

# Zero-Shot Reinforcement Learning from Low Quality Data

NEURIPS 2024

Scott Jeen <sup>$\alpha$</sup> , Tom Bewley <sup>$\beta$</sup>  & Jonathan M. Cullen <sup>$\alpha$</sup>

<sup>$\alpha$</sup>  University of Cambridge  <sup>$\beta$</sup>  University of Bristol

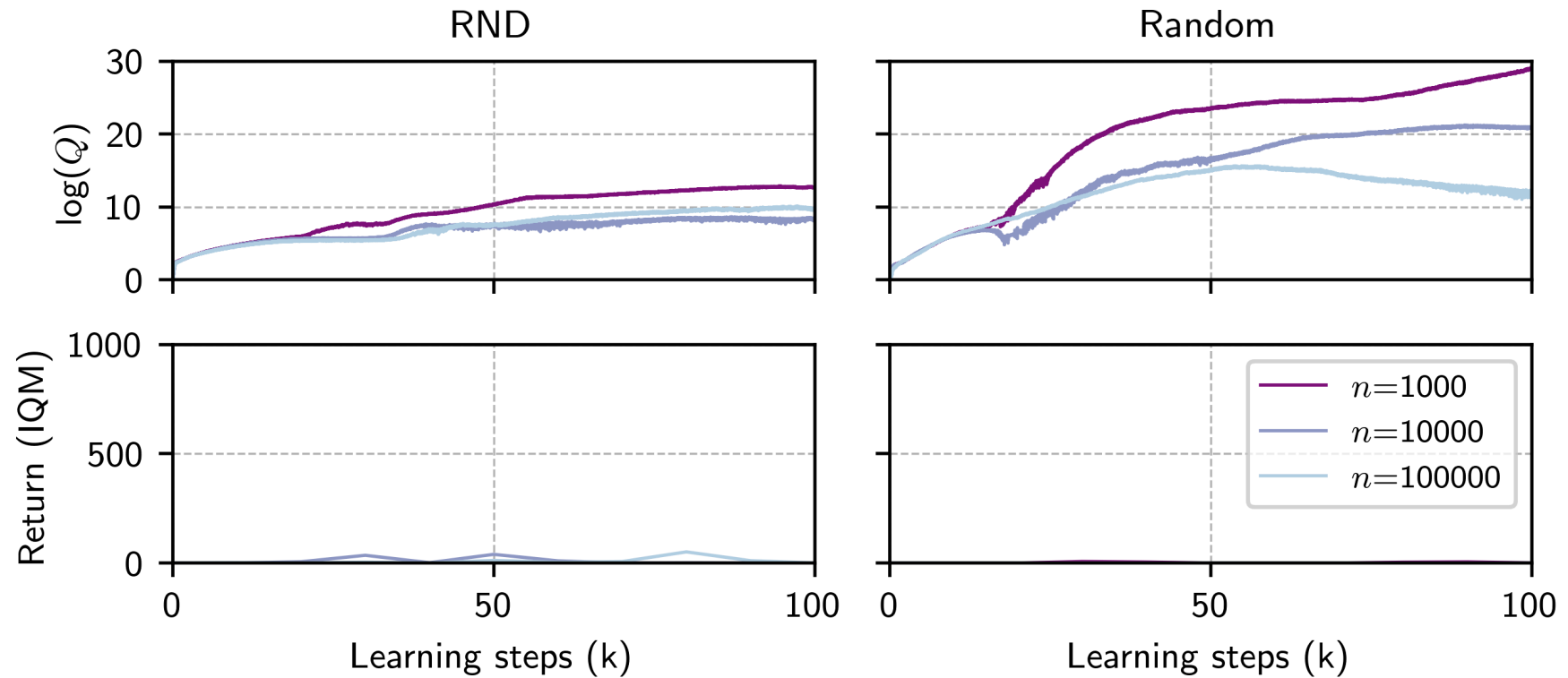
# Motivation

- Training policies to (zero-shot) generalise to unseen tasks in an environment is hard! [1]
- Behaviour Foundation Models (BFMs) based on forward-backward representations (FB) [2] and universal successor features (USF) [3], provide principled mechanisms for performing zero-shot task generalisation
- However, BFMs assumed access to idealised (large & diverse) pre-training datasets that we can't expect for real problems
- **Can we pre-train BFMs on realistic (small & narrow) datasets?**

