E}_G\left[\mathcal{L}(\pi_{(c,h+d)}, \pi^*)|C\right]$$

$$= \min_{\pi_{(c+d,h)}} \mathbb{E}_G\left[\mathcal{L}(\mathbb{E}_{\hat{s}_{t-h-d+1:t-h}}\left[\pi_{(c+d,h)}\Big|C\right], \pi^*)|C\right]$$

$$= \min_{\pi_{(c+d,h)}} \mathbb{E}_G\left[\mathcal{L}(P(\hat{s}^{t-h}_{t-h-d+1:}|s^{t-h-d}_{t-h-d-c:}, a^{t-h-1}_{t-h-d:})\pi_{(c+d,h)}(a_t|\cdots, g^{t-h}_{t-h-d+1:})+\right.$$

$$\left.\sum_{\substack{\hat{s}_{t-h-d+1:t-h} \\ \text{not all } g_{t'}}} P(\hat{s}^{t-h}_{t-h-d+1:}|s_{t-h-d}, a^{t-h-1}_{t-h-d:})\ \pi_{(c+d,h)}(a_t|\cdots, \hat{s}^{t-h}_{t-h-d+1:}), \pi^*)|C\right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G\left[\mathcal{L}(P^d_f\pi_{(c+d,h)}(a_t|\cdots, g^{t-h}_{t-h-d+1:})+\right.$$

$$\left.\sum_{\substack{\hat{s}_{t-h-d+1:t-h} \\ \text{not all } g_{t'}}} P(\hat{s}^{t-h}_{t-h-d+1:}|s_{t-h-d}, a^{t-h-1}_{t-h-d:})\ \pi_{(c+d,h)}(a_t|\cdots, \hat{s}^{t-h}_{t-h-d+1:}), \pi^*)|C\right]$$

where we computed the expectation $\mathbb{E}_{\hat{s}_{t-h-d+1:t-h}}\left[\pi_{(c+d,h)}\Big|C\right]$ by grouping into two terms : one where every $\hat{s}_{t-h-d+1:t-h} = g_{t-h-d+1:t-h}$ and one where there is at least one term $\hat{s}_i$ that is not $g_i$. This grouping was done using the definition of noise in our environment. We introduce the following notation here

$$\hat{P}_{\neq g_{t'}} := \sum_{\substack{\hat{s}_{t-h-d+1:t-h} \\ \text{not all } g_{t'}}} P(\hat{s}^{t-h}_{t-h-d+1:}|s_{t-h-d}, a^{t-h-1}_{t-h-d:}).$$

Similarly,

$$P_{\neq g_{t'}} := \sum_{\substack{s_{t-h-d+1:t-h} \\ \text{not all } g_{t'}}} P(s^{t-h}_{t-h-d+1:}|s_{t-h-d}, a^{t-h-1}_{t-h-d:}).$$

With this notation, we continue our expansion:

$$\min_{\pi_{(c,h+d)}} \mathbb{E}_G\left[\mathcal{L}(\pi_{(c,h+d)}, \pi^*)|C\right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G\left[\mathcal{L}(P^d_f\pi_{(c+d,h)}(a_t|\cdots, g^{t-h}_{t-h-d+1:})+\right.$$

$$\left.\sum_{\substack{\hat{s}_{t-h-d+1:t-h} \\ \text{not all } g_{t'}}} P(\hat{s}^{t-h}_{t-h-d+1:}|s_{t-h-d}, a^{t-h-1}_{t-h-d:})\ \pi_{(c+d,h)}(a_t|\cdots, \hat{s}^{t-h}_{t-h-d+1:}), \pi^*)|C\right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G\left[\mathcal{L}(P^d_f\pi_{(c+d,h)}(a_t|\cdots, g^{t-h}_{t-h-d+1:}) + \hat{P}_{\neq g_{t'}}\ \pi_{(c+d,h)}(a_t|\cdots, \hat{s}^{t-h}_{t-h-d+1:}), \pi^*)|C\right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G\left[P^d_f\mathcal{L}(\pi_{(c+d,h)}(a_t|\cdots, g^{t-h}_{t-h-d+1:}), \pi^*)|C\right]$$

$$+ \mathbb{E}_G\left[\hat{P}_{\neq g_{t'}}\mathcal{L}(\pi_{(c+d,h)}(a_t|\cdots, \hat{s}^{t-h}_{t-h-d+1:}), \pi^*)|C\right]$$

