## Introduction

Recent works such as Decision Transformer (DT, Chen et al. 2021) shows that offline RL problems can be casted as sequence modeling problems and solved by supervised learning methods.

The performance of offline RL however is bottlenecked by the dataset properties and often requires online finetuning for best results.
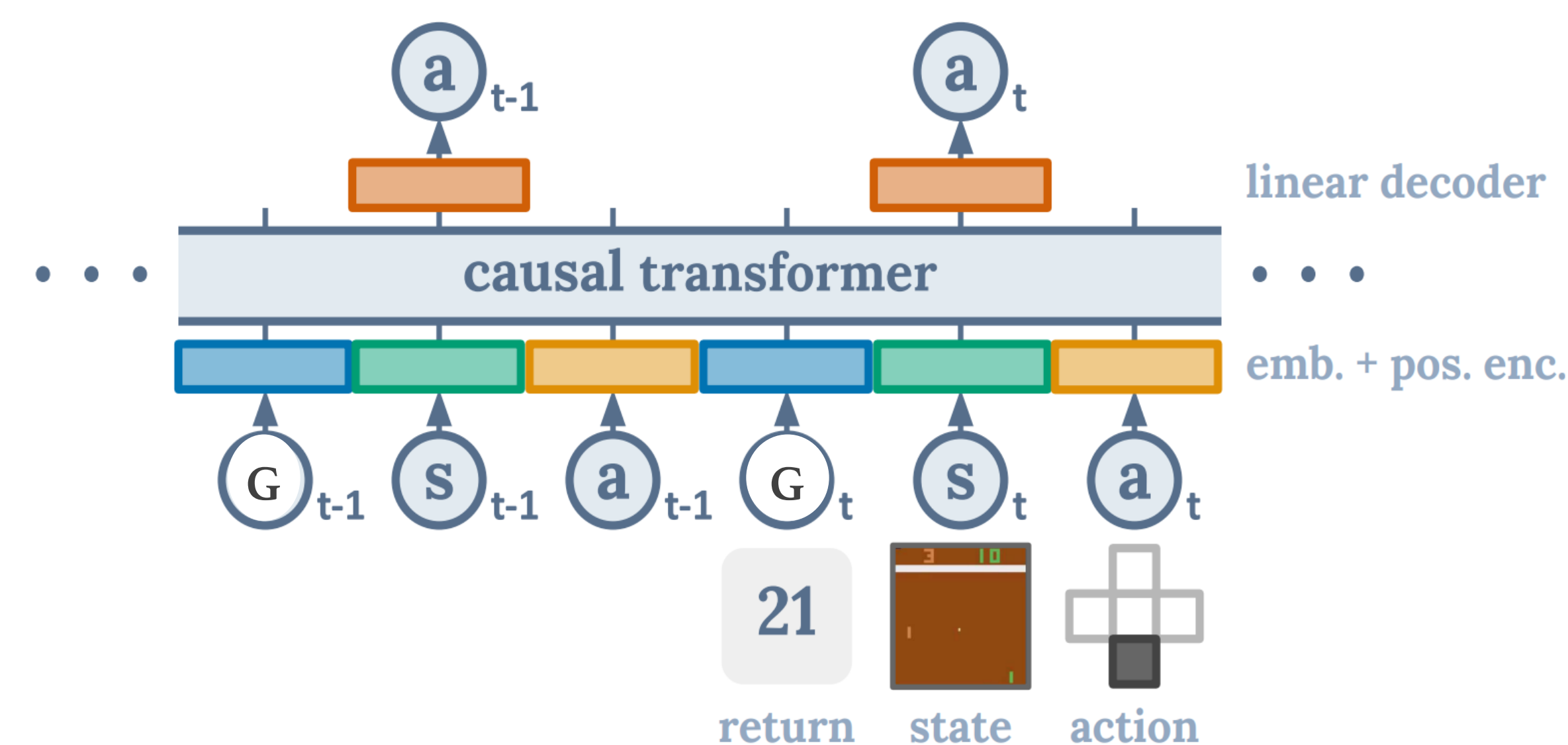
We propose Online Decision Transformers (ODT), an RL algorithm based on supervised sequence modeling that blends offline pretraining with online finetuning in a unified framework.

