

# A Unified Framework for Alternating Offline Model Training and Policy Learning

Shentao Yang<sup>1</sup>, Shujian Zhang<sup>1</sup>, Yihao Feng<sup>2</sup>, Mingyuan Zhou<sup>1</sup>

*<sup>1</sup>The University of Texas at Austin, <sup>2</sup>Salesforce Research*

*October, 2022*

# Proposed Method Sketch

- **Motivation:** model training = MLE  $\neq$  improve policy = model usage.



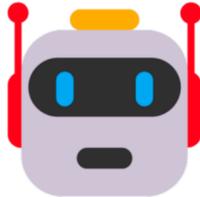
$$\min_{\pi, \hat{P}} C \cdot \sqrt{D_{\pi}(P^*, \hat{P})} \geq \left| J(\pi, P^*) - J(\pi, \hat{P}) \right|$$

# Proposed Method Sketch

- **Motivation:** model training = MLE  $\neq$  improve policy = model usage.



$$\min_{\pi, \hat{P}} C \cdot \sqrt{D_{\pi}(P^*, \hat{P})} \geq \left| J(\pi, P^*) - J(\pi, \hat{P}) \right|$$

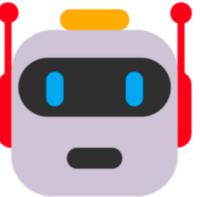
- **Jointly train**  and  to minimize an upper bound of the evaluation error.

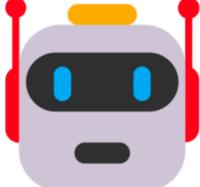
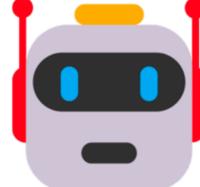
# Proposed Method Sketch

- **Motivation:** model training = MLE  $\neq$  improve policy = model usage.



$$\min_{\pi, \hat{P}} C \cdot \sqrt{D_{\pi}(P^*, \hat{P})} \geq \left| J(\pi, P^*) - J(\pi, \hat{P}) \right|$$

- **Jointly train**  and  to minimize an upper bound of the evaluation error.

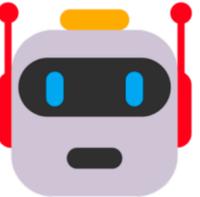
- Fixed ,   $\approx$   only on state-actions visited by .

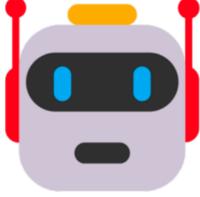
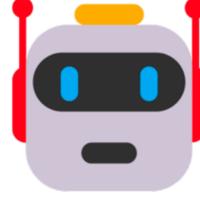
# Proposed Method Sketch

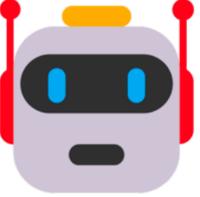
- **Motivation:** model training = MLE  $\neq$  improve policy = model usage.



$$\min_{\pi, \hat{P}} C \cdot \sqrt{D_{\pi}(P^*, \hat{P})} \geq \left| J(\pi, P^*) - J(\pi, \hat{P}) \right|$$

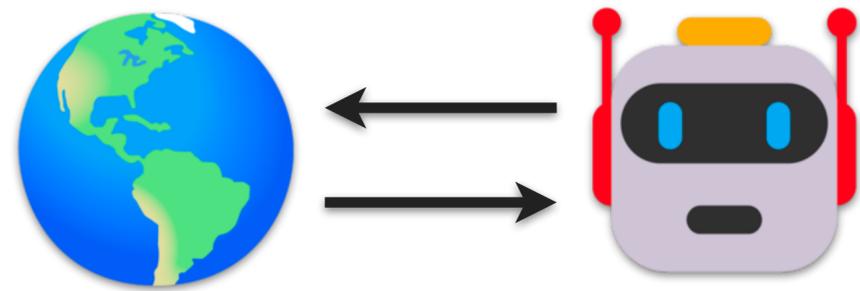
- **Jointly train**  and  to minimize an upper bound of the evaluation error.

- Fixed ,   $\approx$   only on state-actions visited by .

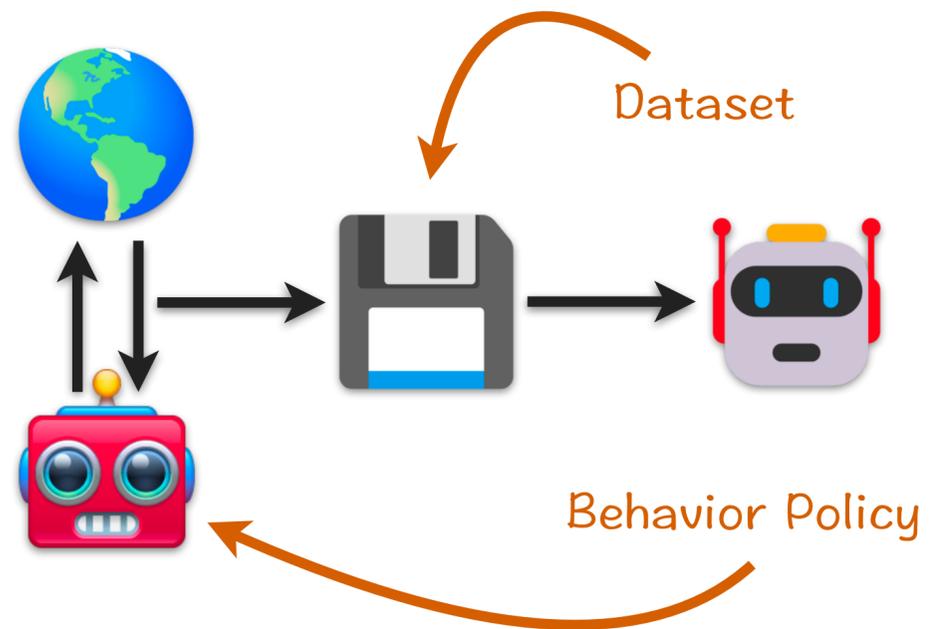
- Fixed , optimize  with a regularization based on .

# Background

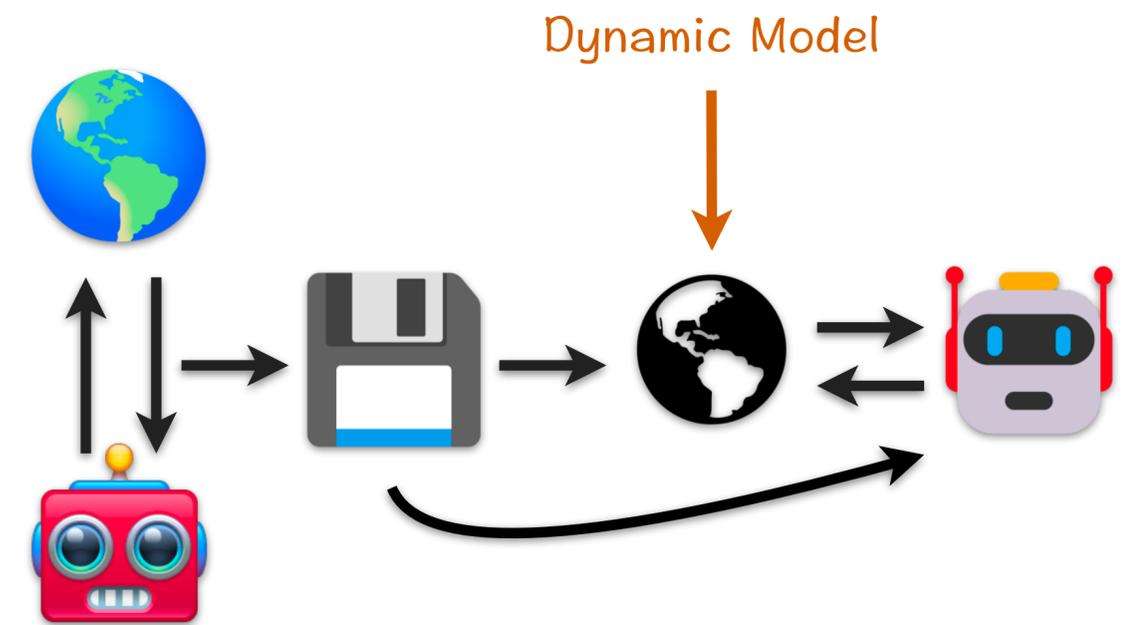
- Offline RL: learn policy from **static** datasets.



(a) Classical RL



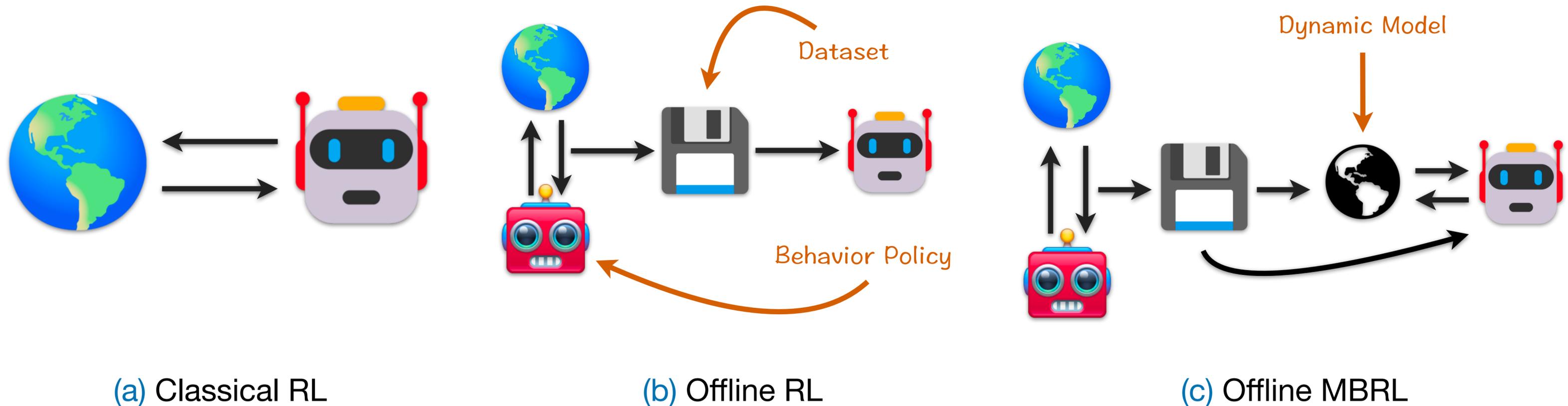
(b) Offline RL



(c) Offline MBRL

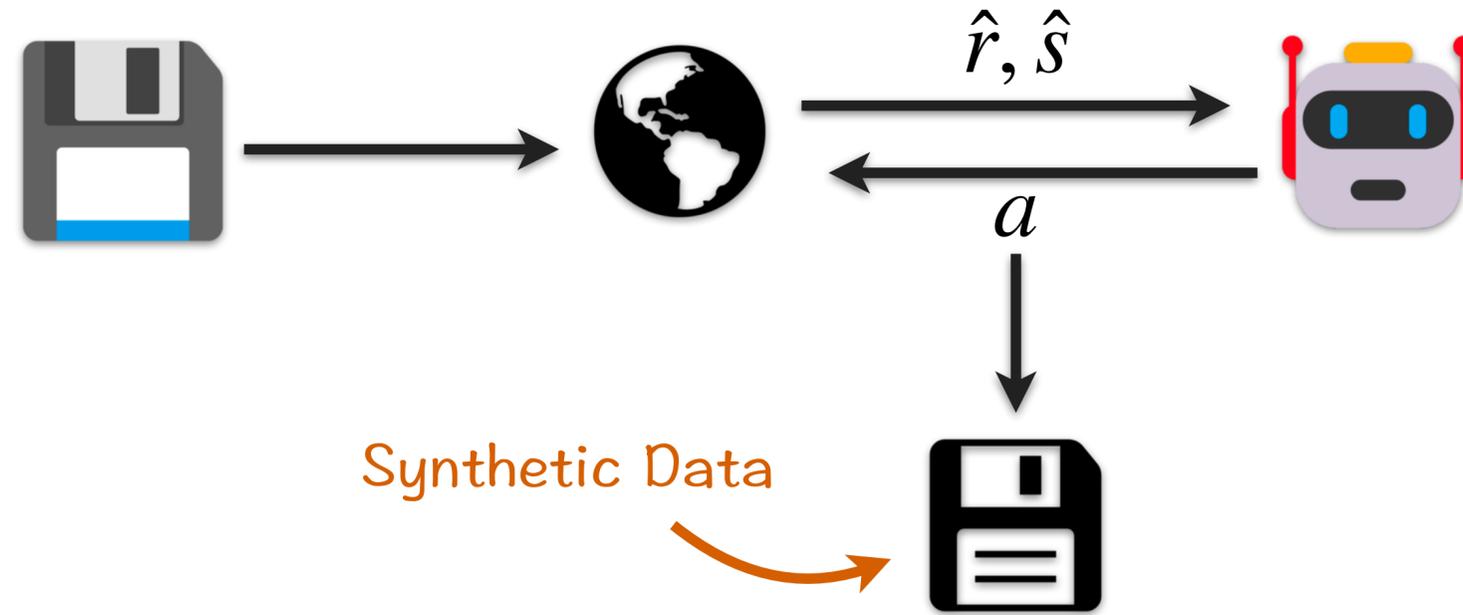
# Background

- Offline RL: learn policy from **static** datasets.
- Offline Model-Based RL (Offline MBRL): learn dynamic from static datasets.