where we got the inequality using the fact that $\mathcal{L}$ is a convex function and, thus, convex in each argument. Next, we take the expectation over $s_{t-h-d+1:t-h}$ by grouping the terms into two: one where every $s_{t-h-d+1:t-h} = g_{t-h-d+1:t-h}$ and one where there is at least one term $s_i \neq g_i$. Then, with some suppression of notation in the expression of $\pi_{(c+d,h)}$ and $G' := G \setminus \{a_t, s_{t-h-d+1:t-h}\}$:

$$\min_{\pi_{(c,h+d)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h+d)}, \pi^*) | C \right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_{a_t}$$

$$\left[ P_f^d \ P_f^d \mathbb{E}_{G'} \left[ \mathcal{L}(\pi_{(c+d,h)}(...\hat{s}_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h} \right] \right.$$

$$+ \ P_{\neq g_{t'}} \ P_f^d \mathbb{E}_{G'} \left[ \mathcal{L}(\pi_{(c+d,h)}(...\hat{s}_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h} \right]$$

$$+ \ P_f^d \hat{P}_{\neq g_{t'}} \mathbb{E}_{G'} \left[ \mathcal{L}(\pi_{(c+d,h)}(...\hat{s}_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h} \right]$$

$$+ \ P_{\neq g_{t'}} \hat{P}_{\neq g_{t'}} \mathbb{E}_{G'} \left[ \mathcal{L}(\pi_{(c+d,h)}(..., \hat{s}_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h} \right]$$

Now, we group all the terms into two - one representing where the learner's simulation matches the reality and one where it does not. Continuing from where we left off, first, define $P_{\hat{s}=s}^f :=$ $P(s_{t-h-d+1:}^{t-h} | s_{t-h-d}, a_{t-h-d:}^{t-h-1})$

$$\min_{\pi_{(c,h+d)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h+d)}, \pi^*) | C \right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_{a_t}$$

$$\left[ P_f^d \ P_f^d \mathbb{E} \left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h} = s_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h} \right] \right.$$

$$+ \ P_{\neq g_{t'}} \ P_{\hat{s}=s}^f \mathbb{E} \left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d+1:}^{t-h} = s_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h} \right]$$

$$+ \ P_f^d \hat{P}_{\neq g_{t'}} \mathbb{E} \left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h} = s_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h} \right]$$

$$+ \ P_{\neq g_{t'}} P_f^d \mathbb{E} \left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d+1:}^{t-h} = g_{t-h-d+1:}^{t-h} \neq s_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h} \right]$$

$$+ \ P_{\neq g_{t'}} P(\hat{s}_{t-h-d+1:}^{t-h} \neq s_{t'} | s_{t-h-d}, a_{t-h-d:}^{t-h-1})$$

$$\mathbb{E} \left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d+1:}^{t-h} \neq s_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h}), \pi^*) \Big| ..., s_{t-h-d+1:}^{t-h} \neq g_{t-h-d+1:}^{t-h} \right] | C, a_t \left] \right] | C \right]$$

For the match terms, we use the fact that $P_f^d \leq 1$ and $P_{\hat{s}=s}^f = P(\hat{s}_{t-h-d+1:}^{t-h} = s_{t-h-d+1:}^{t-h} | s_{t-h-d}, a_{t-h-d:}^{t-h-1}) \leq 1$. For the mismatch terms, we use the definition of $\epsilon_f$ and Lemma 5. Then, we continue:

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] \qquad \text{(Simulation matches reality)}$$

$$+ P_{\neq g_{t'}} \left[ P_f^d \epsilon_f + \hat{P}_{\neq g_{t'}, s_{t'}} \epsilon_f \right] + P_f^d \left[ \hat{P}_{\neq g_{t'}} \epsilon_f \right]. \qquad \text{(Simulation does not match reality)}$$

We simplify the mismatch terms further:

$$\leq P_{\neq g_{t'}} \left[ P_f^d \epsilon_f + (1 - P_f^d) \epsilon_f \right] + P_f^d \left[ \hat{P}_{\neq g_{t'}} \epsilon_f \right]$$

$$= P_{\neq g_{t'}} \epsilon_f + P_f^d \hat{P}_{\neq g_{t'}} \epsilon_f$$

$$= P_{\neq g_{t'}} \epsilon_f + P_f^d \left[ (1 - P_f^d) \right] \epsilon_f$$

$$= (1 - P_f^d) \epsilon_f + P_f^d \left[ (1 - P_f^d) \right] \epsilon_f$$

$$= \epsilon_f \cdot \left[ 1 - P_f^d + P_f^d - P_f^{2d} \right]$$

$$= \epsilon_f \cdot \left[ 1 - P_f^{2d} \right].$$

Next, we simplify the match terms by using the definition of $\alpha_b$:

$$\min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] = \min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*) | C \right] - \alpha_b.$$

