

Directly Forecasting Belief for Reinforcement Learning with Delays

Qingyuan Wu* ¹ Yuhui Wang* ² Simon Sinong Zhan* ³
Yixuan Wang ³ Chung-Wei Lin ⁴ Chen Lv ⁵ Qi Zhu ³ Jürgen Schmidhuber ^{2, 6} Chao Huang ¹

¹University of Southampton

²Gen AI, KAUST

³Northwestern University

⁴National Taiwan University

⁵Nanyang Technological University

⁶The Swiss AI Lab IDSIA/USI/SUPSI

*Equal Contribution



TL;DR

We present the Directly Forecasting Belief Transformer (DFBT) for delayed RL, which can effectively reduce the compounding errors and improve performance. Specifically,

- We present DFBT, a novel directly forecasting belief method that effectively addresses compounding errors in recursively generated belief.
- We propose DFBT-SAC, a novel delayed RL method that further improves the learning efficiency via multi-step bootstrapping on the DFBT.
- We theoretically demonstrate that our DFBT significantly reduces compounding errors compared to the existing recursively forecasting belief approach.
- We empirically demonstrate that our DFBT method effectively forecasts state sequences with significantly higher prediction accuracy compared to baselines.
- We empirically show that our DFBT-SAC outperforms SOTAs in terms of sample efficiency and performance on the MuJoCo benchmark.

Background

A delay-free RL problem is formalized as an MDP represented as $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \rho, \gamma \rangle$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the dynamic function, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, ρ is the initial state distribution, and $\gamma \in (0, 1)$ is the discount factor.

A delayed RL problem can be formalized as an augmented MDP. For instance, a delayed RL problem with constant delays Δ represented as $\langle \mathcal{X}, \mathcal{A}, \mathcal{P}_\Delta, \mathcal{R}_\Delta, \rho_\Delta, \gamma \rangle$, where

- Augmented state space $\mathcal{X} := \mathcal{S} \times \mathcal{A}^\Delta$
- Action space \mathcal{A}
- Delayed dynamic $\mathcal{P}_\Delta(x_{t+1}|x_t, a_t) := \mathcal{P}(s_{t-\Delta+1}|s_{t-\Delta}, a_{t-\Delta})\delta_{a_t}(a'_t) \prod_{i=1}^{\Delta-1} \delta_{a_{t-i}}(a'_{t-i})$
- Delayed reward function $\mathcal{R}_\Delta(x_t, a_t) := \mathbb{E}_{s_t \sim b(\cdot|x_t)} [\mathcal{R}(s_t, a_t)]$;
- Initial augmented state distribution $\rho_\Delta = \rho \prod_{i=1}^{\Delta} \delta_{a_{-i}}$
- Discount factor $\gamma \in (0, 1)$
- Belief representation $b : \mathcal{X} \times \mathcal{S} \rightarrow [0, 1]$

Research Problem

Specifically, belief representation is defined as follows:

$$b(s_t|x_t) := \int_{\mathcal{S}^\Delta} \prod_{i=0}^{\Delta-1} \mathcal{P}(s_{t-\Delta+i+1}|s_{t-\Delta+i}, a_{t-\Delta+i}) ds_{t-\Delta+i+1}.$$

The belief representation can be viewed as the recursive forward prediction of the dynamics \mathcal{P} . With the belief representation, the agent can directly learn in the original state space \mathcal{S} .

However, the recursive process is evidently affected by the error accumulation of the approximate dynamic function across Δ steps.

The compounding errors grow exponentially with the delays Δ . This fundamental limitation of such recursive methodology for belief forecasting leads to significant performance degradation, especially in environments with long-delayed signals.