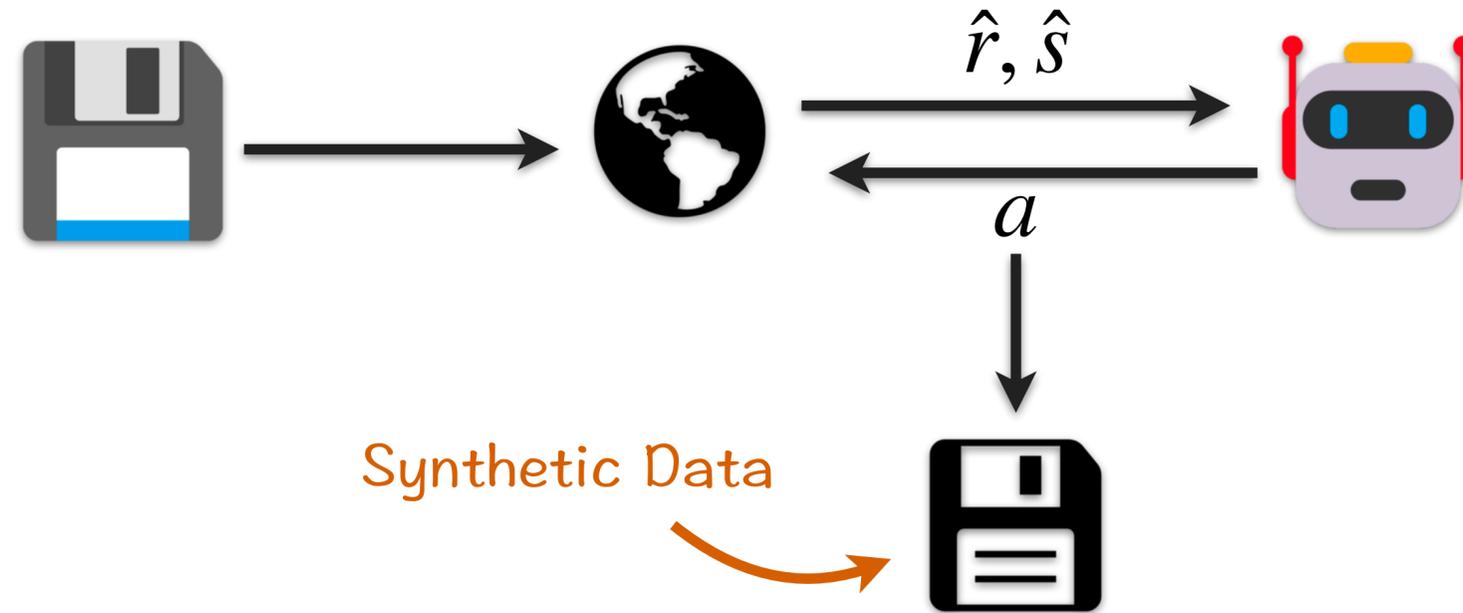
# Background

- Benefits of offline MBRL



# Background

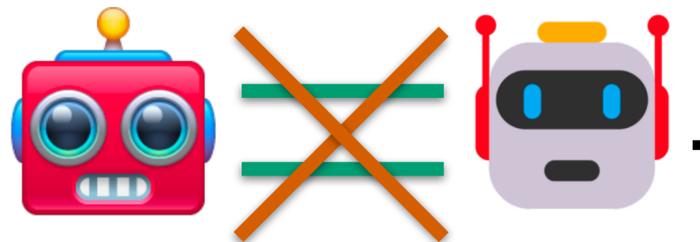
- Benefits of offline MBRL



- Offline model-free RL

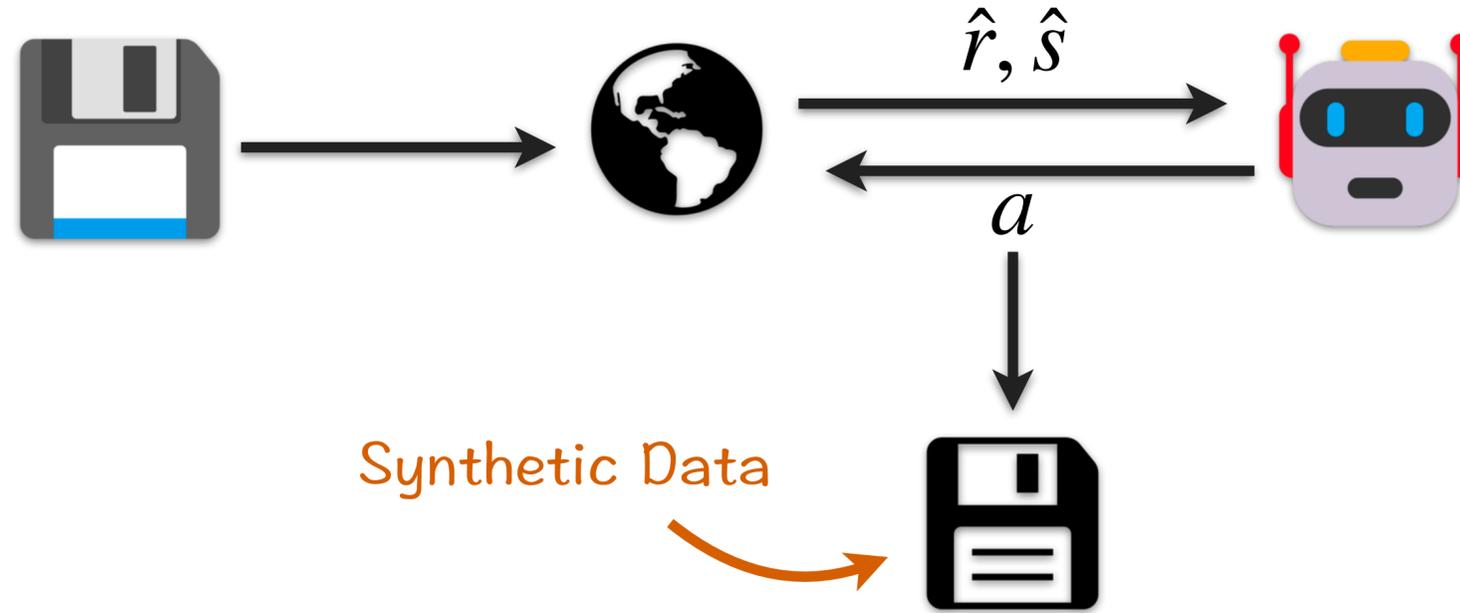
- Only know reward and next state at state-actions **within** the dataset.

- Off-policy issue



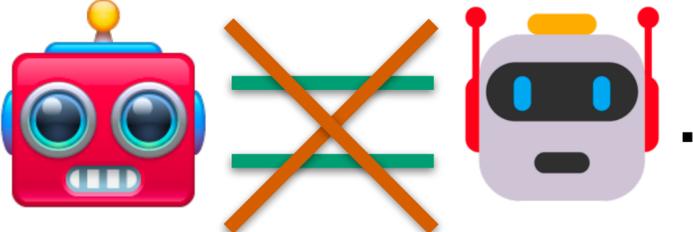
# Background

- Benefits of offline MBRL



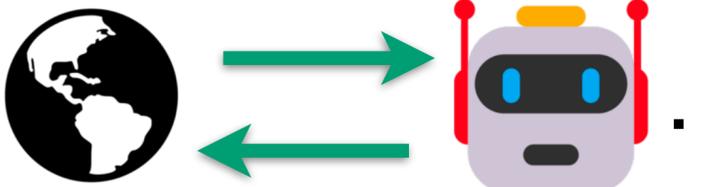
- Offline model-free RL

- Only know reward and next state at state-actions **within** the dataset.

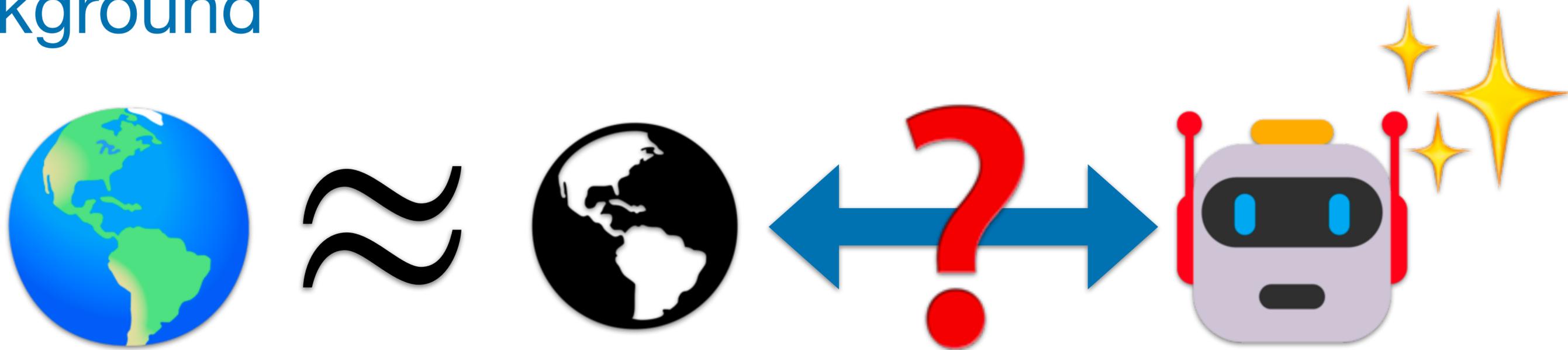
- Off-policy issue  .
- The diagram shows a red robot icon on the left and a purple robot icon on the right. Two horizontal green lines connect them, and a large orange 'X' is drawn over these lines, indicating that the two robots are not interacting or that the data is off-policy.

- Offline model-based RL

- Estimate reward and next state at **new** state-actions.