Substituting these two terms back in, we conclude. $\qquad \square$

Now, we prove the lower bound:

**Proposition 1** (Consistency-Reactivity Inequalities - Lower Bound). Let $\mathcal{L}$ be a non-linear, convex loss function. Let $\mathcal{S}^+ \subset \{s_{t-k:t}\}$ be the states both models observe and let $\mathcal{S}^- := \{s_{t-k:t}\} \setminus \mathcal{S}^+$. Let $C := \{a_{t-h-d:t-1}\} \cup \mathcal{S}^+$, $G := \{a_t, z_{t-k:t}\} \cup \mathcal{S}^-$. Then, we can bound the expected loss of the $(c, h+d)$-policy and the $(c, h)$-policy as:

$$\min_{\pi_h} \mathbb{E}_G \left[ \mathcal{L}(\pi_h, \pi^*)|C \right] \leq \min_{\pi_{h+d}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{h+d}, \pi^*)|C \right] - \alpha_f + \epsilon_b(1 - P_b^{2d}).$$

*Proof.* We proceed in a manner similar to the proof of the upper bound. For ease of notation, we will write $x_{a:}^b$ to mean $x_{a:b}$. Additionally, for greater clarity, we will explicitly include the context length of each model, so $\pi_{(c,h)} = \pi_h$ and $\pi_{(c,h+d)} = \pi_{h+d}$. We start by writing, using Assumption 1,

$$\min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*) \mid C \right]$$

$$= \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(P(g_{t-h-d-c:}^{t-h-c-1}|s_{t-h-c})\pi_{(c+d,h)} + \sum_{\substack{\hat{s}_{t-h-d-c:}^{t-h-c-1}, \\ \text{not all } g_{t'}}} P(\hat{s}_{t-h-d-c:}^{t-h-c-1}|s_{t-h-c})\pi_{(c+d,h)}^t, \pi^*) \Big| C \right].$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(P_b^d \pi_{(c+d,h)} + \sum_{\substack{\hat{s}_{t-h-d-c:}^{t-h-c-1}, \\ \text{not all } g_{t'}}} P(\hat{s}_{t-h-d-c:}^{t-h-c-1}|s_{t-h-c})\pi_{(c+d,h)}^t, \pi^*) \Big| C \right].$$

We introduce the following notation here

$$\hat{P}_{\neq g_{t'}} := \sum_{\substack{\hat{s}_{t-h-d-c:t-h-c-1} \\ \text{not all } g_{t'}}} P(\hat{s}_{t-h-d-c:}^{t-h-c-1}|s_{t-h-c}).$$

Similarly,

$$P_{\neq g_{t'}} := \sum_{\substack{s_{t-h-d-c:t-h-c-1} \\ \text{not all } g_{t'}}} P(s_{t-h-d-c:}^{t-h-c-1}|s_{t-h-c}).$$

With this notation, we continue our expansion:

$$\min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*)|C \right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G \Big[ \mathcal{L}(P_b^d \pi_{(c+d,h)}(a_t|s_{t-h-c:}^{t-h}, g_{t-h-d-c:}^{t-h-c-1}, a_{t-h:}^{t-1}) +$$

$$\hat{P}_{\neq g_{t'}} \pi_{(c+d,h)}(a_t|s_{t-h-c:}^{t-h}, \hat{s}_{t-h-d-c:}^{t-h-c-1}, a_{t-h:}^{t-1}), \pi^*) \mid C \Big]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G \Big[ P_b^d \mathcal{L}(\pi_{(c+d,h)}(a_t|..., g_{t-h-d-c:}^{t-h-c-1}), \pi^*) \mid C \Big] +$$

$$\mathbb{E}_G \Big[ \hat{P}_{\neq g_{t'}} \mathcal{L}( \pi_{(c+d,h)}(a_t|..., \hat{s}_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1}), \pi^*) \mid C \Big]$$

where we got the inequality using the fact that $\mathcal{L}$ is a convex function. Next, we take the expectation over $s_{t-h-d-c:t-h-c-1}$ by grouping the terms into two: one where every $s_{t-h-d-c:t-h-c-1} = g_{t-h-d-c:t-h-c-1}$ and one where there is at least one term $s_i \neq g_i$. Then, again suppressing some terms inside the expression of $\pi_{(c+d,h)}$:

$$\min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*) | C \right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_{a_t}$$

$$\left[ P_b^d \ P_b^d \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(..., \hat{s}_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1} = s_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1} \right] \right.$$

$$+ P_{\neq g_{t'}} \ P_b^d \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(..., \hat{s}_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1} \neq s_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1} \right]$$

$$+ P_b^d \hat{P}_{\neq g_{t'}} \ \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(..., \hat{s}_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1} = s_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1} \right]$$

$$+ \left. P_{\neq g_{t'}} \hat{P}_{\neq g_{t'}} \ \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(..., \hat{s}_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1} \right] \Big| C, a_t \right] | C \right]$$

