



Forecasting in Offline Reinforcement Learning for Non-stationary Environments



Suzan Ece Ada
Boğaziçi University
University of Tübingen



Georg Martius
University of Tübingen



Emre Uğur
Boğaziçi University



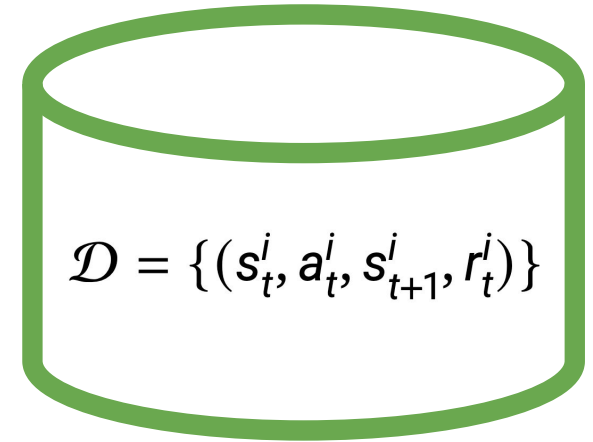
Erhan Öztop
Osaka University
Özyeğin University



Offline Reinforcement Learning (RL)

- Learning a policy from previously collected datasets without real-world interaction.
- The dataset $\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$ consists of MDP tuples generated by a behavior policy $\pi_\beta(\mathbf{a}|\mathbf{s})$.
- Avoids expensive and risky data collection.

Aim: Improve upon the behavior policy used for data collection.



⊘ Assumes stationarity or full observability during training and at test time.

Non-stationary Environments

Each component of the underlying MDP or POMDP can evolve over **time.**

[Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards continual reinforcement learning: A review and perspectives., 2022]

$$(\mathcal{S}(t), \mathcal{A}(t), \mathcal{T}(t), \mathcal{R}(t), \mathbf{x}(t), \mathcal{O}(t))$$

Standard Offline RL Fails with Non-stationary Perturbations

❌ Standard offline RL assume stationarity or full observability during training and at test time.

Robust offline RL addresses settings with **perturbed states at test-time** (e.g., **Gaussian or adversarial noise**). The real world is non-stationary.

How can we handle non-stationary perturbations?

We need test-time adaptation without re-training and presupposing any specific pattern of future non-stationarity during training.



Image generated by Gemini

Forecasting in Offline Reinforcement Learning for Non-stationary Environments

Forecasting in Non-stationary Offline RL (FORL) tackles the foundational challenge of additive episodic biases in the observation function.

🌟 AI Assistants (LLMs)	💊 Healthcare & Finance	🤖 Industrial Robots
Semantic or contextual bias	Data may be withheld (privacy/regulations).	Sensors drift. Robots require daily calibration offsets.

Standard **noise-driven** or **parametric** state-estimation techniques, which typically rely on **smoothly varying or randomly perturbed functions**, cannot reliably identify **persistent, episode-wide offsets** that are **not available after episode terminates**.

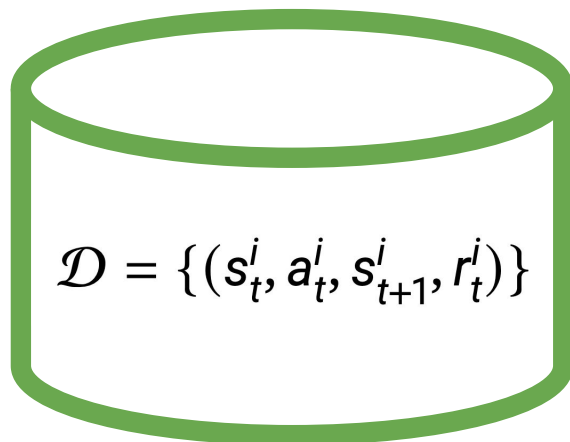
Problem Statement

Training

Offline dataset $\mathcal{D} = \{(s_t^i, a_t^i, s_{t+1}^i, r_t^i)\}$ from a stationary environment.

No knowledge of how offsets evolve in the future during training.

Train an offline RL policy π_ϕ and a FORL *diffusion model* (**FORL-DM**) using \mathcal{D} .



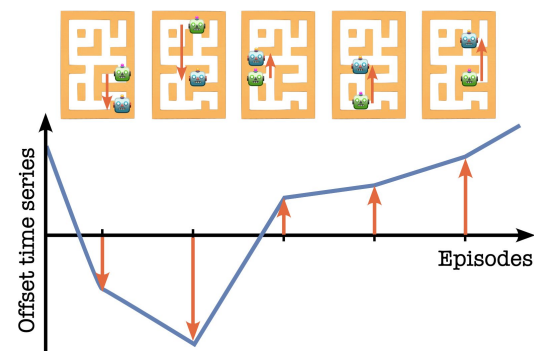
Test Environment

At test time, the agent faces an infinite sequence of POMDPs $\{\hat{\mathcal{M}}_j\}_{j=1}^\infty$. Each POMDP \mathcal{M}_j is described by a 7-tuple [Kaelbling et al., 1998]:

$$\hat{\mathcal{M}}_j = (\mathcal{S}, \mathcal{A}, \mathcal{O}_j, \mathcal{T}, \mathcal{R}, \rho_0, \mathbf{x}_j),$$

where \mathcal{S}, \mathcal{A} , transition function \mathcal{T} and the reward function \mathcal{R} remain identical to the training MDP.