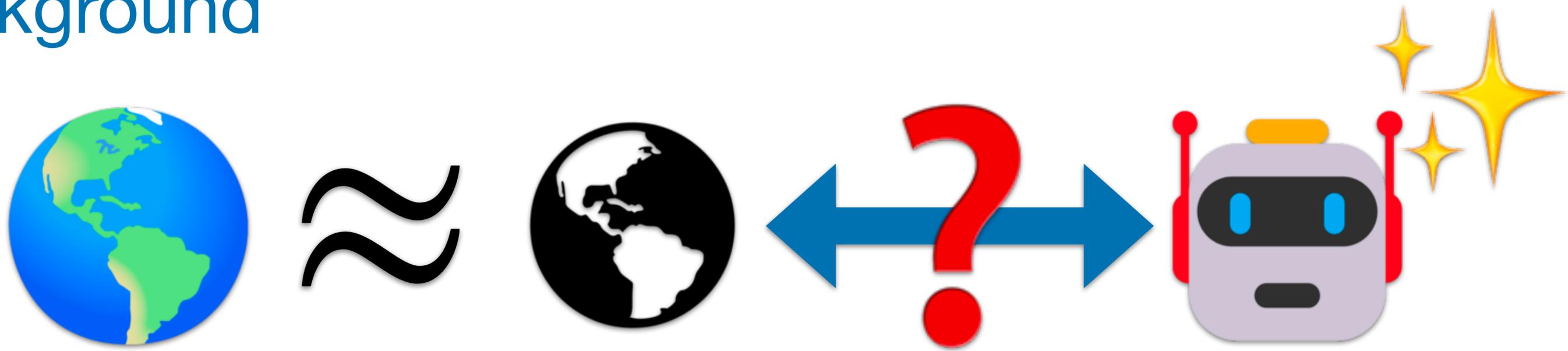
- $\approx$  on-policy  .
- The diagram shows a globe icon on the left and a purple robot icon on the right. Two horizontal green arrows connect them, one pointing from the globe to the robot and one pointing from the robot to the globe, indicating that the robot is interacting with the environment.

# Background



- Most offline MBRL: **pre-train** a **fixed** dynamic model on  .
  - Objective: MLE — “simply a mimic of the world.”
  - Usage: improve the policy.

# Background



- Most offline MBRL: **pre-train** a **fixed** dynamic model on  .
  - Objective: MLE — “simply a mimic of the world.”
  - Usage: improve the policy.
- Objective mismatch: **model training**  $\neq$  **model usage**.
  - Especially when  is limited and  is hard to learn.

# Proposed Method: Bounding the Evaluation Error

- A tractable upper bound for the evaluation error

$$\left| J(\pi, P^*) - J(\pi, \hat{P}) \right| \leq C \cdot \sqrt{D_\pi(P^*, \hat{P})}, \quad \text{with}$$

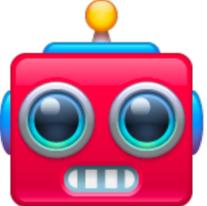
$$D_\pi(P^*, \hat{P}) \triangleq \mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \text{KL} \left( P^*(s' | s, a) \pi_b(a' | s') \parallel \hat{P}(s' | s, a) \pi(a' | s') \right) \right],$$

# Proposed Method: Bounding the Evaluation Error

- A tractable upper bound for the evaluation error

$$\left| J(\pi, P^*) - J(\pi, \hat{P}) \right| \leq C \cdot \sqrt{D_\pi(P^*, \hat{P})}, \quad \text{with}$$

$$D_\pi(P^*, \hat{P}) \triangleq \mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \text{KL} \left( P^*(s' | s, a) \pi_b(a' | s') \parallel \hat{P}(s' | s, a) \pi(a' | s') \right) \right],$$

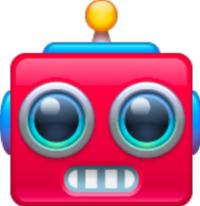
- $\pi_b$  is the behavior policy  .

# Proposed Method: Bounding the Evaluation Error

- A tractable upper bound for the evaluation error

$$\left| J(\pi, P^*) - J(\pi, \hat{P}) \right| \leq C \cdot \sqrt{D_\pi(P^*, \hat{P})}, \quad \text{with}$$

$$D_\pi(P^*, \hat{P}) \triangleq \mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \text{KL} \left( P^*(s' | s, a) \pi_b(a' | s') \parallel \hat{P}(s' | s, a) \pi(a' | s') \right) \right],$$

- $\pi_b$  is the behavior policy  .

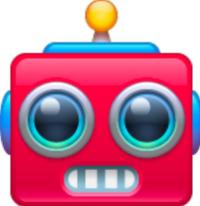
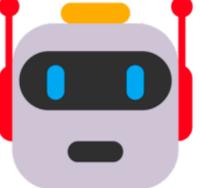
- $d_{\pi_b, \gamma}^{P^*}$  is the offline-data distribution  .

# Proposed Method: Bounding the Evaluation Error

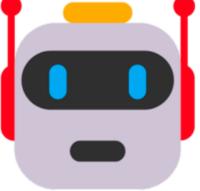
- A tractable upper bound for the evaluation error

$$\left| J(\pi, P^*) - J(\pi, \hat{P}) \right| \leq C \cdot \sqrt{D_\pi(P^*, \hat{P})}, \quad \text{with}$$

$$D_\pi(P^*, \hat{P}) \triangleq \mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \text{KL} \left( P^*(s' | s, a) \pi_b(a' | s') \parallel \hat{P}(s' | s, a) \pi(a' | s') \right) \right],$$

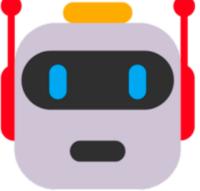
- $\pi_b$  is the behavior policy  .
- $d_{\pi_b, \gamma}^{P^*}$  is the offline-data distribution  .
- $\omega(s, a) \triangleq \frac{d_{\pi, \gamma}^{P^*}(s, a)}{d_{\pi_b, \gamma}^{P^*}(s, a)}$  is the density ratio between  and visitation freq. of  .

# Proposed Method: Model Training

- Fix , we train the model  by

$$\ell(\hat{P}) \triangleq - \mathbb{E}_{(s,a,s') \sim d_{\pi_b}^{P^*}} \left[ \omega(s, a) \log \left\{ \hat{P}(s' | s, a) \right\} \right] = D_{\pi}(P^*, \hat{P}) - C', \quad \text{with } C' \text{ a constant to } \hat{P}.$$

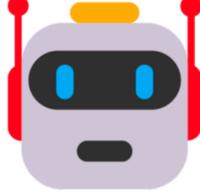
# Proposed Method: Model Training

- Fix , we train the model  by

$$\ell(\hat{P}) \triangleq - \mathbb{E}_{(s,a,s') \sim d_{\pi_b}^{P^*}} \left[ \omega(s, a) \log \left\{ \hat{P}(s' | s, a) \right\} \right] = D_{\pi}(P^*, \hat{P}) - C', \quad \text{with } C' \text{ a constant to } \hat{P}.$$

- $(s, a, s')$  is one transition in .

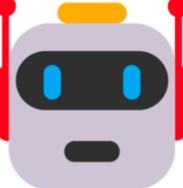
# Proposed Method: Model Training

- Fix , we train the model  by

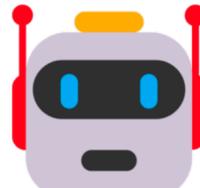
$$\ell(\hat{P}) \triangleq - \mathbb{E}_{(s,a,s') \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \log \left\{ \hat{P}(s' | s, a) \right\} \right] = D_{\pi}(P^*, \hat{P}) - C', \quad \text{with } C' \text{ a constant to } \hat{P}.$$

- $(s, a, s')$  is one transition in .
- Given  $\omega(s, a)$ , a stable **weighted MLE** objective.

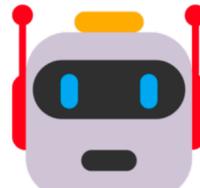
# Proposed Method: Policy Learning

- A lower-bound of  performance:  $J(\pi, \hat{P}) - C \cdot \sqrt{D_{\pi}(P^*, \hat{P})}$ .

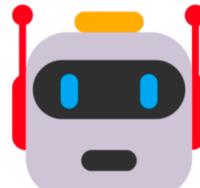
# Proposed Method: Policy Learning

- A lower-bound of  performance:  $J(\pi, \hat{P}) - C \cdot \sqrt{D_{\pi}(P^*, \hat{P})}$ .
- Fix , empirically helpful to construct the regularizer by:

# Proposed Method: Policy Learning

- A lower-bound of  performance:  $J(\pi, \hat{P}) - C \cdot \sqrt{D_{\pi}(P^*, \hat{P})}$ .
- Fix , empirically helpful to construct the regularizer by:
  - Removing the  $\sqrt{\cdot}$ .

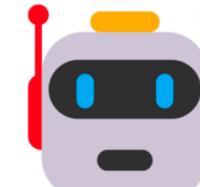
# Proposed Method: Policy Learning

- A lower-bound of  performance:  $J\left(\pi, \widehat{P}\right) - C \cdot \sqrt{D_{\pi}(P^*, \widehat{P})}$ .

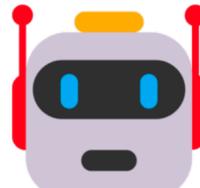
- Fix , empirically helpful to construct the regularizer by:

- Removing the  $\sqrt{\cdot}$ .
- Applying a further relaxation

$$D_{\pi}(P^*, \widehat{P}) \leq C'' \cdot \text{KL}\left(P^*(s' | s, a) \pi_b(a' | s') d_{\pi_b, \gamma}^{P^*}(s, a) \parallel \widehat{P}(s' | s, a) \pi(a' | s') d_{\pi_b, \gamma}^{P^*}(s) \pi(a | s)\right)$$

- Stronger regularizer: regularizes  at both  $s$  and  $s'$ .

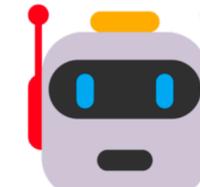
# Proposed Method: Policy Learning

- A lower-bound of  performance:  $J\left(\pi, \widehat{P}\right) - C \cdot \sqrt{D_{\pi}(P^*, \widehat{P})}$ .

- Fix , empirically helpful to construct the regularizer by:

- Removing the  $\sqrt{\cdot}$ .
- Applying a further relaxation

$$D_{\pi}(P^*, \widehat{P}) \leq C'' \cdot \text{KL}\left(P^*(s' | s, a) \pi_b(a' | s') d_{\pi_b, \gamma}^{P^*}(s, a) \parallel \widehat{P}(s' | s, a) \pi(a' | s') d_{\pi_b, \gamma}^{P^*}(s) \pi(a | s)\right)$$

- Stronger regularizer: regularizes  at both  $s$  and  $s'$ .
- Changing KL-divergence to Jensen-Shannon divergence.

# Proposed Method: Density-Ratio Training

- Fixed-point style method, ~~saddle-point optimization~~.

# Proposed Method: Density-Ratio Training

- Fixed-point style method, ~~saddle-point optimization~~.
- A simple MSE objective:

$$\mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s, a) \right] = \gamma \mathbb{E}_{\substack{(s, a, s') \sim d_{\pi_b, \gamma}^{P^*} \\ a' \sim \pi(\cdot | s')}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s', a') \right] + (1 - \gamma) \mathbb{E}_{\substack{s \sim \mu_0(\cdot) \\ a \sim \pi(\cdot | s)}} \left[ Q_{\pi}^{\hat{P}}(s, a) \right].$$

# Proposed Method: Density-Ratio Training

- Fixed-point style method, ~~saddle-point optimization~~.

- A simple MSE objective:

$$\mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s, a) \right] = \gamma \mathbb{E}_{\substack{(s, a, s') \sim d_{\pi_b, \gamma}^{P^*} \\ a' \sim \pi(\cdot | s')}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s', a') \right] + (1 - \gamma) \mathbb{E}_{\substack{s \sim \mu_0(\cdot) \\ a \sim \pi(\cdot | s)}} \left[ Q_{\pi}^{\hat{P}}(s, a) \right].$$

- Based on the “forward” Bellman equation for  $\omega(s, a)$  —not tractable 😭 !

- Use Q-function as test function and  $\sum_{(s', a')}$  on both sides.
- Primal-dual relation between  $\omega(s, a)$  and Q-function in OPE.

# Proposed Method: Density-Ratio Training

- Fixed-point style method, ~~saddle-point optimization~~.