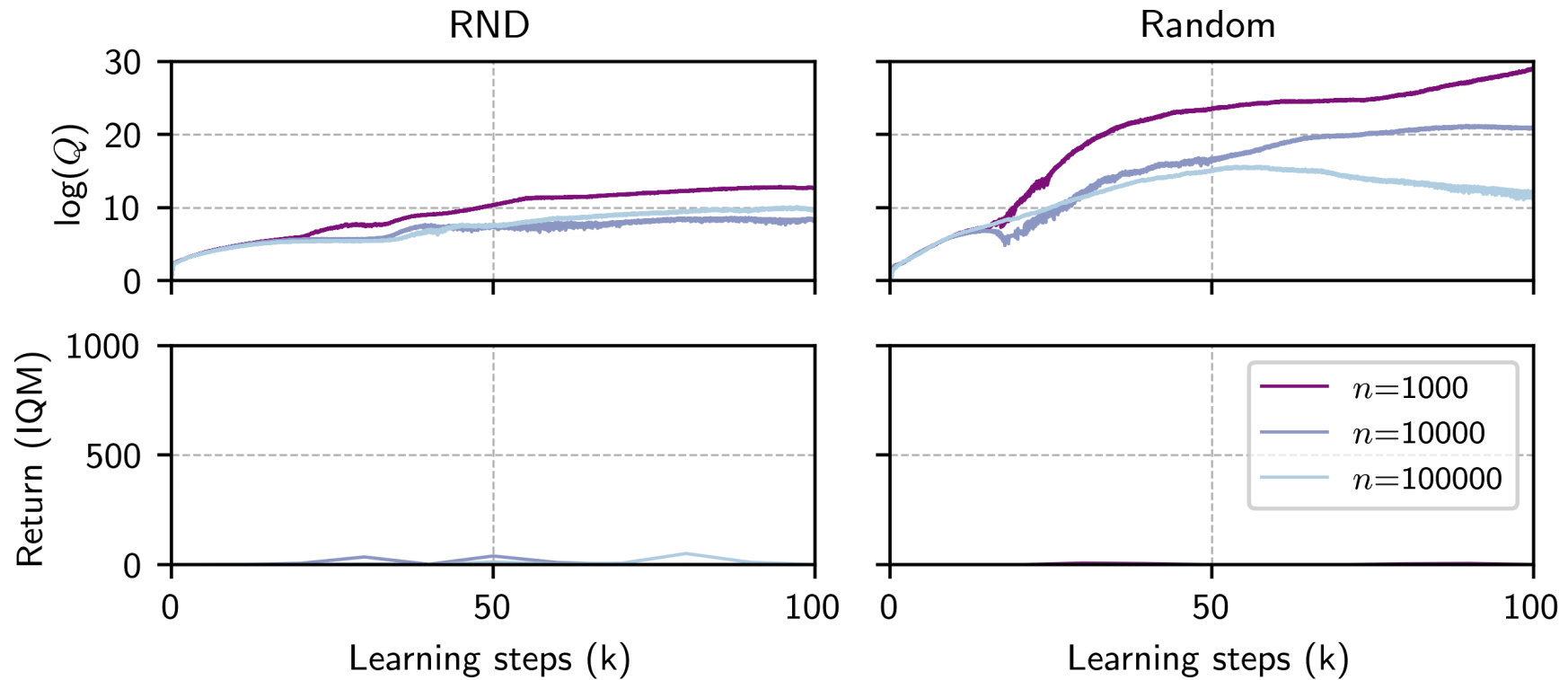
[1] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. *A survey of zero-shot generalisation in deep reinforcement learning*. JAIR 2023

[2] Ahmed Touati, Jérémy Rapin, and Yann Ollivier. *Does zero-shot reinforcement learning exist?* ICLR 2023  
<https://enjeener.io/projects/zero-shot-rl/>