ODT enables stable online learning while retraining the simplicity of sequence modeling.

## Base Model

Decision Transformer (Chen et al. 2021) models a trajectory $\tau$ as (RTG, state, action) sequences.

RTG (return-to-go) $g_t = \sum_{t'=t}^{|\tau|} r_{t'}$



21
return    state    action

DT architecture (Chen et al. 2021)

DT generates return-conditioned policies.

Rollout
1. Specify the desired return $g_1$ and an initial state $s_1$.
2. Generate $a_1$, execute it and then observe $s_2$ and $r_1$.
3. Compute $g_2 = g_1 - r_1$. Now we can generate $a_2$.
4. Repeat until the episode terminates.

## Online Decision Transformer

### Stochastic Policy

$$\pi_\theta(a_t|\mathbf{s}_{-K,t}, \mathbf{g}_{-K,t}) = \mathcal{N}(\mu_\theta(\mathbf{s}_{-K,t}, \mathbf{g}_{-K,t}), \Sigma_\theta(\mathbf{s}_{-K,t}, \mathbf{g}_{-K,t}))$$

Generate action based on recent $K$ states and RTGs

### Max-Ent Sequence Modeling

$$\min_\theta J(\theta) \quad \text{subject to} \quad H_\theta^{\mathcal{T}}[\mathbf{a}|\mathbf{s}, \mathbf{g}] \geqslant \beta$$

- $J(\theta)$ negative log-likelihood of sequence data

$$J(\theta) = \frac{1}{K} \mathbb{E}_{(\mathbf{a},\mathbf{s},\mathbf{g})\sim\mathcal{T}}[-\log\pi_\theta(\mathbf{a}|\mathbf{s},\mathbf{g})]$$
$$= \frac{1}{K} \mathbb{E}_{(\mathbf{a},\mathbf{s},\mathbf{g})\sim\mathcal{T}}[-\sum_{k=1}^K \log\pi_\theta(a_k|\mathbf{s}_{-K,k}, \mathbf{g}_{-K,k})]$$

simple supervised learning, no return optimization

- $H_\theta^{\mathcal{T}}[\mathbf{a}|\mathbf{s},\mathbf{g}]$  sequence-level policy entropy

$$H_\theta^{\mathcal{T}}[\mathbf{a}|\mathbf{s},\mathbf{g}] = \frac{1}{K} \mathbb{E}_{(\mathbf{s},\mathbf{g})\sim\mathcal{T}}[H[\pi_\theta(\mathbf{a}|\mathbf{s},\mathbf{g})]]$$
$$= \frac{1}{K} \mathbb{E}_{(\mathbf{s},\mathbf{g})\sim\mathcal{T}}[\sum_{k=1}^K H[\pi_\theta(a_k|\mathbf{s}_{-K,k}, \mathbf{g}_{-K,k})]]$$

- $\beta$ -dim(action)

### Offline Pretraining + Online Finetuning

---
**Algorithm 1: Online Decision Transformer**

---
1 **Input:** offline data $\mathcal{T}_{\text{offline}}$, rounds $R$, exploration RTG $g_{\text{online}}$, buffer size $N$, gradient iterations $I$, pretrained policy $\pi_\theta$
2 **Intialization:** Replay buffer $\mathcal{T}_{\text{replay}} \leftarrow$ top $N$ trajectories in $\mathcal{T}_{\text{offline}}$.
3 **for** $round = 1, \ldots, R$ **do**
4      // use randomly sampled actions
     Trajectory $\tau \leftarrow$ Rollout using $\mathcal{M}$ and $\pi_\theta(\cdot|\mathbf{s}, \mathbf{g}(g_{\text{online}}))$.
5      $\mathcal{T}_{\text{replay}} \leftarrow \{\mathcal{T}_{\text{replay}} \setminus \{\text{the oldest trajectory}\}\} \bigcup \{\tau\}$.
6      $\pi_\theta \leftarrow$ Finetune ODT on $\mathcal{T}_{\text{replay}}$ for $I$ iterations via Algorithm 2.