- A simple MSE objective:

$$\mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s, a) \right] = \gamma \mathbb{E}_{\substack{(s, a, s') \sim d_{\pi_b, \gamma}^{P^*} \\ a' \sim \pi(\cdot | s')}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s', a') \right] + (1 - \gamma) \mathbb{E}_{\substack{s \sim \mu_0(\cdot) \\ a \sim \pi(\cdot | s)}} \left[ Q_{\pi}^{\hat{P}}(s, a) \right].$$

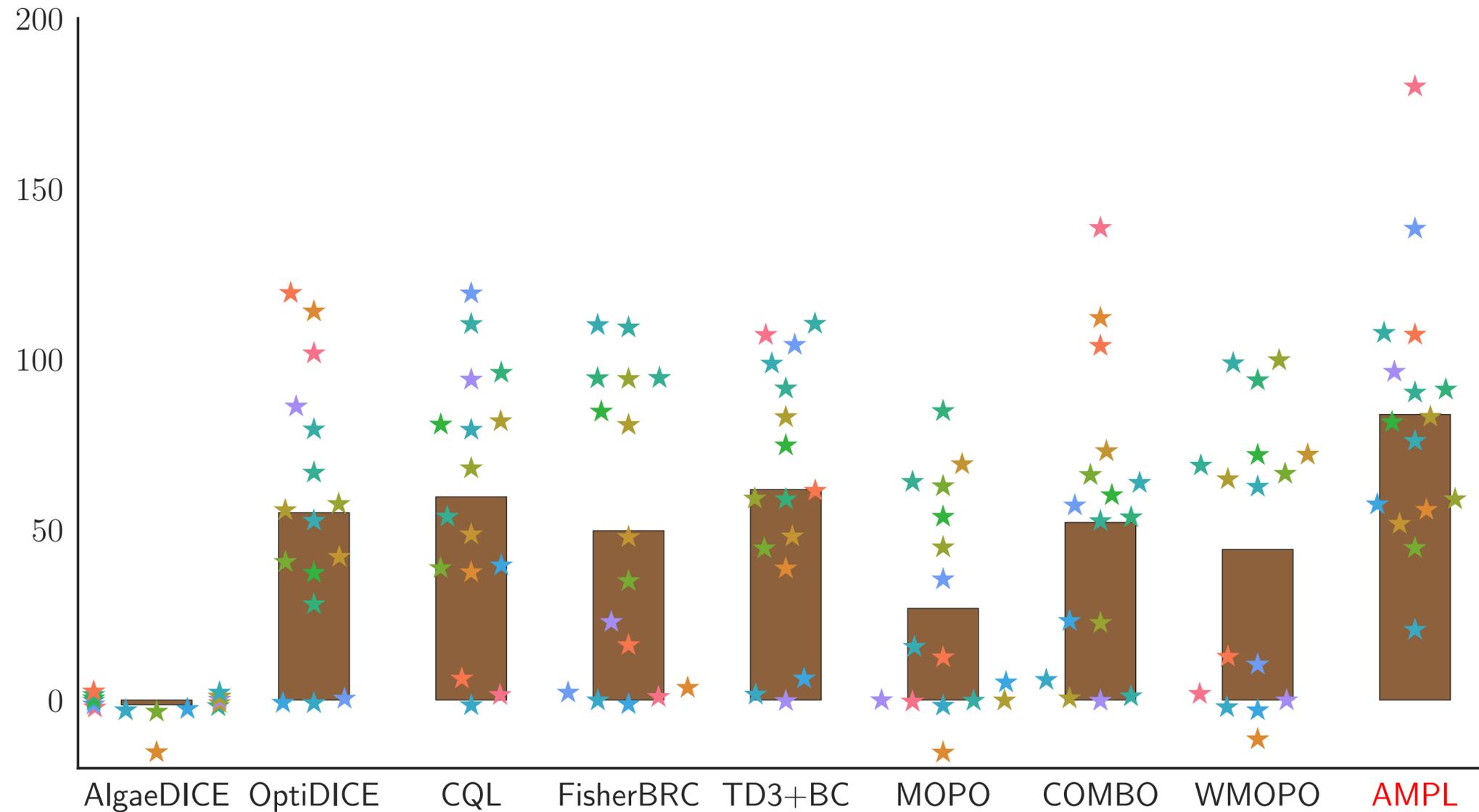
- Based on the “forward” Bellman equation for  $\omega(s, a)$  —not tractable 😭 !

- Use Q-function as test function and  $\sum_{(s', a')}$  on both sides.

- Primal-dual relation between  $\omega(s, a)$  and Q-function in OPE.

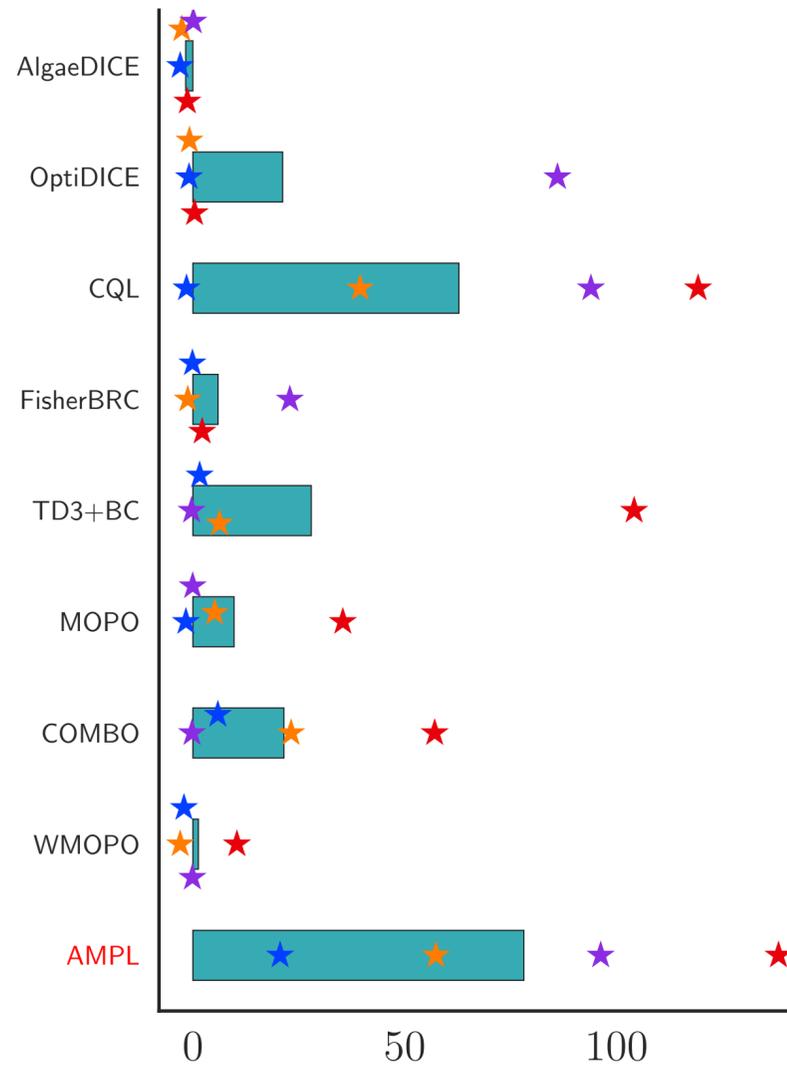
- Only requires samples from  and the initial state-distribution.

# Results: Main Method



- Our offline Alternating Model-Policy Learning (AMPL) performs well on D4RL tasks.

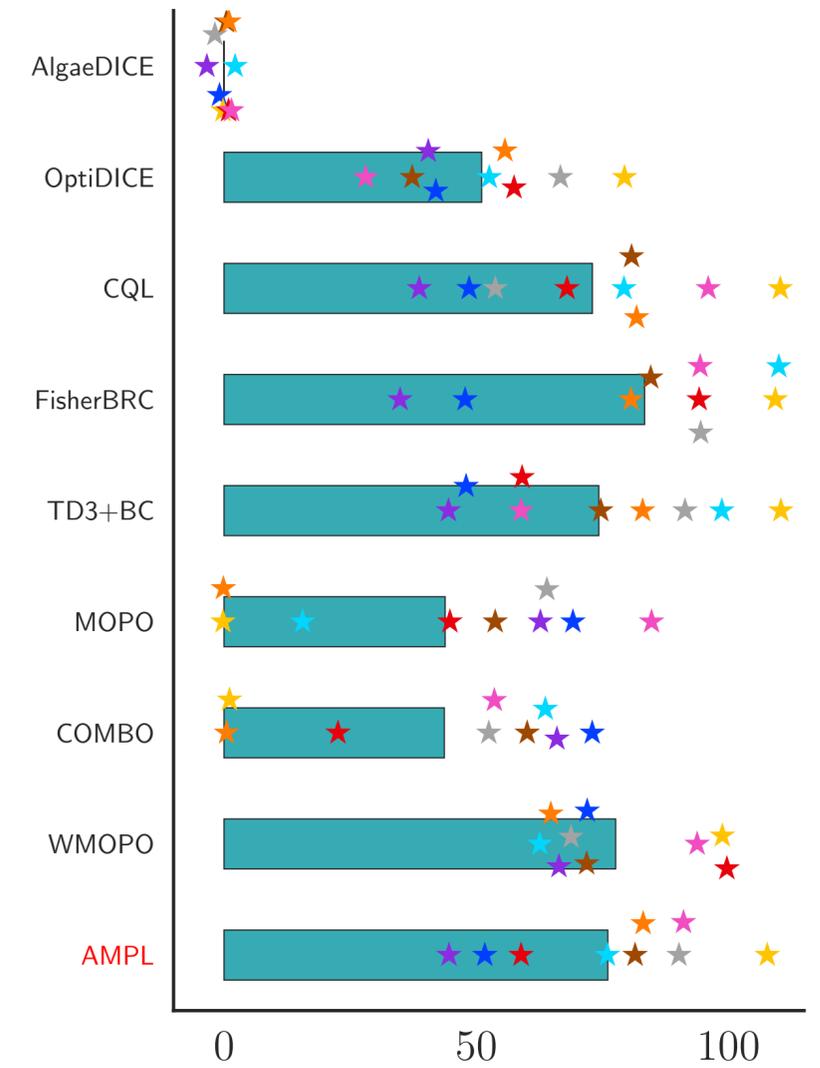
# Results: Main Method



(a) Adroit



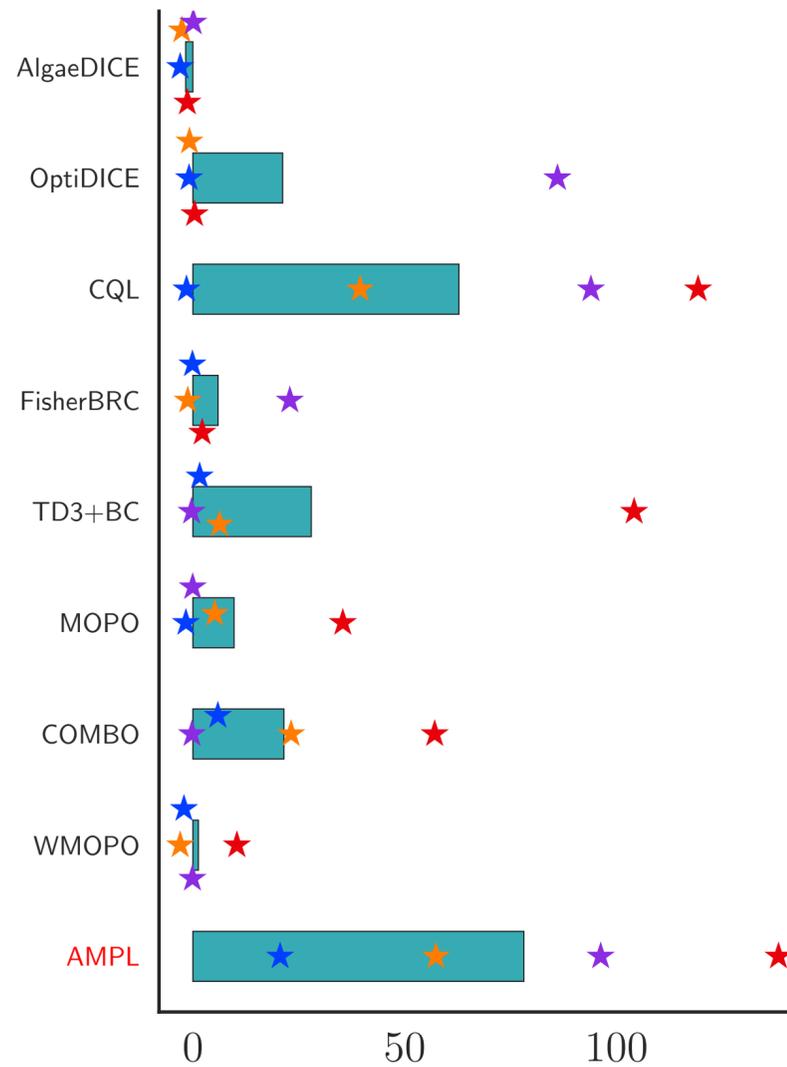
(b) Maze2D



(c) MuJoCo

- Learn well on the MuJoCo datasets.