\mathbf{x} is the **deterministic** observation function ($\mathbf{x} : \mathcal{S} \rightarrow \mathcal{O}$) [Bonet, 2012, Khetarpal et al., 2022] that evolves.

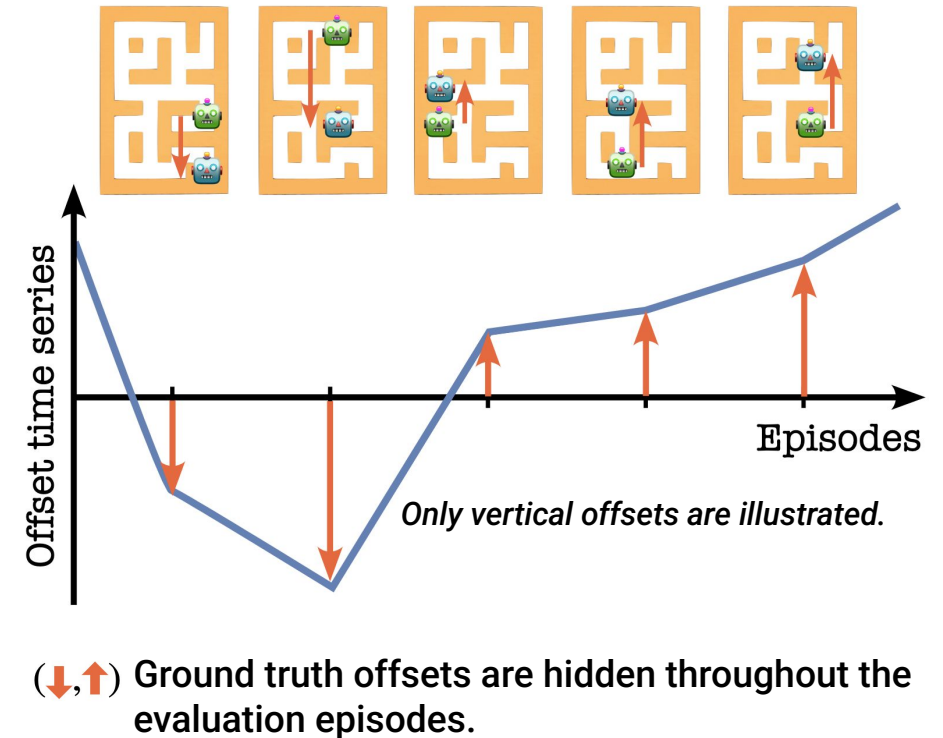


Test Environment: Time-varying Unknown Process

The non-stationarity is driven by a time-varying, unknown process:

- Each **episode** j is characterized by an **unknown** offset b^j added to the true state: $o_t = s_t + b^j$.
- The initial state distribution ρ_0 is uniform, providing no information about the offset.
- The offsets originate from **real-world time-series** $(\dots, b^{j-1}, b^j, b^{j+1}, \dots)$
- The agent only sees the sequence of biased observations $\{o_t\}$.

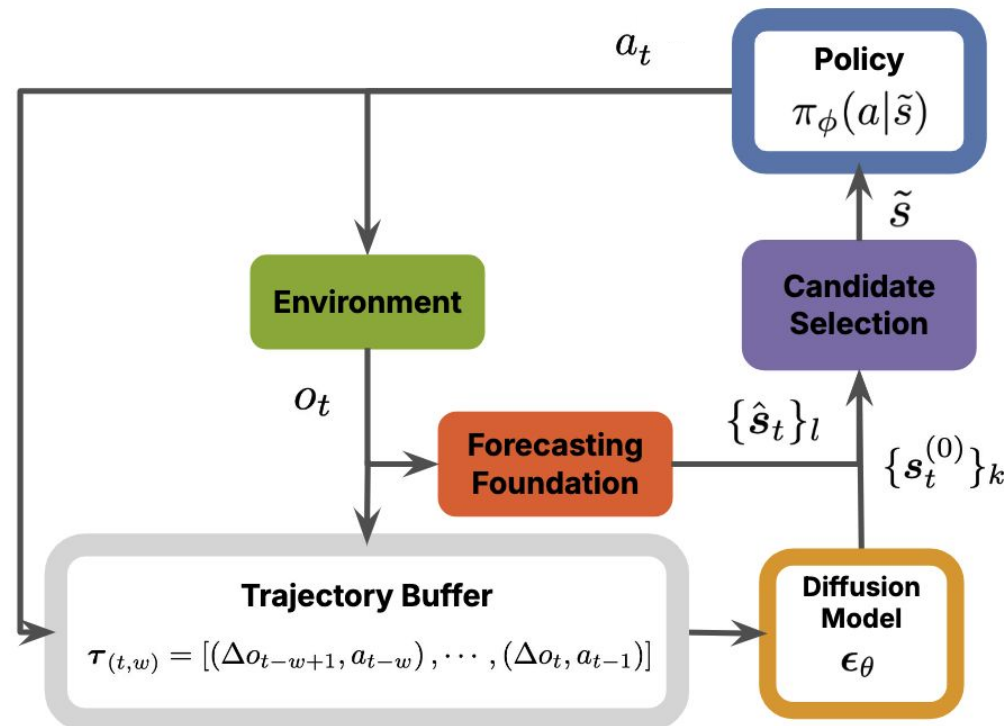
The agent must act effectively from the onset of each episode, despite being "lost."