---
**Algorithm 2: ODT Training**

---
1 **Input:** model parameters $\theta$, replay buffer $\mathcal{T}_{\text{replay}}$, training iterations $I$, context length $K$, batch size $B$
2 Compute the trajectory sampling probability $p(\tau) = |\tau|/\sum_{\tau\in\mathcal{T}} |\tau|$.
3 **for** $t = 1, \ldots, I$ **do**
4      Sample $B$ trajectories out of $\mathcal{T}_{\text{replay}}$ according to $p$.
5      **for** *each sampled trajectory $\tau$* **do**
       // Hindsight Return Relabeling
6        $\mathbf{g} \leftarrow$ the RTG sequence computed by the true rewards: $\mathbf{g}_t = \sum_{j=t}^{|\tau|} r_j, 1 \leqslant t \leqslant |\tau|$.
7        $(\mathbf{a}, \mathbf{s}, \mathbf{g}) \leftarrow$ a length $K$ sub-trajectory uniformly sampled from $\tau$.
8      $\theta \leftarrow$ one gradient update using the sampled $\{(\mathbf{a}, \mathbf{s}, \mathbf{g})\}$s.
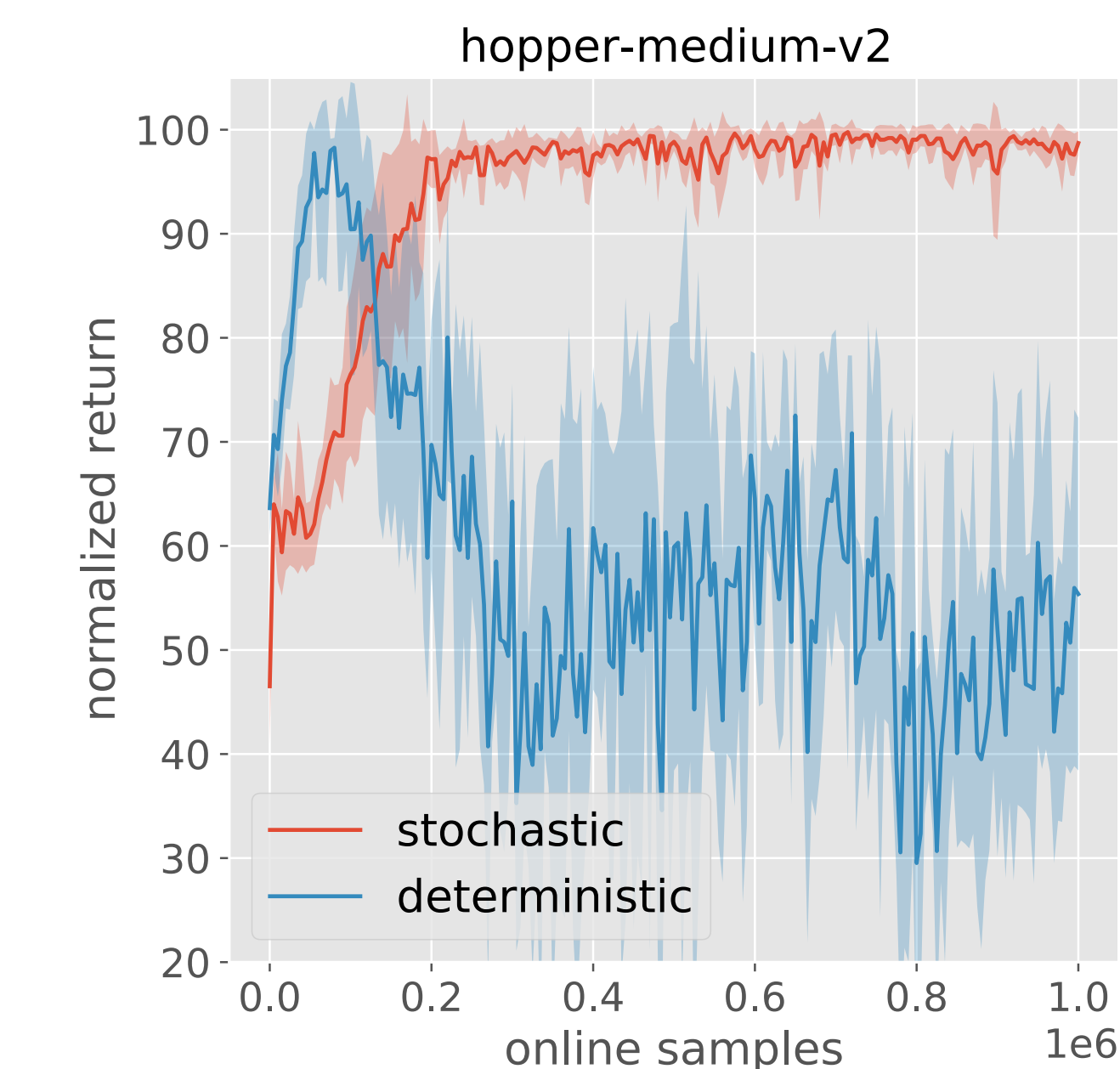
---