# Out-of-distribution Value Overestimation in BFM

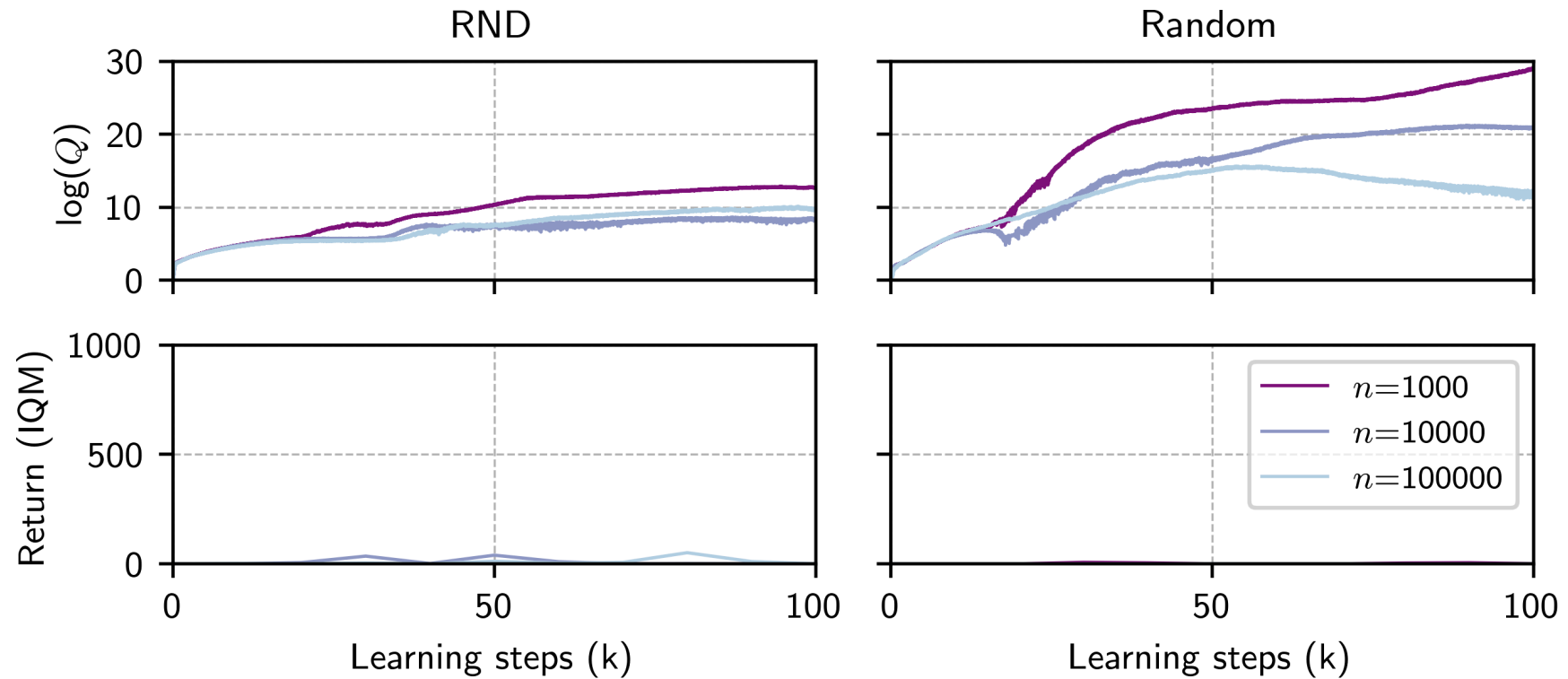


# Out-of-distribution Value Overestimation in BFM



$$\mathcal{L}_{\text{FB}} = \mathbb{E}_{(s_t, a_t, s_{t+1}, s_+) \sim \mathcal{D}, z \sim \mathcal{Z}} [(F(s_t, a_t, z)^\top B(s_+) - \gamma \bar{F}(s_{t+1}, \pi_z(s_{t+1}), z)^\top \bar{B}(s_+))^2 - 2F(s_t, a_t, z)^\top B(s_{t+1})]$$

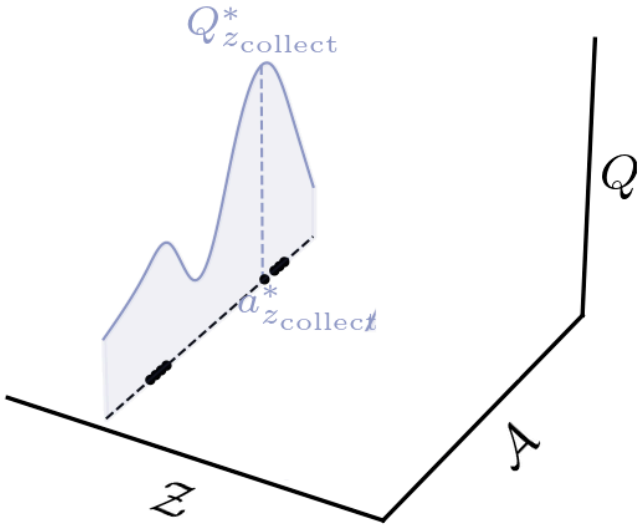
# Out-of-distribution Value Overestimation in BFM