Overview of FORL Framework at Test-time

Forecast offsets over next P episodes (zero-shot).

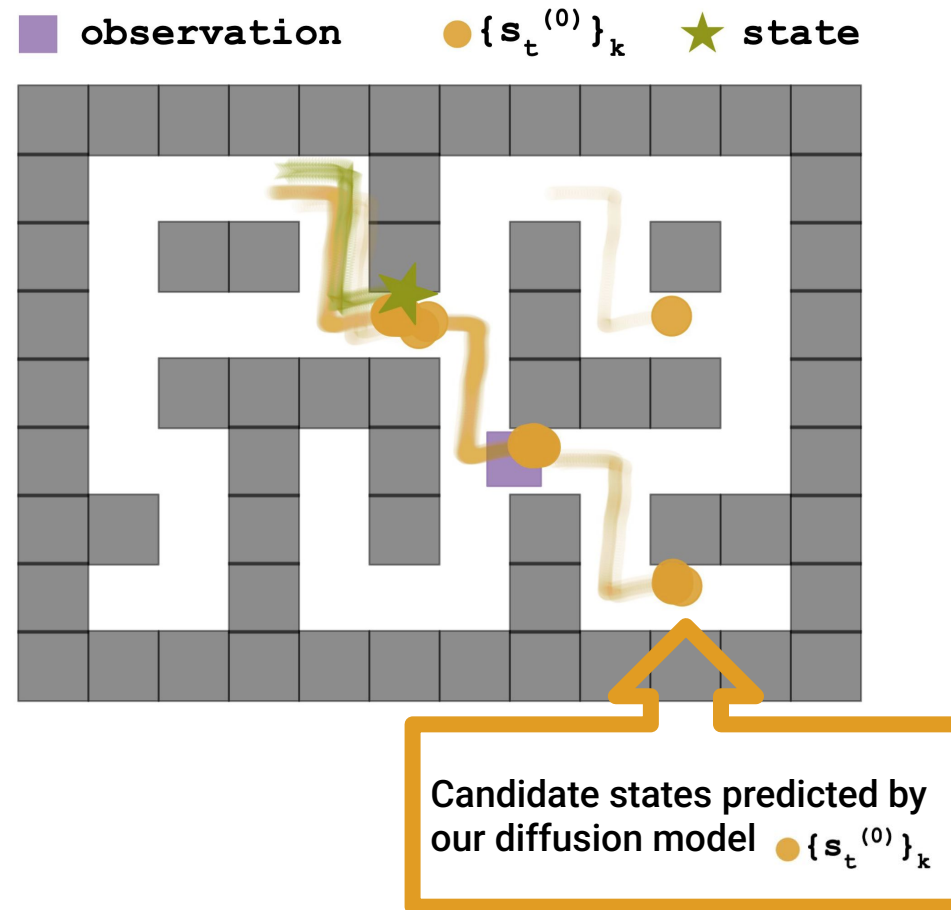
- The observations are processed by both the trajectory buffer and the time-series **forecasting foundation** module.
- Observation changes and actions sampled from the **policy**, $(\Delta o_t, a_t)$ are stored in the trajectory buffer.
- The **diffusion model** generates candidate states $\{s_t^{(0)}\}_k$ conditioned on $\tau_{(t,w)}$.
- The **candidate selection** module then generates the estimated \tilde{s}_t .



Candidate State Generation: Diffusion Model

How can an agent infer its location when the state is subject to a large unknown offset?

- An agent observes a sequence of its own actions and the resulting changes in its (biased) observations.
- This action-effect history, $(\Delta o, a)$, is **invariant** to the offset b^j .