hindsight return relabeling - use observed return instead of target return
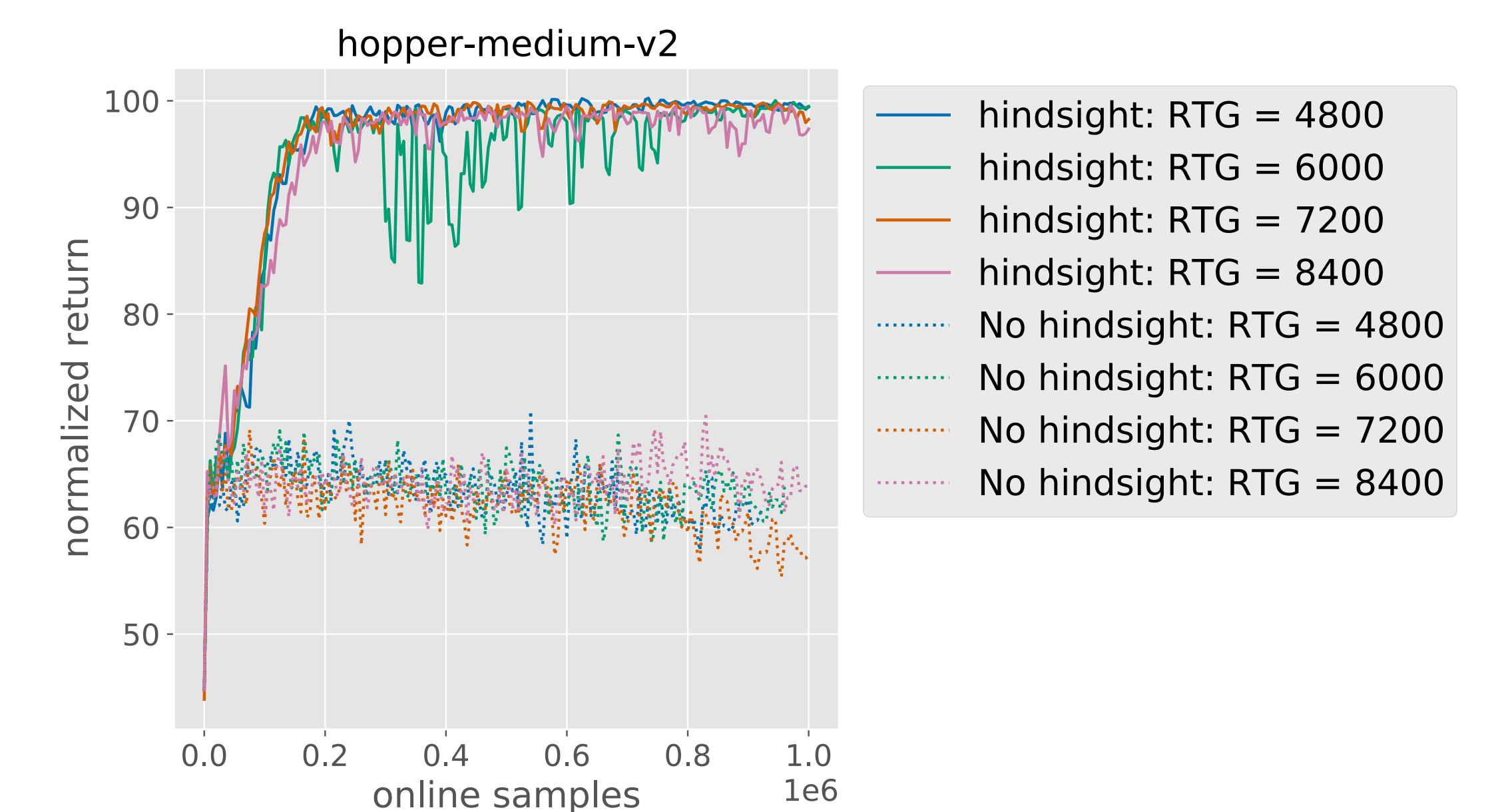
## Benchmark Comparison