Now, we group all the terms into two - one representing where the learner's simulation matches the reality and one where it does not. Continuing from where we left off and defining $P_{\hat{s}=s}^b :=$ $P(\hat{s}_{t-h-d-c:}^{t-h-c-1} = s_{t-h-d-c:}^{t-h-c-1} | s_{t-h-c})$:

$$\min_{\pi_{(c,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}, \pi^*) | C \right]$$

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_{a_t}$$

$$\left[ P_b^d \ P_b^d \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1} \right] \right.$$

$$+ P_{\neq g_{t'}} \ P_{\hat{s}=s}^b \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d-c:}^{t-h-c-1} = s_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1} \right]$$

$$+ P_b^d \hat{P}_{\neq g_{t'}} \ \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., s_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1} \right]$$

$$+ \left. P_{\neq g_{t'}} P_b^{2d} \ \mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d-c:}^{t-h-c-1} = g_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1} \right] \Big| C, a_t \right] | C \right]$$

$$+ P_{\neq g_{t'}} P(\hat{s}_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1}, s_{t-h-d-c:}^{t-h-c-1} \mid s_{t-h-c})$$

$$\mathbb{E}\left[ \mathcal{L}(\pi_{(c+d,h)}(a_t | ..., \hat{s}_{t-h-d-c:}^{t-h-c-1} \neq s_{t-h-d-c:}^{t-h-c-1}), \pi^*) \Big| ..., s_{t-h-d-c:}^{t-h-c-1} \neq g_{t-h-d-c:}^{t-h-c-1} \right] \Big| C, a_t \right] | C \right]$$

For the match terms, we use the fact that $P_b^d \leq 1$ and $P(\hat{s}_{t-h-d-c:}^{t-h-c-1} = s_{t-h-d-c:}^{t-h-c-1} | s_{t-h-c:}^{t-h}) \leq 1$. For the mismatch terms, we use the definition of $\epsilon_b$ and Lemma 5. Then, we continue:

$$\leq \min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right] \qquad \text{(Simulation matches reality)}$$

$$+ P_{\neq g_{t'}} \left[ P_b^d \epsilon_b + \hat{P}_{\neq g_{t'}, s_{t'}} \epsilon_b \right] + P_b^d \left[ \hat{P}_{\neq g_{t'}} \epsilon_b \right]. \qquad \text{(Simulation does not match reality)}$$

We simplify the mismatch terms further:

$$\leq P_{\neq g_{t'}} \left[ P_b^d \epsilon_b + (1 - P_b^d)\epsilon_b \right] + P_b^d \left[ \hat{P}_{\neq g_{t'}} \epsilon_b \right]$$

$$= P_{\neq g_{t'}} \epsilon_b + P_b^d \hat{P}_{\neq g_{t'}} \epsilon_b$$

$$= P_{\neq g_{t'}} \epsilon_b + P_b^d \left[ (1 - P_b^d) \right] \epsilon_b$$

$$= (1 - P_b^d)\epsilon_b + P_b^d \left[ (1 - P_b^d) \right] \epsilon_b$$

$$= \epsilon_b \cdot \left[ 1 - P_b^d + P_b^d - P_b^{2d} \right]$$

$$= \epsilon_b \cdot \left[ 1 - P_b^{2d} \right].$$

Next, we simplify the match terms by using the definition of $\alpha_f$ which follows from Proposition 7 in a manner similar to the definition of $\alpha_b$:

$$\min_{\pi_{(c+d,h)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c+d,h)}, \pi^*) \Big| C \right]$$

$$= \min_{\pi_{(c,h+d)}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h+d)}, \pi^*) | C \right] - \alpha_f. \qquad \text{(Proposition 7)}$$

We substitute these terms back in to get the desired bound. □

## C.4 DISCUSSION ON ASSUMPTION 1

We assumed that $c + h < k$ so that the larger action chunk model can condition on more states from the distant past that are temporally correlated with $a_t$. In the case where $c + h \geq k$, the larger action chunk model will not get any advantage (so, $\alpha = 0$) and can only suffer from not having observed the recent past states. However, this is a very unlikely scenario because we expect human demonstrators to have a large memory horizon *i.e.* $k$ is expected to be large.