Candidate State Generation: Diffusion Model

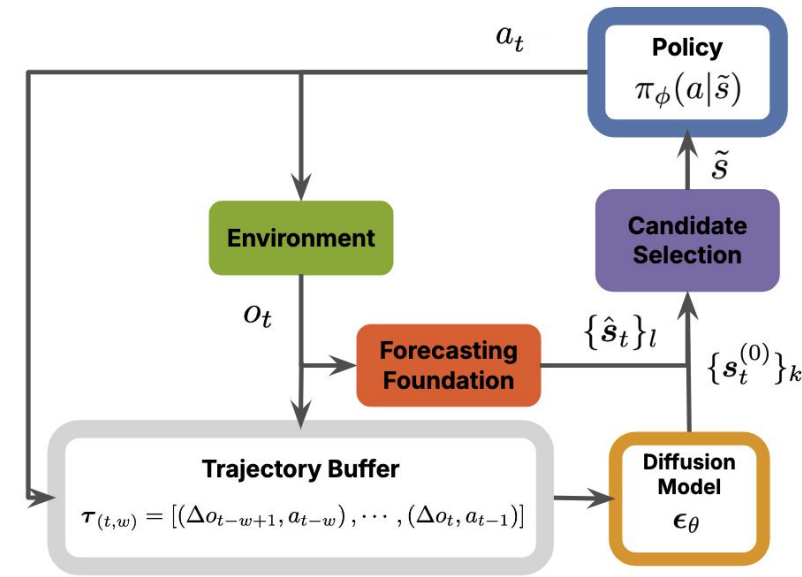
We model the **multi-modal distribution** of plausible states using a conditional diffusion model.

1. Define the **agent's recent history** as a trajectory $\tau_{(t,w)}$ of the past w action-effect Δo_t pairs:

$$\tau_{(t,w)} = [(\Delta o_{t-w+1}, a_{t-w}), \dots, (\Delta o_t, a_{t-1})]$$

2. Train the conditional diffusion model FORL-DM to generate plausible states s_t conditioned on this history using the **stationary offline RL dataset**.

Goal: Learn $p(\mathbf{s}_t \mid \tau_{(t,w)})$



At test time, this model acts as a powerful inference engine, generating a set of k plausible state candidates $\{s_t^{(0)}\}_k$ from the agent's real-time, biased observations and actions.

Forecasting Foundation Model

We use a probabilistic, zero-shot **foundation model** (Lag-Llama [Rasul et al., 2023]) to forecast future offsets based on past history.

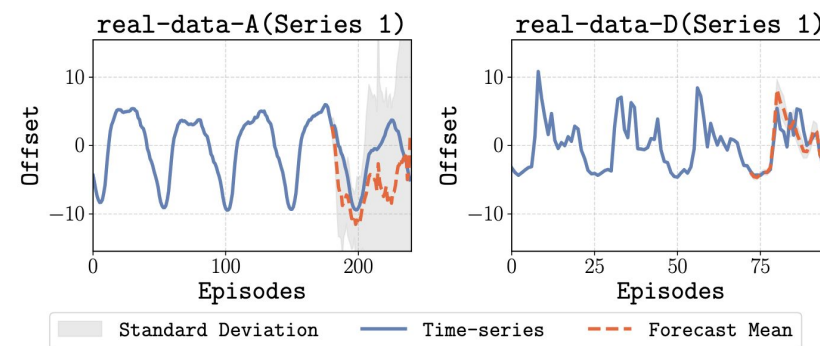
- **Input:** A context window of C past, known offsets:

$$\{b^{j-C}, \dots, b^{j-1}\}$$

- **Output:** A probabilistic forecast of I samples for the next P episodes:

$$\{\hat{b}_I^j, \dots, \hat{b}_I^{j+P-1}\}$$

The model is **univariate**, and each dimension of the offset vector b^j is forecast independently.



Rasul, Kashif, et al. "Lag-Llama: Towards foundation models for time series forecasting."

Candidate Selection: Dimension-wise Closest Match (DCM)

DCM is a **simple** and **light-weight** approach to adaptively fuse samples from 2 sources.