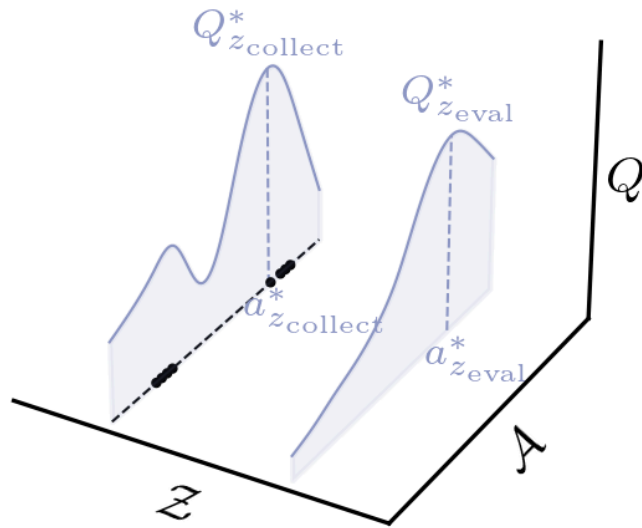
$$\mathcal{L}_{\text{FB}} = \mathbb{E}_{(s_t, a_t, s_{t+1}, s_+) \sim \mathcal{D}, z \sim \mathcal{Z}} [(F(s_t, a_t, z)^\top B(s_+) - \gamma \bar{F}(s_{t+1}, \underbrace{\pi_z(s_{t+1})}_{\text{OOD}}, z)^\top \bar{B}(s_+))^2 - 2F(s_t, a_t, z)^\top B(s_{t+1})]$$