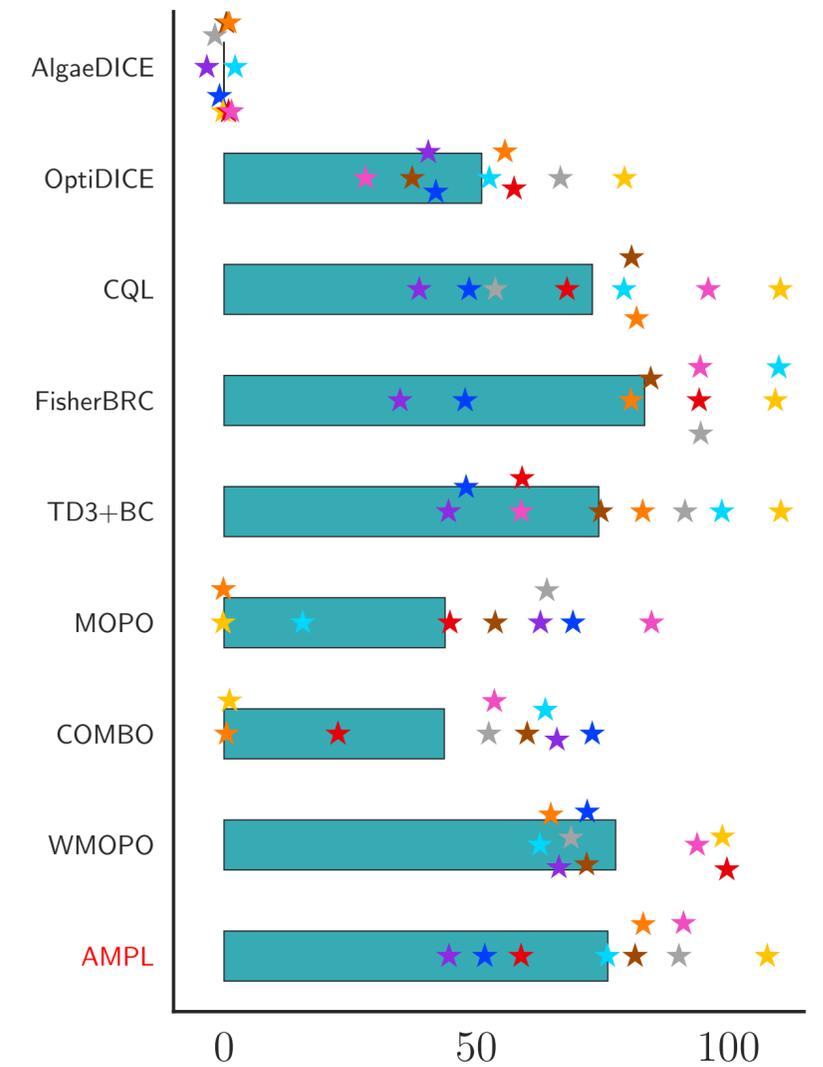
# Results: Main Method



(a) Adroit



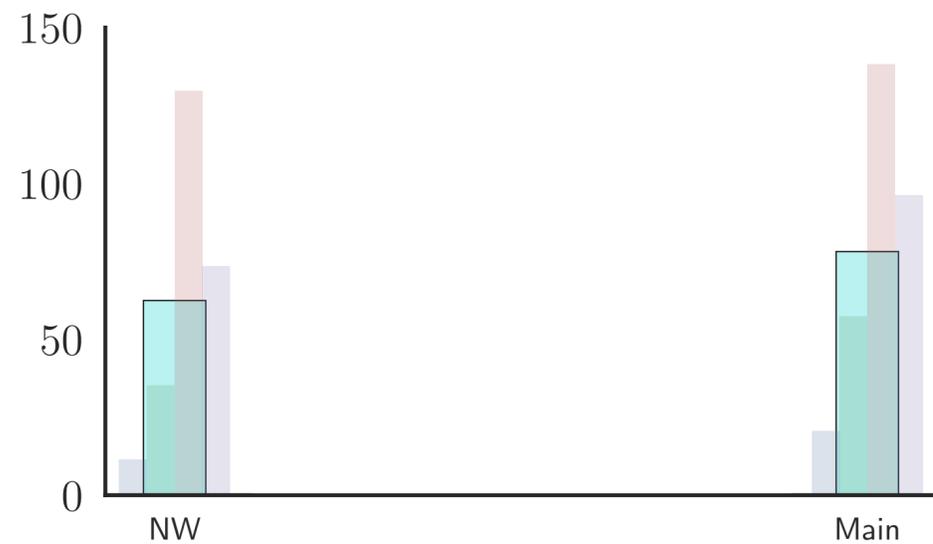
(b) Maze2D



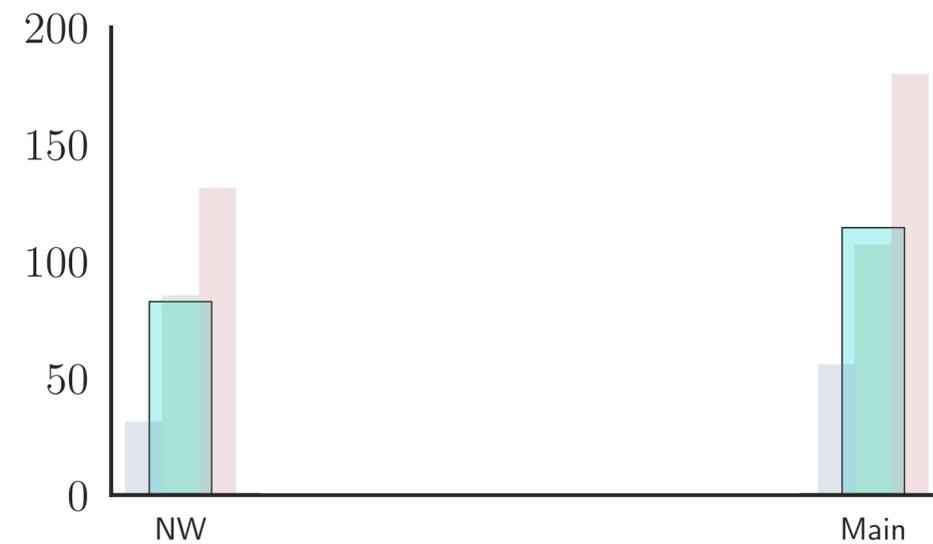
(c) MuJoCo

- Learn well on the MuJoCo datasets.
- Robust and good results on the challenging Adroit and Maze2D datasets.

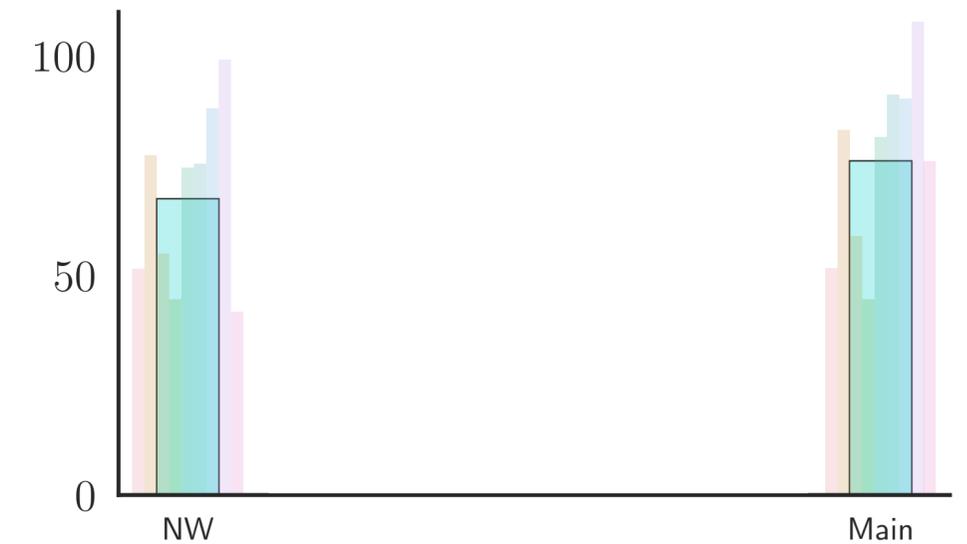
# Ablation Study I: Does weighted model (re)training help?



(a) Adroit



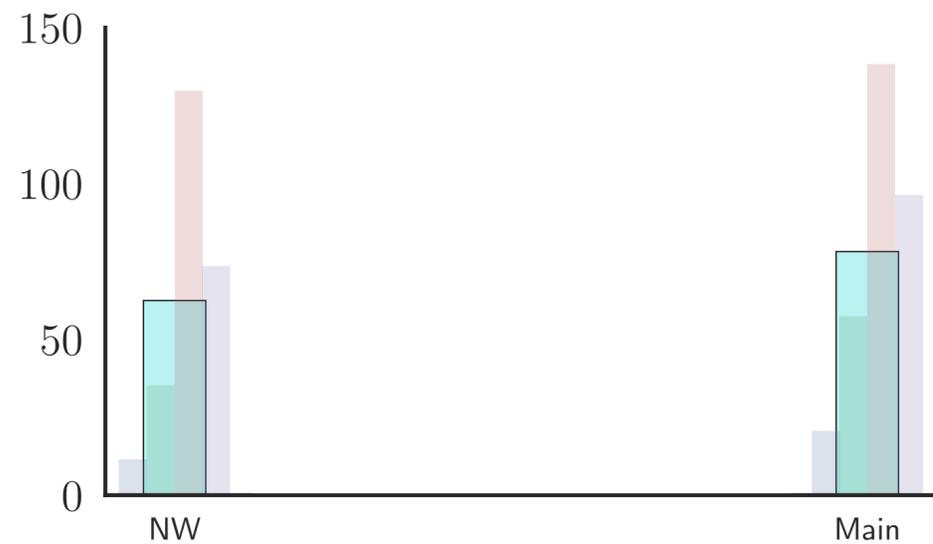
(b) Maze2D



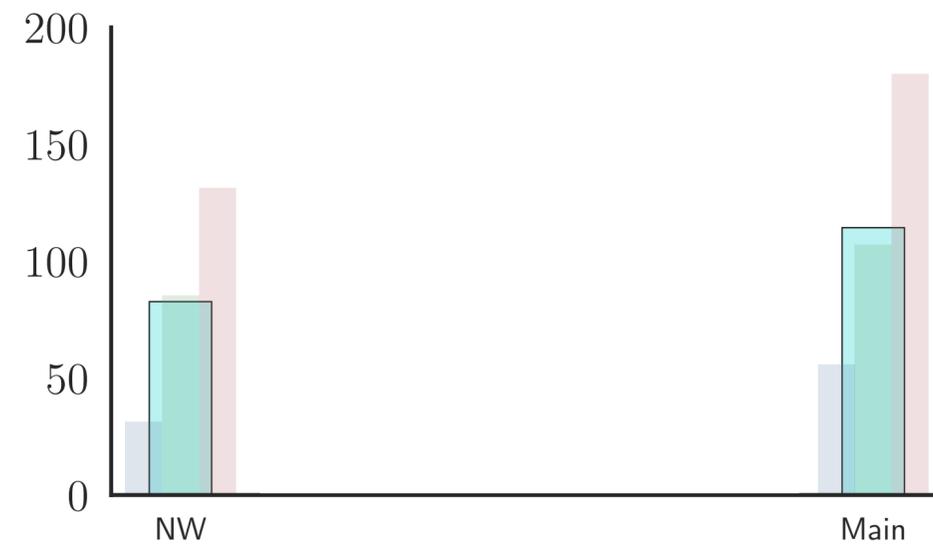
(c) MuJoCo

- Variant: training  only at the beginning using MLE — No Weights (NW).