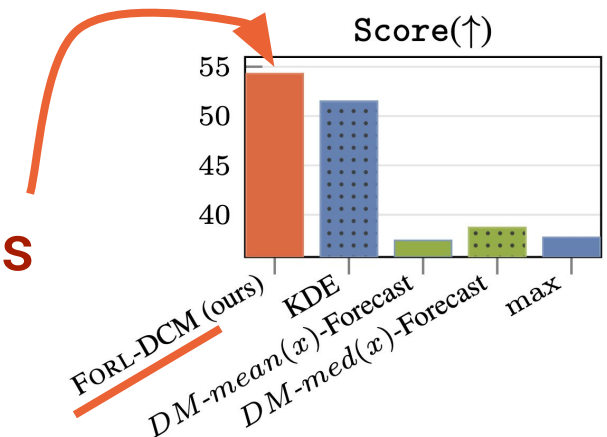
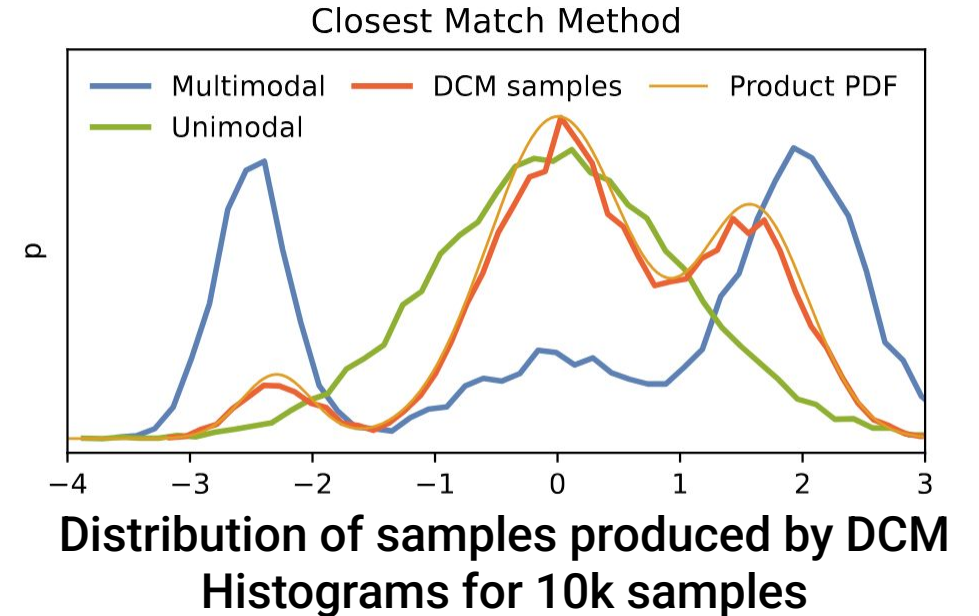
Let $\mathcal{D}_{\text{Diffusion}} = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$, $\mathcal{D}_{\text{Forecaster}} = \{\mathbf{y}_1, \dots, \mathbf{y}_l\}$, where $\mathbf{x}_i, \mathbf{y}_j \in \mathbb{R}^n$. Then DCM constructs $\mathbf{z} \in \mathbb{R}^n$ by

$$\mathbf{z}_d = \mathbf{y}_{j^*(d),d} \quad \text{where} \quad j^*(d) = \arg \min_j \left(\min_i |\mathbf{x}_{i,d} - \mathbf{y}_{j,d}| \right),$$

where $d = 1 \dots n$ are the dimensions.

For each dimension d , we choose the sample from $\mathcal{D}_{\text{Forecaster}}$ that has a closest sample in $\mathcal{D}_{\text{Diffusion}}$.

Highest performance compared to standard approaches
No hyperparameters & No fallback mechanism



Experimental Setup and Baselines

FORL is plug and play for different offline RL algorithms.

FORL- π vs. baselines:

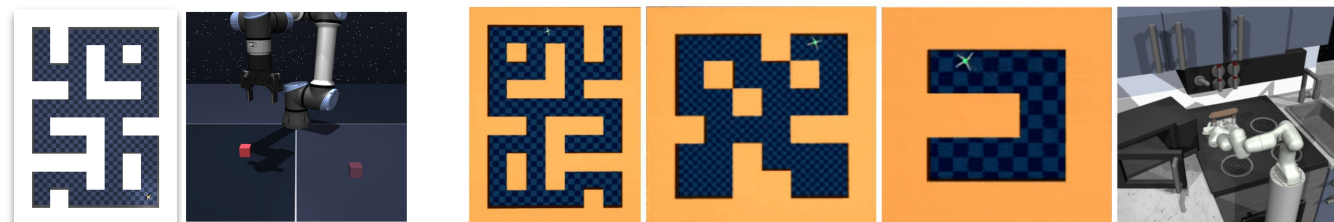
$$\pi \in \left\{ \begin{array}{l} \text{DQL [Wang et al., 2023], TD3BC [Fujimoto, 2021],} \\ \text{RORL [Yang, 2022], IQL [Kostrikov et al., 2021],} \\ \text{FQL [Park et al., 2025b]} \end{array} \right\}$$

- π
- $\pi\text{-Lag}$ [Rasul et al. 2023]
- $\pi^{(*)}\text{-DMBP}$ [Yang et al., 2024] - Lag [Rasul et al. 2023]
 $(*) \pi \in \{\text{DQL, TD3BC, RORL, IQL}\}$