Method

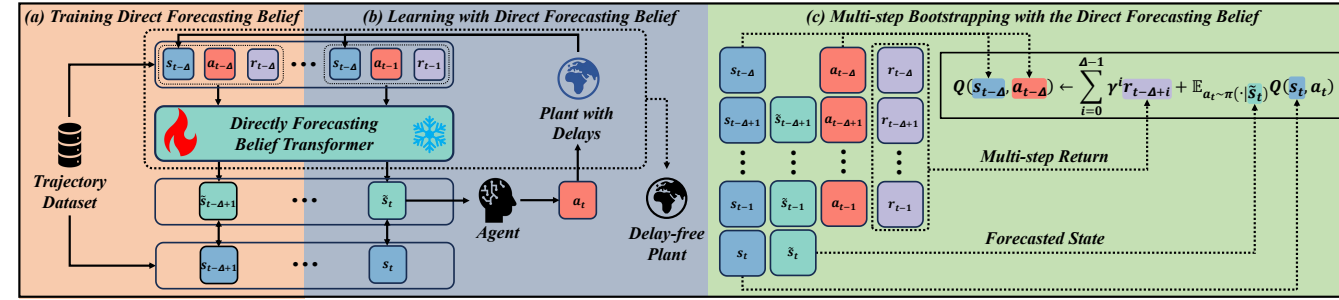


Figure 1. Pipeline of DFBT-SAC.

Given sub-trajectory with Δ timesteps $\{s_{t-\Delta+i}, a_{t-\Delta+i}, r_{t-\Delta+i}\}_{i=0}^{\Delta}$. We reform the representation of the augmented state to Δ tokens for sequence modeling: $x_t^{\text{tokens}} = \{s_{t-\Delta}, a_{t-\Delta+i}, r_{t-\Delta+i}\}_{i=0}^{\Delta-1}$. Then, DFBT predicts the unobserved Δ states $\{s_{t-\Delta+i}\}_{i=1}^{\Delta}$ via autoregressive modeling with loss:

$$\nabla_{\theta} \left[\sum_{i=1}^{\Delta} \left[-\log b_{\theta}^{(i)}(s_{t-\Delta+i}|x_t^{\text{tokens}}) \right] \right], \quad (1)$$

where $b_{\theta}^{(i)}(\cdot|x_t)$ represents the i -th prediction. The critic of DFBT-SAC is multi-step bootstrapped on the states predicted by the DFBT. Specifically, the critic Q_{ψ} parameterized by ψ is updated via:

$$\nabla_{\psi} \left[\frac{1}{2} (Q_{\psi}(s_{t-\Delta}, a_{t-\Delta}) - \mathbb{Y})^2 \right], \quad (2)$$

where N -step ($N \leq \Delta$) temporal difference target \mathbb{Y} is defined as:

$$\mathbb{Y} := \sum_{i=0}^{N-1} \gamma^i r_{t-\Delta+i} + \gamma^N \mathbb{E}_{\substack{a \sim \pi(\cdot|s_{t-\Delta+N}) \\ s_{t-\Delta+N} \sim b_{\theta}^{(N)}(\cdot|x_t^{\text{tokens}})}} [Q(s_{t-\Delta+N}, a) + \log \pi(a|s_{t-\Delta+N})].$$

Main Theoretical Results

Performance Difference of Recursively Forecasting Belief

[Theorem 5.5] For the delay-free policy π and the delayed policy π_{Δ} . Given any $x_t \in \mathcal{X}$, the performance difference $I^{\text{recursive}}(x_t)$ of the recursively forecasting belief b_{θ} can be bounded as follows, respectively. For deterministic delays Δ , we have

$$|I^{\text{recursive}}(x_t)| \leq |I_{\Delta}^{\text{true}}(x_t)| + L_V \underbrace{\frac{1 - L_P^{\Delta}}{1 - L_P} \epsilon_P}_{\text{compounding errors}}.$$

For stochastic delays $\delta \sim d_{\Delta}(\cdot)$, we have

$$|I^{\text{recursive}}(x_t)| \leq \mathbb{E}_{\delta \sim d_{\Delta}(\cdot)} \left[|I_{\delta}^{\text{true}}(x_t)| + L_V \underbrace{\frac{1 - L_P^{\delta}}{1 - L_P} \epsilon_P}_{\text{compounding errors}} \right].$$