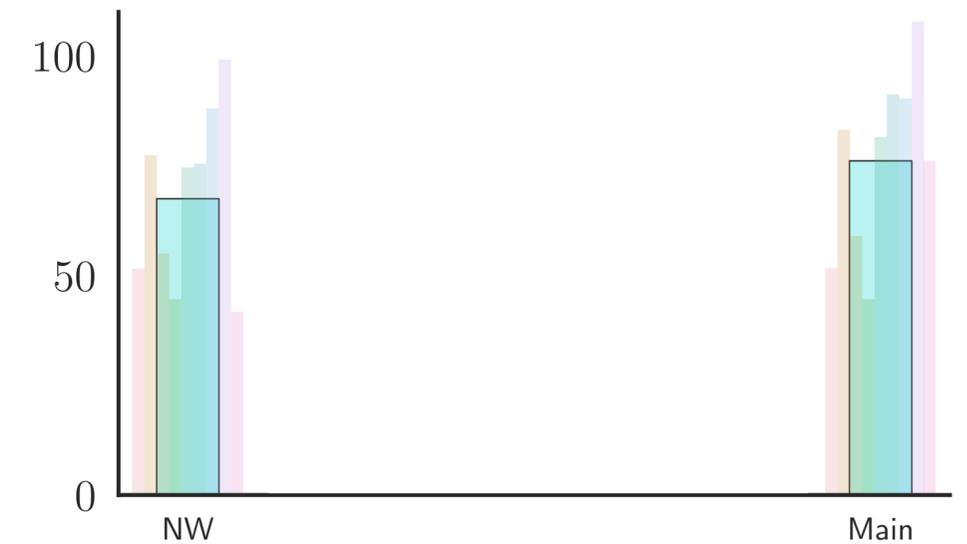
# Ablation Study I: Does weighted model (re)training help?



(a) Adroit



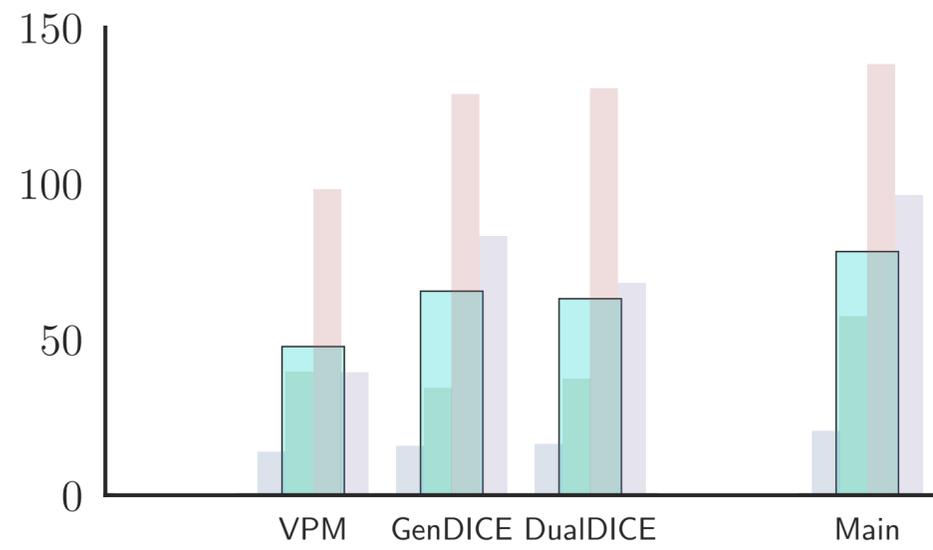
(b) Maze2D



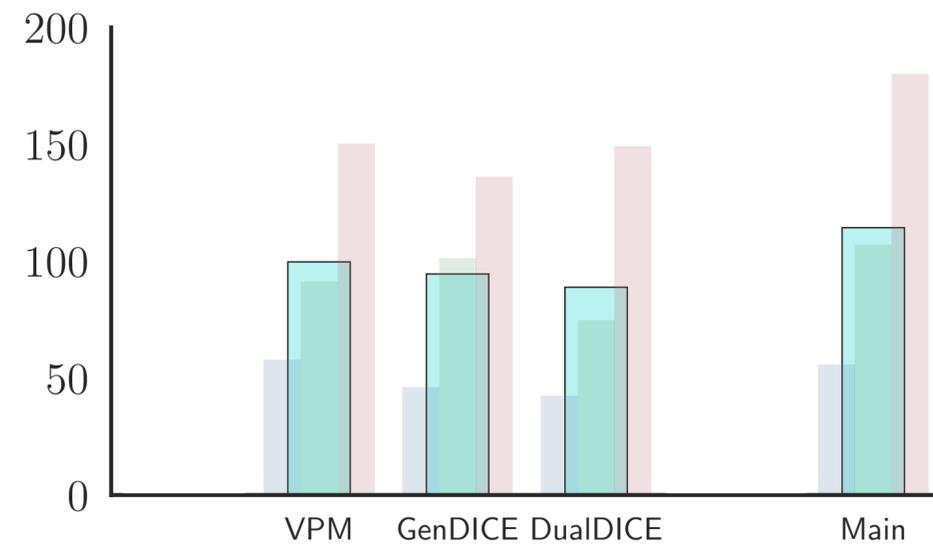
(c) MuJoCo

- Variant: training  only at the beginning using MLE — No Weights (NW).
- On all three domains, the NW variant generally underperforms the main method.

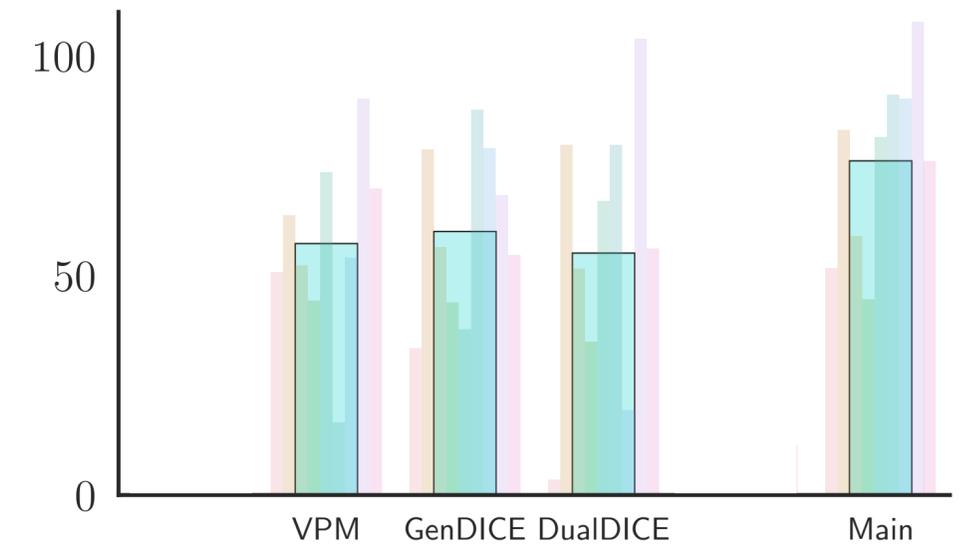
# Ablation Study II: Other density-ratio estimation methods?



(a) Adroit



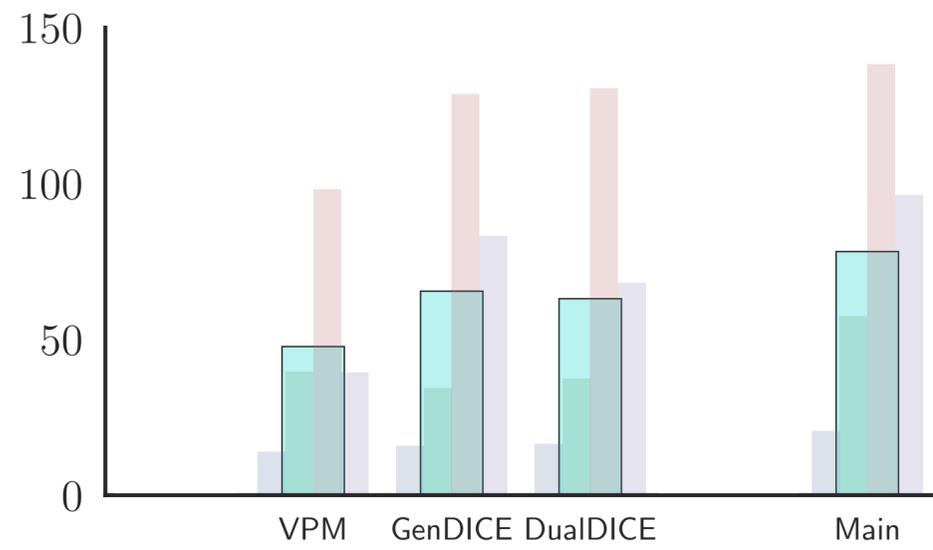
(b) Maze2D



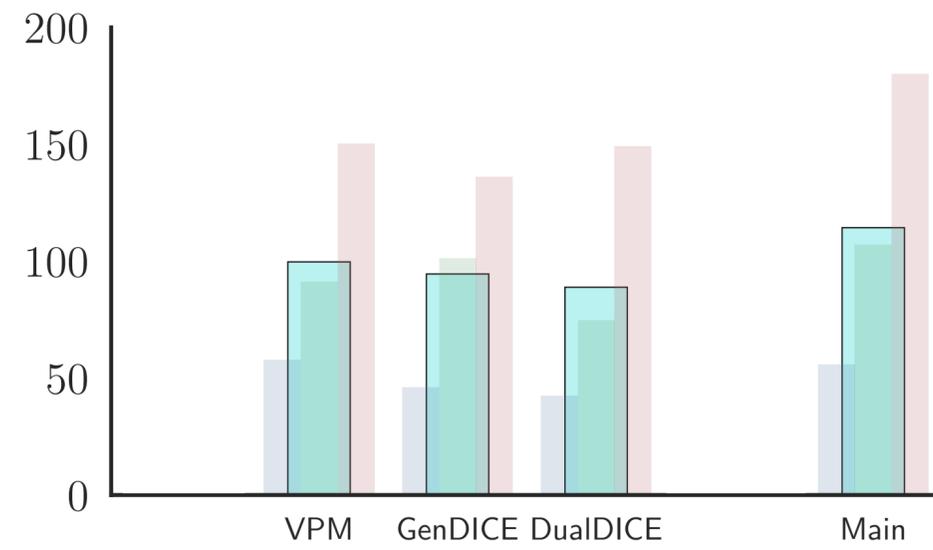
(c) MuJoCo

- Variant:  $\omega(s, a)$  is estimated by VPM, GenDICE, and DualDICE.