# *Conservative* BFM

# Conservative BFM

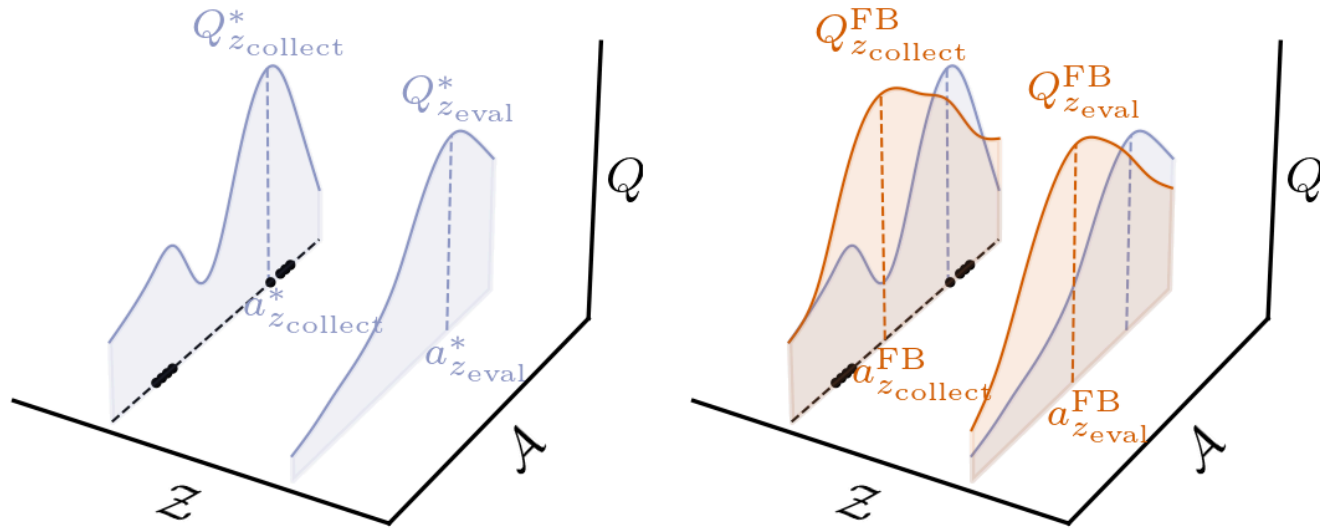


# Conservative BFM

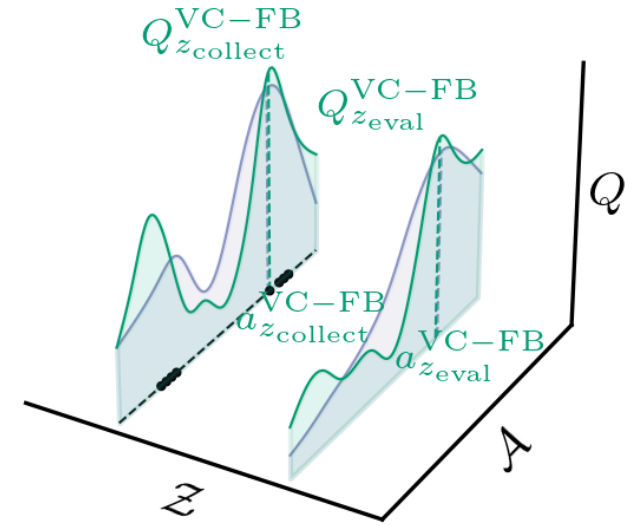
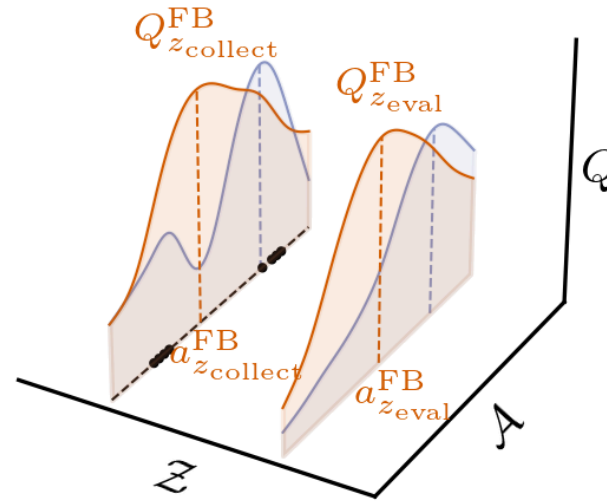
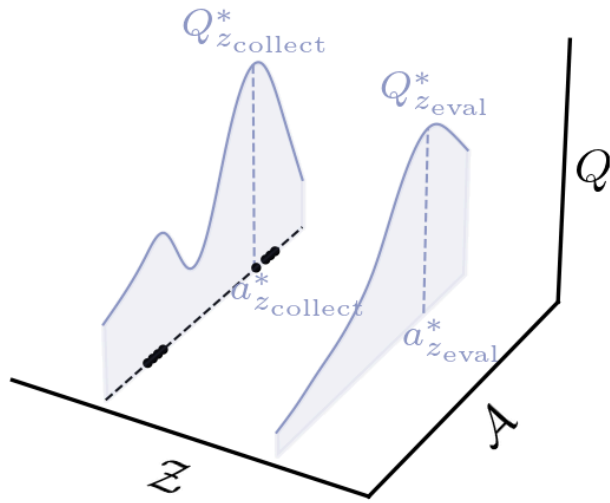




# Conservative BFM



# Conservative BFM



$$\mathcal{L}_{\text{VC-FB}} = \alpha \cdot (\mathbb{E}_{s \sim \mathcal{D}, a \sim \mu(a|s), z \sim \mathcal{Z}} [F(s, a, z)^\top z] - \mathbb{E}_{(s,a) \sim \mathcal{D}, z \sim \mathcal{Z}} [F(s, a, z)^\top z] - \mathcal{H}(\mu)) + \mathcal{L}_{\text{FB}}$$

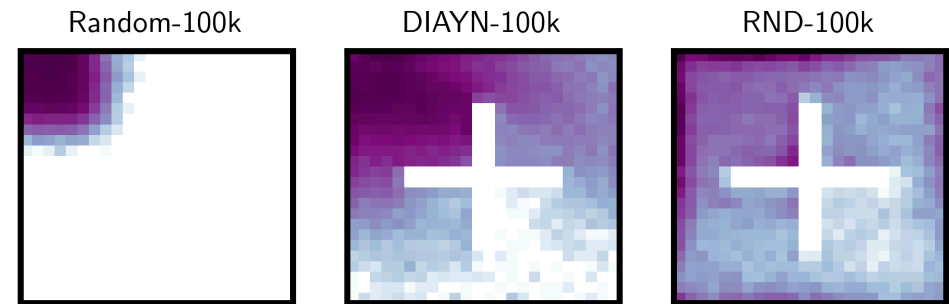
$$\mathcal{L}_{\text{MC-FB}} = \alpha \cdot (\mathbb{E}_{s \sim \mathcal{D}, a \sim \mu(a|s), z \sim \mathcal{Z}, s_+ \sim \mathcal{D}} [F(s, a, z)^\top B(s_+)] - \mathbb{E}_{(s,a) \sim \mathcal{D}, z \sim \mathcal{Z}, s_+ \sim \mathcal{D}} [F(s, a, z)^\top B(s_+)] - \mathcal{H}(\mu)) + \mathcal{L}_{\text{FB}}$$