Main Experimental Results

Belief Errors Comparison

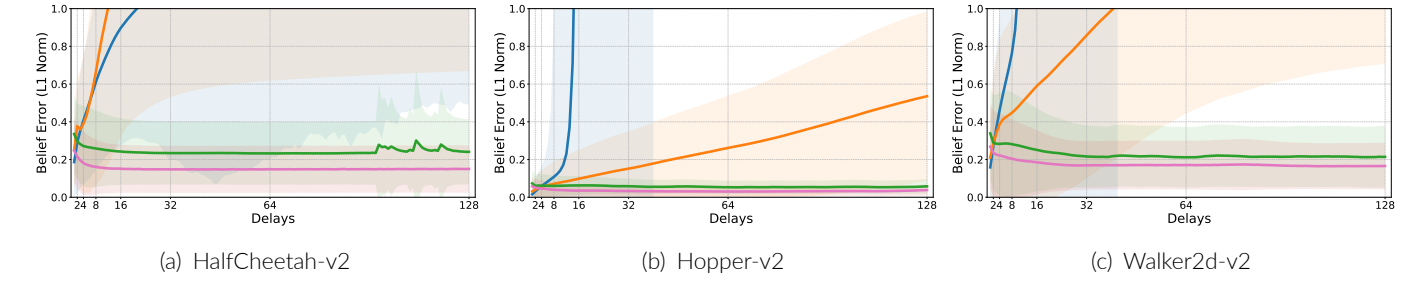


Figure 2. Belief errors comparison.

Performance Comparison

The best performance is underlined, the best belief-based method is in **red**.

Table 1. Performance on MuJoCo with Deterministic Delays.

Task	Delays	Augmentation-based			Belief-based			
		A-SAC	BPQL	ADRL	DATS	D-Dreamer	D-SAC	DFBT-SAC (ours)
HalfCheetah-v2	8	0.10±0.01	0.40±0.04	0.44±0.03	0.08±0.01	0.08±0.01	0.12±0.06	0.35±0.12
	32	0.02±0.02	0.40±0.03	0.26±0.04	0.11±0.04	0.08±0.00	0.08±0.02	0.42±0.03
	128	0.04±0.06	0.08±0.13	0.14±0.02	0.10±0.08	0.15±0.05	0.09±0.04	0.41±0.03
Hopper-v2	8	0.61±0.31	0.87±0.09	0.95±0.16	0.41±0.31	0.11±0.01	0.16±0.05	0.77±0.18
	32	0.11±0.02	<u>0.89±0.14</u>	0.73±0.20	0.07±0.04	0.11±0.05	0.11±0.01	0.68±0.20
	128	0.04±0.01	0.08±0.02	0.07±0.01	0.08±0.01	0.09±0.03	0.06±0.01	0.20±0.03
Walker2d-v2	8	0.44±0.26	<u>1.07±0.02</u>	0.97±0.10	0.13±0.05	0.11±0.06	0.09±0.05	0.99±0.03
	32	0.10±0.02	<u>0.37±0.25</u>	0.16±0.08	0.02±0.03	0.08±0.05	0.08±0.02	0.64±0.10
	128	0.06±0.00	0.07±0.03	0.08±0.01	0.02±0.02	0.08±0.05	0.11±0.06	0.40±0.08

Table 2. Performance on MuJoCo with Stochastic Delays.

Task	Delays	Augmentation-based			Belief-based			
		A-SAC	BPQL	ADRL	DATS	D-Dreamer	D-SAC	DFBT-SAC (ours)
HalfCheetah-v2	$U(1, 8)$	0.09±0.01	0.21±0.07	0.17±0.07	0.09±0.03	0.02±0.01	0.03±0.01	0.37±0.12
	$U(1, 32)$	0.01±0.00	<u>0.33±0.07</u>	0.23±0.02	0.11±0.04	0.02±0.00	0.01±0.01	0.31±0.16
	$U(1, 128)$	0.01±0.01	0.03±0.03	0.15±0.02	0.16±0.03	0.16±0.00	0.02±0.00	0.39±0.04
Hopper-v2	$U(1, 8)$	0.17±0.05	0.20±0.04	0.18±0.04	0.04±0.01	0.07±0.05	0.14±0.04	0.86±0.18
	$U(1, 32)$	0.05±0.01	0.07±0.09	0.05±0.01	0.05±0.01	0.04±0.01	0.03±0.01	0.43±0.21
	$U(1, 128)$	0.03±0.01	0.04±0.01	0.04±0.02	0.05±0.00	0.03±0.01	0.03±0.00	0.14±0.01
Walker2d-v2	$U(1, 8)$	0.36±0.24	0.40±0.32	0.41±0.15	0.07±0.01	0.07±0.05	0.12±0.04	1.11±0.10
	$U(1, 32)$	0.12±0.03	0.16±0.04	0.11±0.05	0.09±0.04	0.12±0.04	0.05±0.02	0.67±0.15
	$U(1, 128)$	0.06±0.01	0.06±0.06	0.04±0.02	0.10±0.04	0.15±0.07	0.03±0.04	0.30±0.13