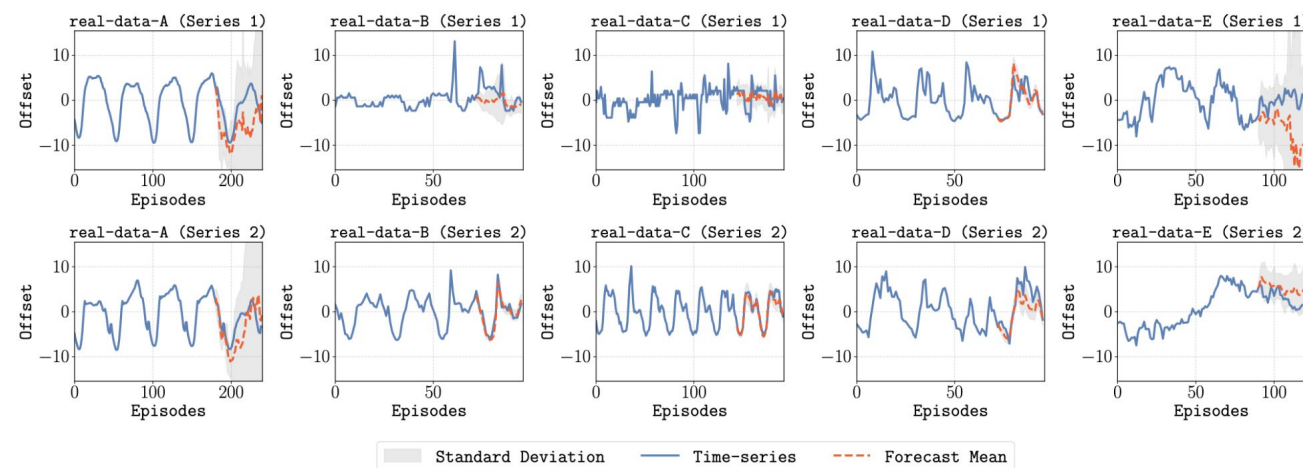
Offline RL Datasets

OGBench [Park et al., 2025a]

D4RL [Fu et al., 2020]



Time-series offsets from energy and finance data GluonTS [Alexandrov et al., 2020]