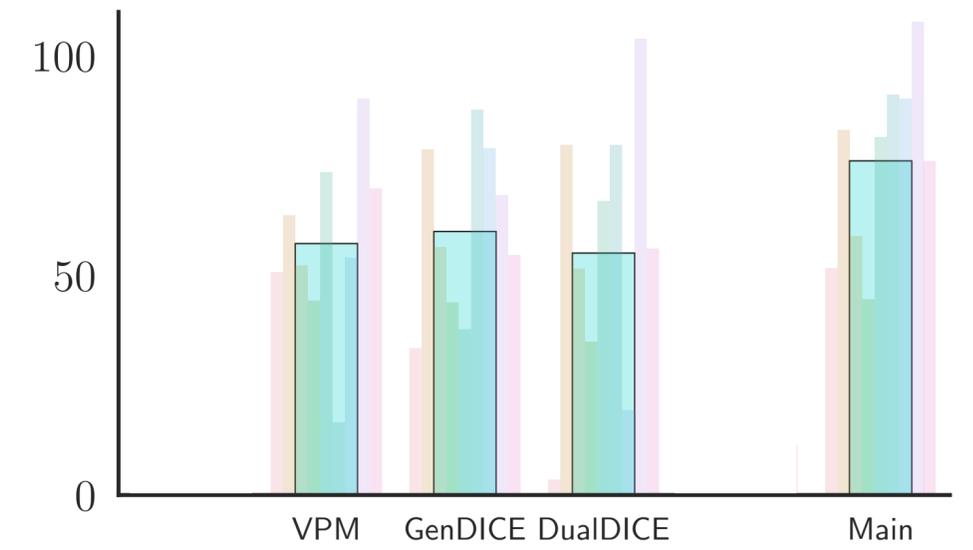
# Ablation Study II: Other density-ratio estimation methods?



(a) Adroit



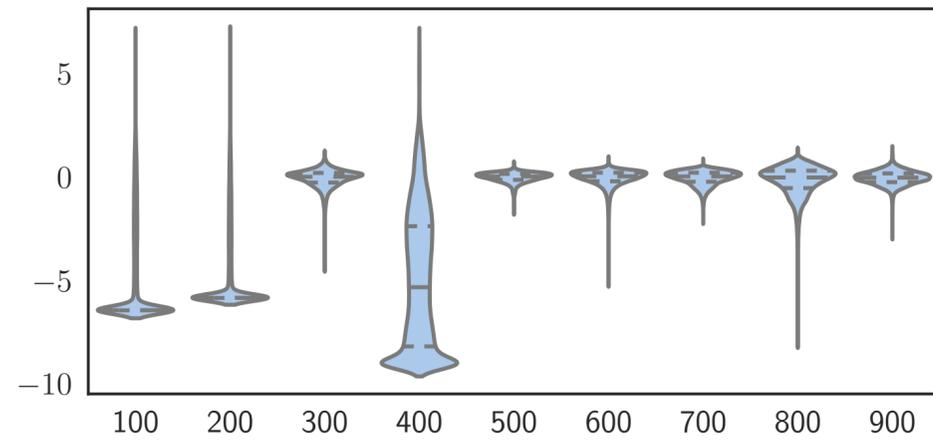
(b) Maze2D



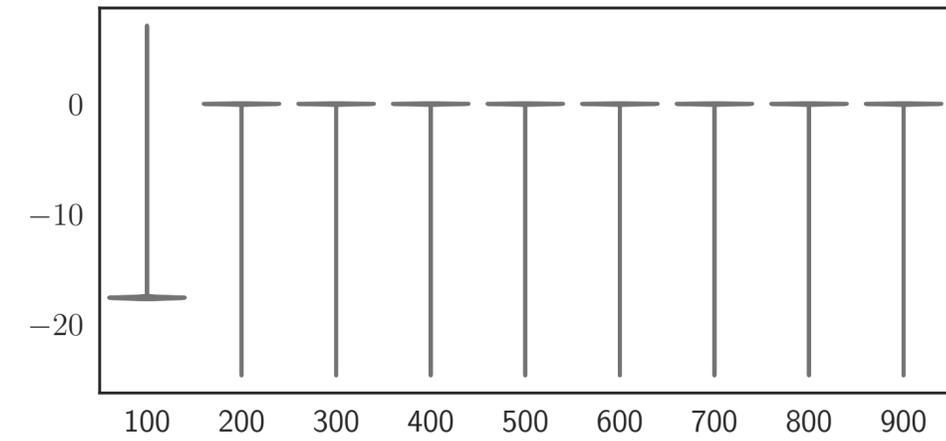
(c) MuJoCo

- Variant:  $\omega(s, a)$  is estimated by VPM, GenDICE, and DualDICE.
- On all three domains, these three variants generally underperform our method.

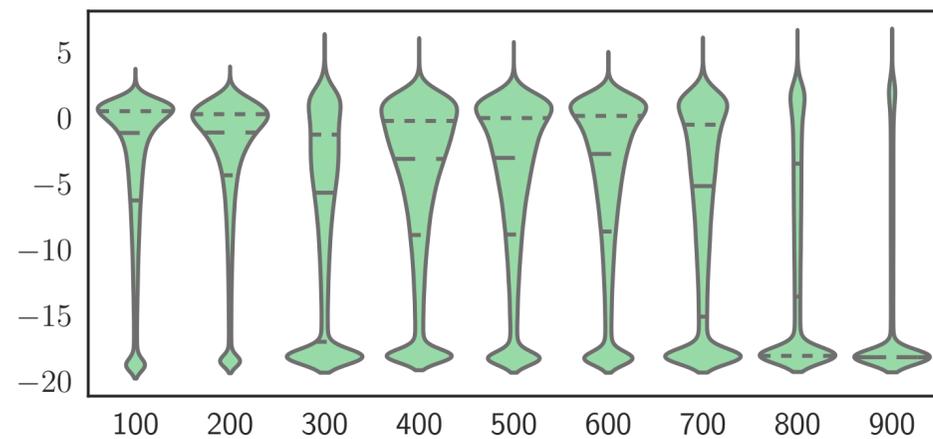
# Ablation Study II: Other density-ratio estimation methods?



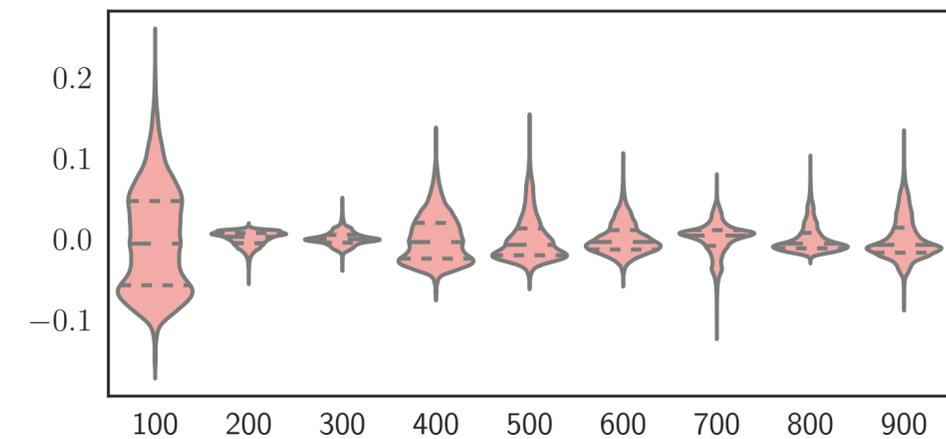
(a) VPM



(b) GenDICE



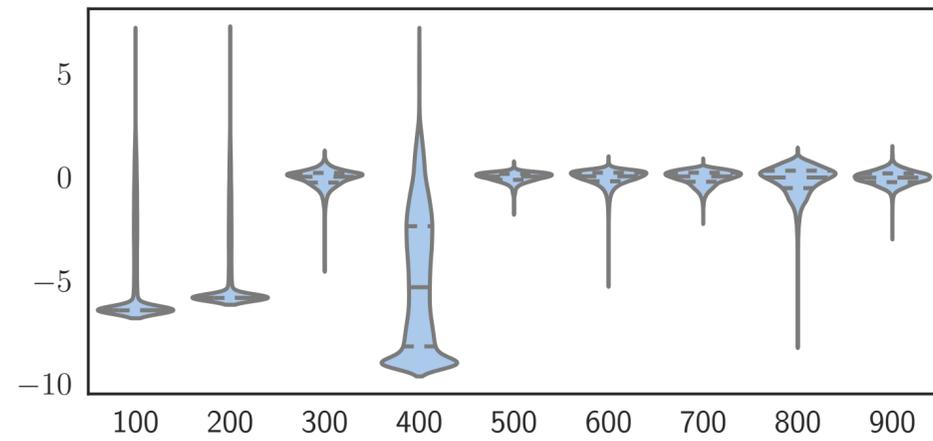
(c) DualDICE



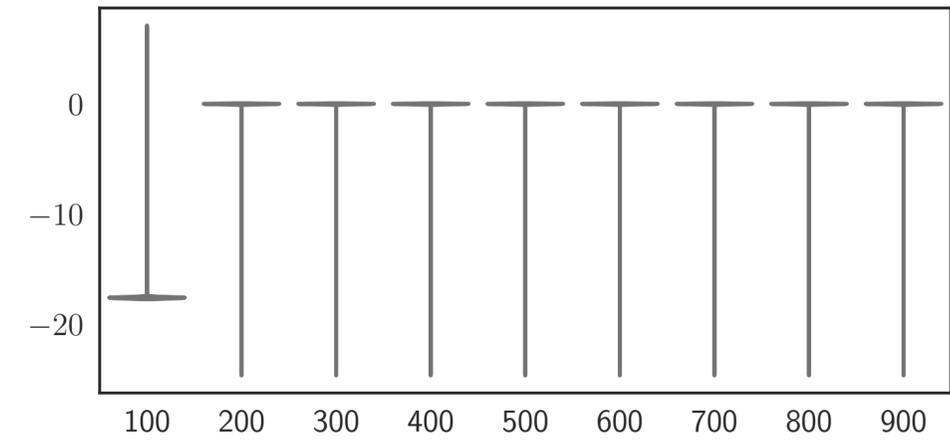
(d) Ours

- Distribution plot of  $\log(\omega(s, a))$  during the training process, on “walker2d-medium-replay.”

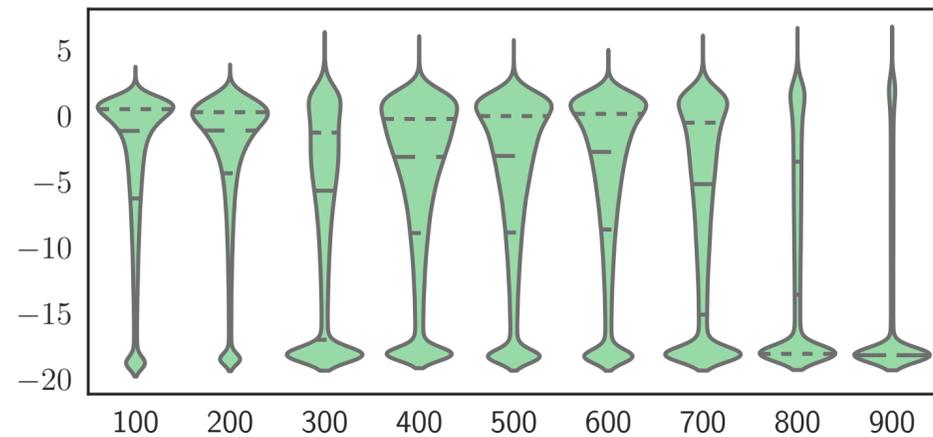
# Ablation Study II: Other density-ratio estimation methods?