## C.5 PROOF OF COROLLARY 2 AND COROLLARY 3

Now, we prove Corollary 2 as a direct consequence of the Consistency-Reactivity Inequalities. In a near-deterministic environment, $P_f$ is close to 1. This is because, conditioned on the state and action at time $t' - 1$, we can confidently infer the state visited at time $t'$ as the environment lacks noise.

**Corollary 2 (Restated)** In a near-deterministic environment, if $a_t$ is temporally dependent on at least one state in $\{s_{t-h-c-d:t-h-c-1}\}$ and $\epsilon_f$ is finite,

$$\min_{\pi_{h+d}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{h+d}, \pi^*) | C \right] < \min_{\pi_h} \mathbb{E}_G \left[ \mathcal{L}(\pi_h, \pi^*) | C \right]$$

*Proof.* This follows from Proposition 1 by taking $P_f \approx 1$, $\epsilon_f$ not large and $\alpha_b > 0$ (since $a_t$ is temporally dependent on at least one state in $\{s_{t-h-c-d:t-h-c-1}\}$ and the states in $\mathcal{T}_b$ cannot be deterministically inferred from the context of $\pi_h$). □

Lastly, we prove Corollary 3. We assume that, in a highly stochastic environment, $P_b$ is small. This is because, for any time step $t'$, the agent could reach $s_{t'+1}$ from many different states at time $t'$ due to the noise in the environment.

**Corollary 3 (Restated)** In a highly stochastic environment, if temporal dependency decreases over the number of time steps *i.e.* $\alpha_f > \epsilon_b$, then

$$\min_{\pi_{h+d}} \mathbb{E}_G \left[ \mathcal{L}(\pi_{h+d}, \pi^*) | C \right] > \min_{\pi_h} \mathbb{E}_G \left[ \mathcal{L}(\pi_h, \pi^*) | C \right]$$

*Proof.* Starting from Proposition 1, with $P_b$ small (since the environment is highly stochastic), we get

$$\min_{\pi_{(c,h)}^t} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}^t - \pi_E) | C \right] \leq \min_{\pi_{(c,h+d)}^t} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h+d)}^t - \pi_E^t) | C \right] - \alpha_f + \epsilon_b.$$

Since temporal dependency reduces as number of steps grow, $a_t$ is more temporally dependent on the recent past states than on the distant past, so $\alpha_f > \epsilon_b$. With this observation, we get:

$$\min_{\pi_{(c,h)}^t} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h)}^t - \pi_E) | C \right] < \min_{\pi_{(c,h+d)}^t} \mathbb{E}_G \left[ \mathcal{L}(\pi_{(c,h+d)}^t - \pi_E^t) | C \right].$$

□

## C.6 CLOSED-LOOP VERSUS OPEN-LOOP IN HIGHLY STOCHASTIC AND NEAR-DETERMINISTIC ENVIRONMENTS

The Consistency-Reactivity Inequalities allow us to make an even stronger statement when we compare strictly closed-loop policies with open-loop ones. Consider the same set-up as before with $h = 0$. Thus, $\pi_{(c,0)}$ represents a closed-loop policy whereas $\pi_{(c,d)}$ represents an open-loop one. We can compare these policies' divergences with the expert across the entire trajectory in the limiting cases of the environment stochasticity.

**Corollary 9.** *In a highly stochastic environment, if temporal dependency decreases such that $\alpha_f > \epsilon_b$ at all time steps, then divergence between the closed-loop policy over the full trajectory is lower than that between the open-loop policy and the expert. In a near deterministic environment, if there is at least one time step $t$ such that $a_t$ depends on some state in $s_{t-d-c:t-c-1}$, then the divergence between the closed-loop policy over the full trajectory is greater than that between the open-loop policy and the expert.*

*Proof.* At any arbitrary time step $t$, the chunks of the two policies may be aligned in one of two ways:

Case 1: $\pi_{(c,1)}$ is executing $a_t$ as the first action in its action chunk and $\pi_{(c,1+d)}$ is also executing $a_t$ as the first action in its action chunk.

Case 2: $\pi_{(c,1)}$ is executing $a_t$ as the first action in its action chunk and $\pi_{(c,1+d)}$ is also executing $a_t$ as the $k$-th action, where $k \in (1, 1+d]$ in its action chunk.

Using the Consistency-Reactivity Inequalities, in Case 1, both policies have equal divergence.

However, in case 2, using Corollary 3, we know that the closed-loop policy will outperform the open-loop one in the first setting of the statement and open-loop will outperform in the second. From this, we can conclude the divergence across the full trajectory. $\qquad\square$