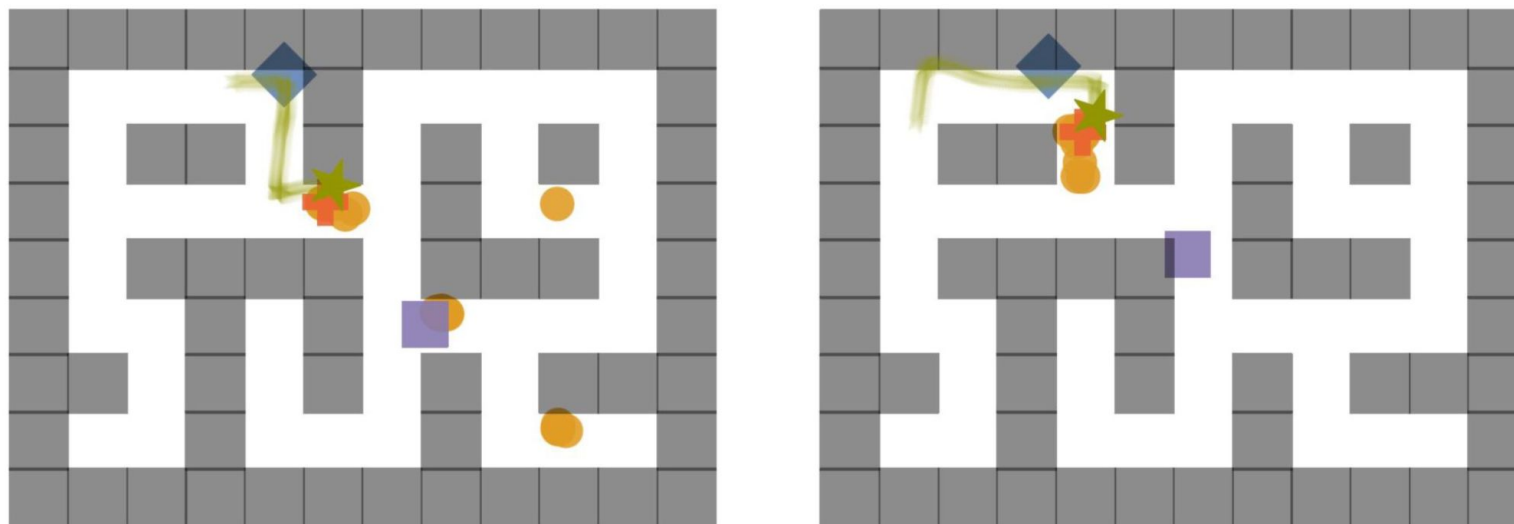


State Prediction Error

Visualization of predicted states and **true state** ★

+ FORL is the closest to the true state

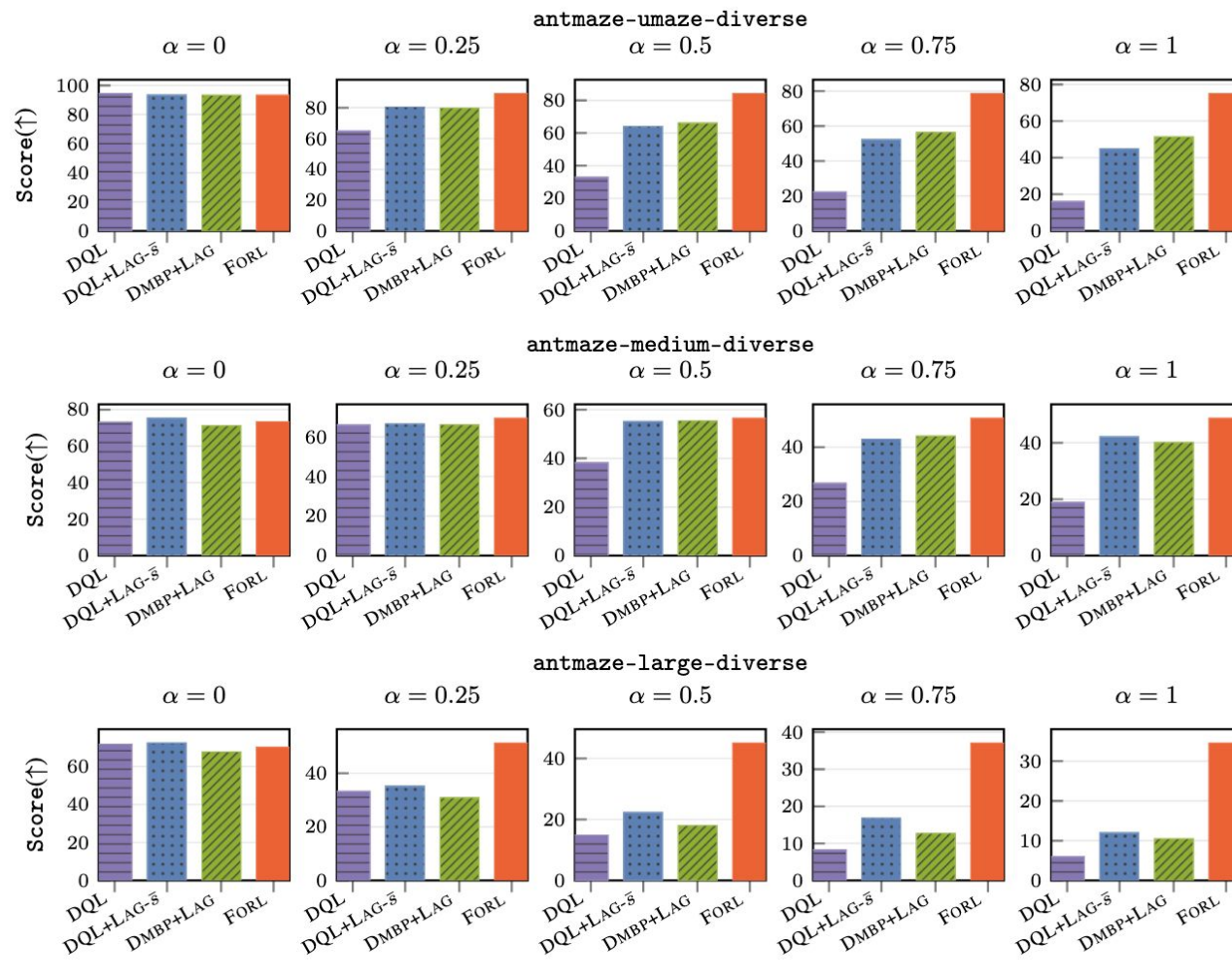
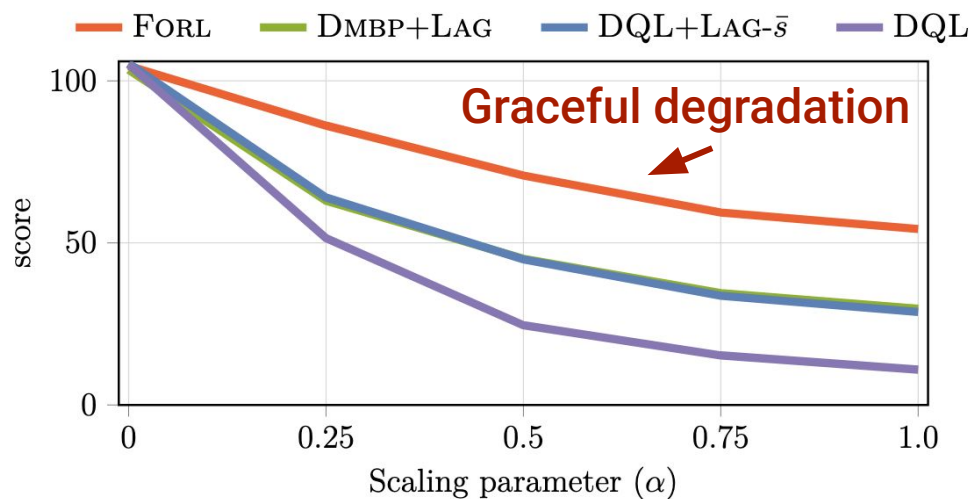
● $\{s_t^{(0)}\}_k$ ■ observation ◆ DQL+LAG- \bar{s} + FORL ★ state