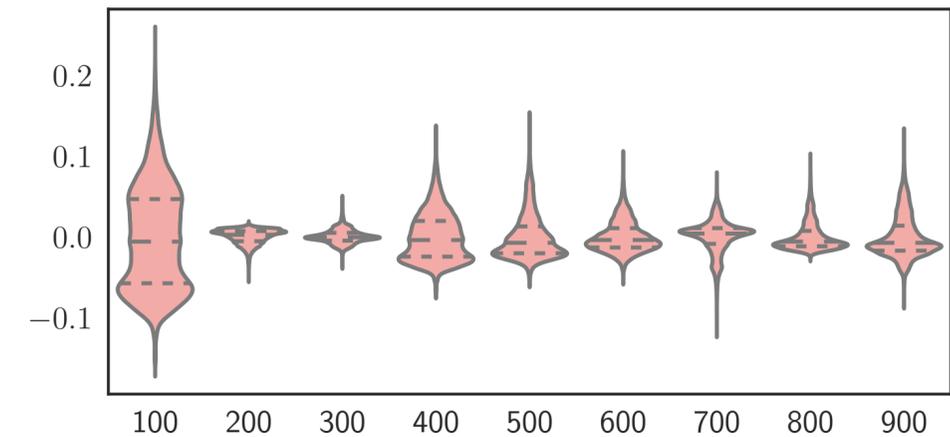
(a) VPM



(b) GenDICE



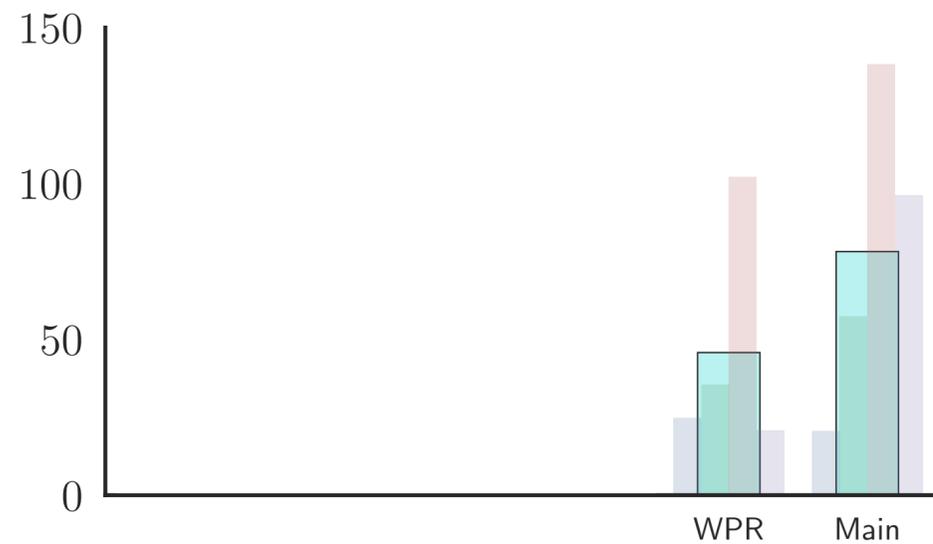
(c) DualDICE



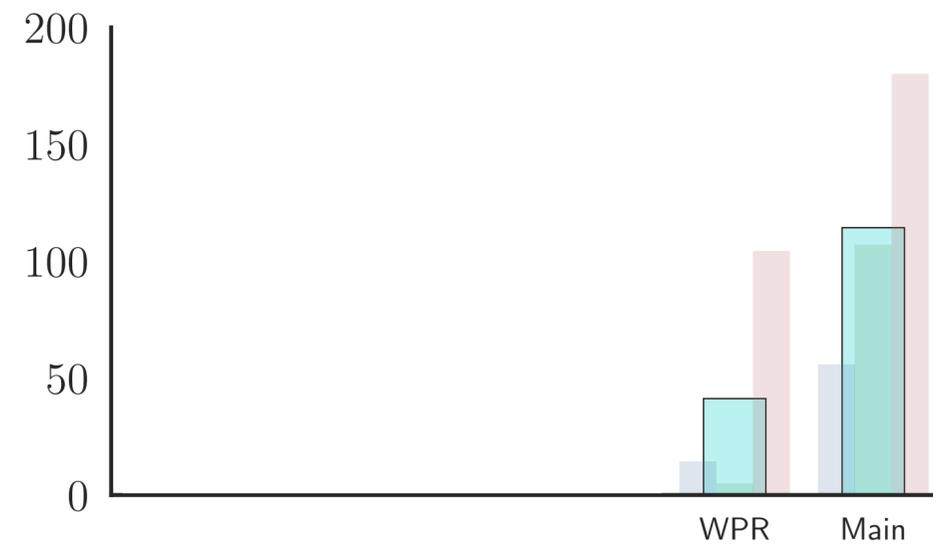
(d) Ours

- Distribution plot of  $\log(\omega(s, a))$  during the training process, on “walker2d-medium-replay.”
- Three alternatives can be unstable to provide good density-ratio for  training.

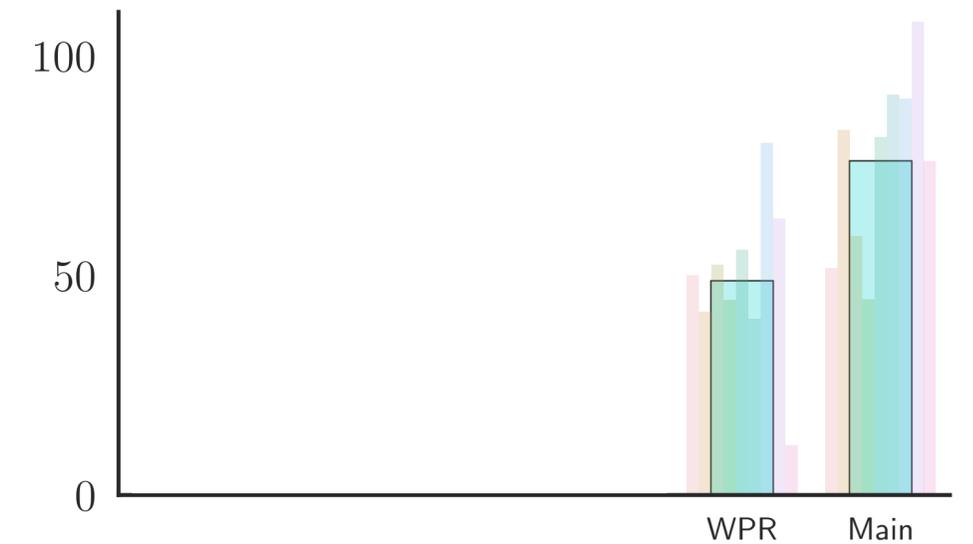
# Ablation Study III: A weighted policy regularizer?



(a) Adroit



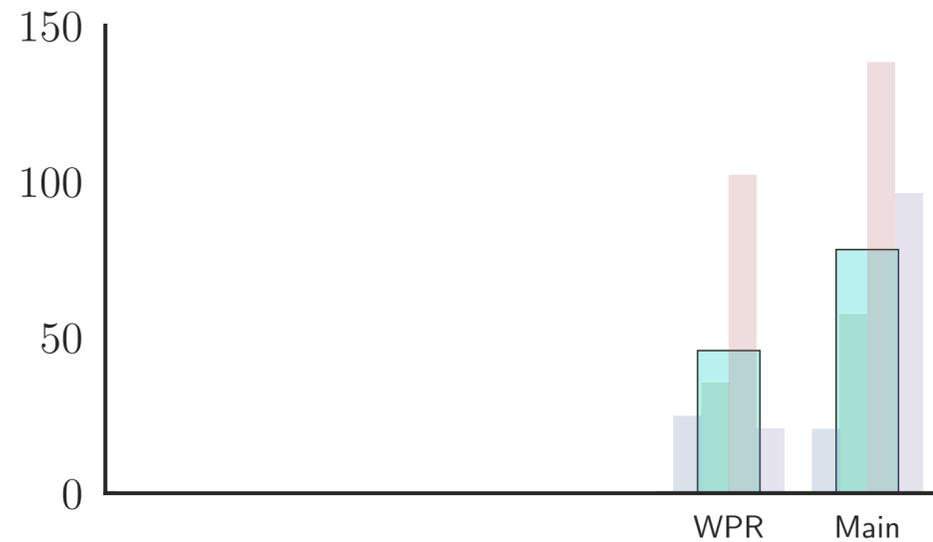
(b) Maze2D



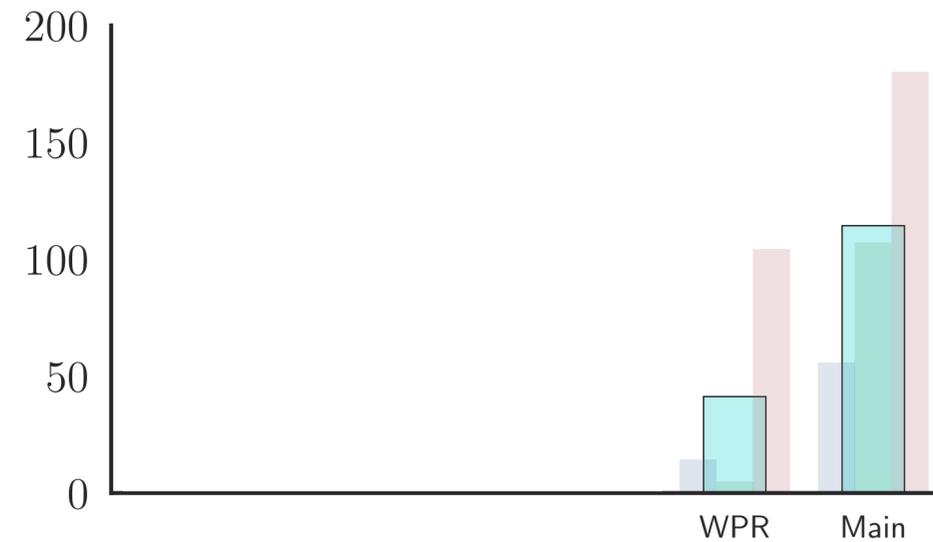
(c) MuJoCo

- Variant: policy regularizer is weighted by the density ratio  $\omega(s, a)$  (WPR).

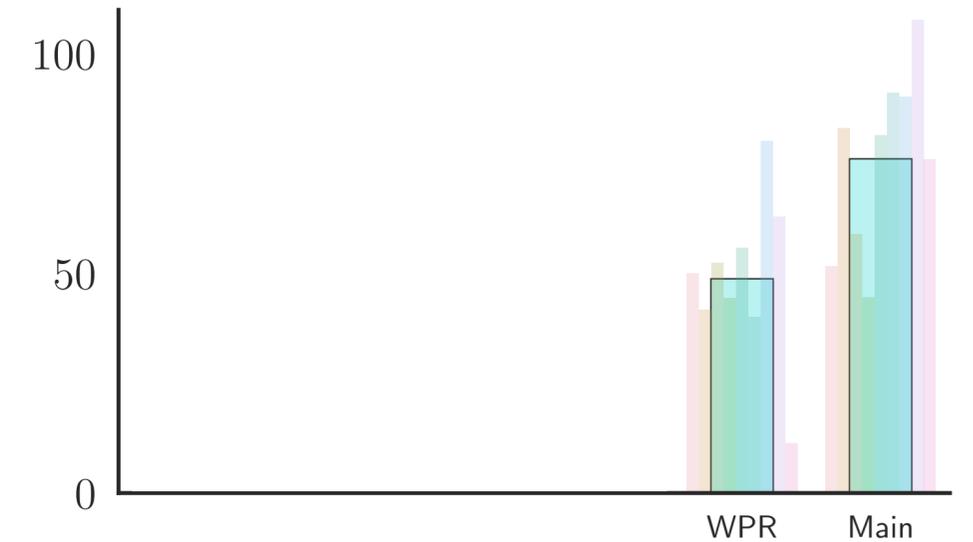
# Ablation Study III: A weighted policy regularizer?



(a) Adroit

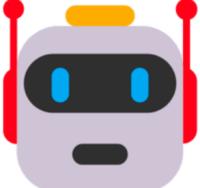


(b) Maze2D



(c) MuJoCo

- Variant: policy regularizer is weighted by the density ratio  $\omega(s, a)$  (WPR).

- Additional instability in training   $\implies$  underperform!

# Summary

- **Goal:** close the mismatched model objectives in offline MBRL.
- **Method:** offline **A**lternating **M**odel-**P**olicy **L**earning.

*QR code for the full paper!*



*QR code for the GitHub Repo!*