| dataset | ODT (offline) | ODT (0.2m) | $\delta_{\text{ODT}}$ | IQL (offline) | IQL (0.2m) | $\delta_{\text{IQL}}$ |
|---|---|---|---|---|---|---|
| hopper-medium | $66.95 \pm 3.26$ | $\mathbf{97.54 \pm 2.10}$ | $\mathbf{30.59}$ | $63.81 \pm 9.15$ | $66.79 \pm 4.07$ | $2.98$ |
| hopper-medium-replay | $86.64 \pm 5.41$ | $88.89 \pm 6.33$ | $2.25$ | $92.13 \pm 10.43$ | $\mathbf{96.23 \pm 4.35}$ | $\mathbf{4.10}$ |
| walker2d-medium | $72.19 \pm 6.49$ | $76.79 \pm 2.30$ | $\mathbf{4.60}$ | $79.89 \pm 3.06$ | $\mathbf{80.33 \pm 2.33}$ | $0.44$ |
| walker2d-medium-replay | $68.92 \pm 4.79$ | $\mathbf{76.86 \pm 4.04}$ | $\mathbf{7.94}$ | $73.67 \pm 6.37$ | $70.55 \pm 5.81$ | $-3.12$ |
| halfcheetah-medium | $42.72 \pm 0.46$ | $42.16 \pm 1.48$ | $-0.56$ | $47.37 \pm 0.29$ | $47.41 \pm 0.15$ | $\mathbf{0.04}$ |
| halfcheetah-medium-replay | $39.99 \pm 0.68$ | $40.42 \pm 1.61$ | $\mathbf{0.43}$ | $44.10 \pm 1.14$ | $44.14 \pm 0.3$ | $0.04$ |
| ant-medium | $91.33 \pm 4.13$ | $90.79 \pm 5.80$ | $-0.54$ | $99.92 \pm 5.86$ | $\mathbf{100.85 \pm 2.02}$ | $\mathbf{0.93}$ |
| ant-medium-replay | $86.56 \pm 3.26$ | $\mathbf{91.57 \pm 2.73}$ | $\mathbf{5.01}$ | $91.21 \pm 7.27$ | $91.36 \pm 1.47$ | $0.15$ |
| sum | | $\mathbf{605.02}$ | $\mathbf{49.72}$ | | $597.66$ | $5.56$ |
| antmaze-umaze | $53.10 \pm 4.21$ | $\mathbf{88.5 \pm 5.88}$ | $\mathbf{35.4}$ | $87.1 \pm 2.81$ | $\mathbf{89.5 \pm 5.43}$ | $2.4$ |
| antmaze-umaze-diverse | $50.20 \pm 6.69$ | $\mathbf{56.00 \pm 5.69}$ | $\mathbf{7.99}$ | $64.4 \pm 8.95$ | $\mathbf{56.8 \pm 6.42}$ | $-7.6$ |
| sum | | $\mathbf{144.5}$ | $\mathbf{43.39}$ | | $146.3$ | $-5.2$ |

Baseline: Implicit Q-Learning (IQL, Kostrikov 2021)
Absolute performance: ODT is comparable
Finetuning Gain: ODT is much better

## Ablation Study



Stochasticity is important to enable stable performance improvement in online training



Hindsight return relabeling is critical for correcting bias in collected data