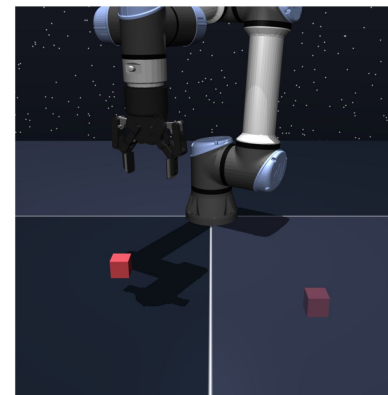
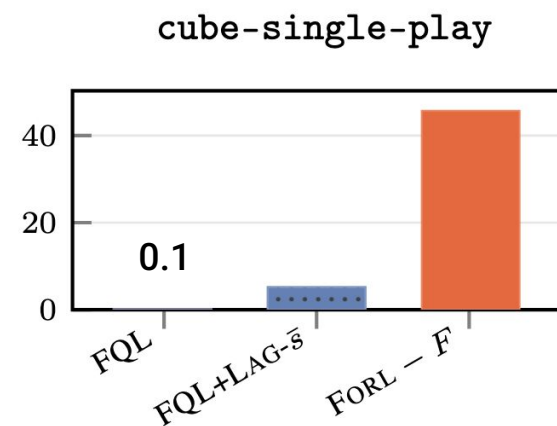
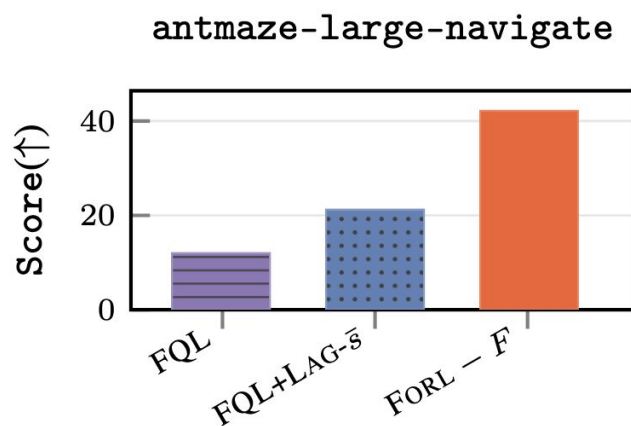
FORL has the lowest prediction error



How gracefully does performance degrade as the offset magnitude α is scaled from 0 (no offset) \rightarrow 1?



OGBench: FORL Outperforms the Baselines

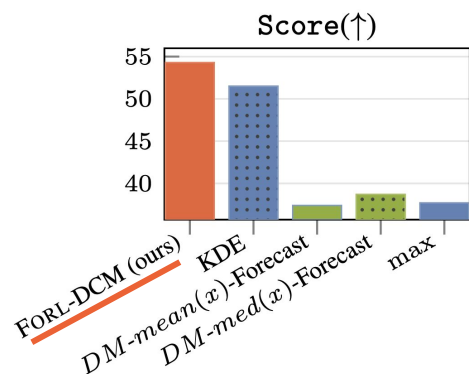




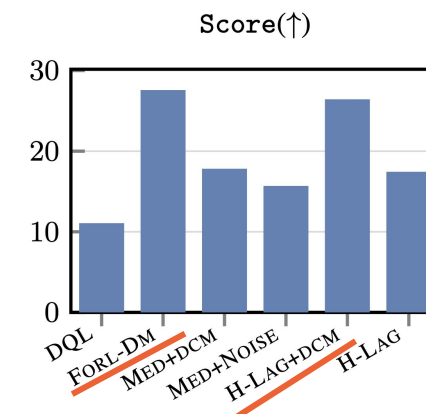
Thank you! We invite you to our project web-page!

<https://sites.google.com/view/forecastrl>

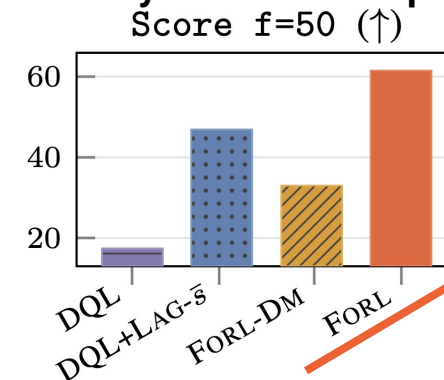
Dimension Wise Closest Match



What if we do not have access to the past offsets?



What if the offset changes are “in-episode” every $f=50$ timesteps?



Can FORL serve as a plug-and-play module for different offline-RL algorithms without retraining?



Suzan Ece Ada
Boğaziçi University
University of Tübingen



Georg Martius
University of Tübingen



Emre Uğur
Boğaziçi University



Erhan Öztop
Osaka University
Özyeğin University