# ExORL Results

## Baselines

- Zero-shot RL: FB, SF-LAP [5]
- Goal-conditioned RL: GC-IQL [6]
- Offline RL: CQL [7]

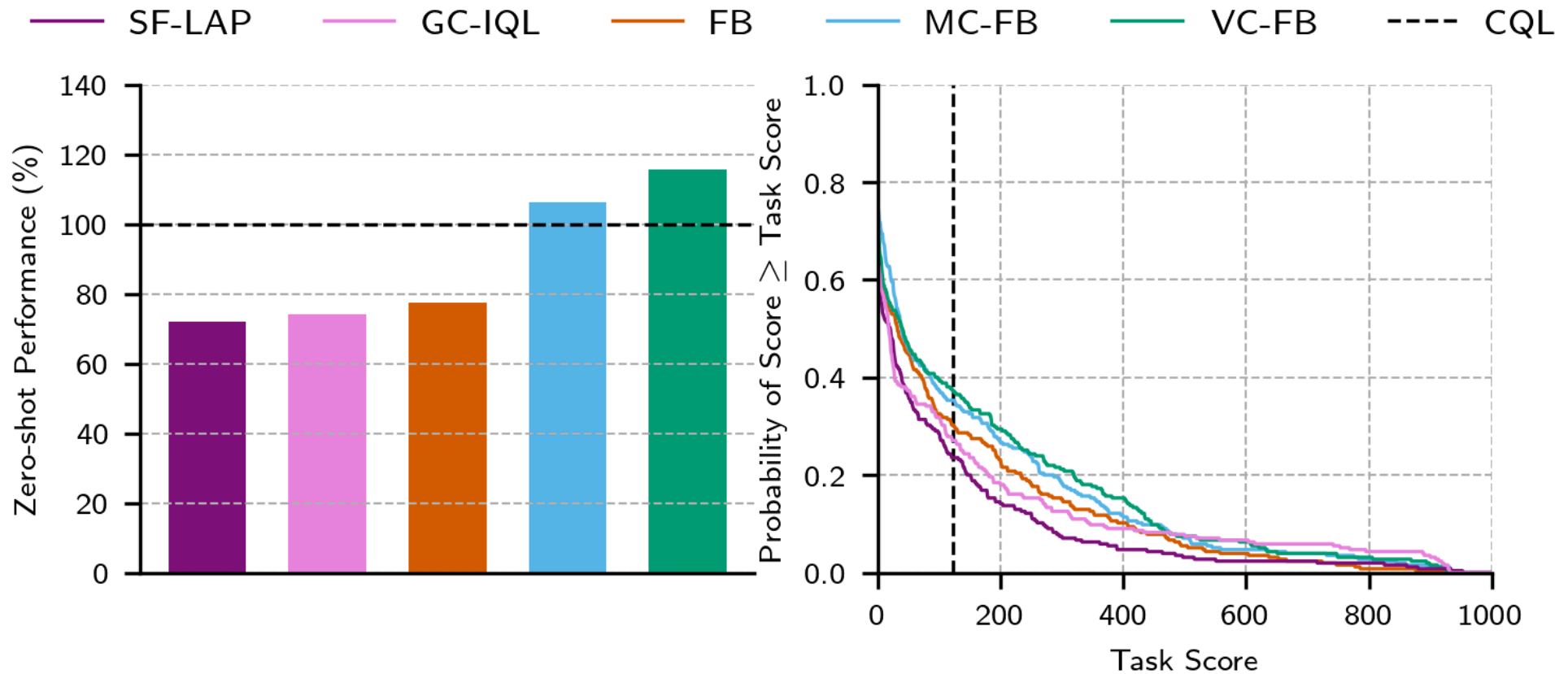
## Datasets



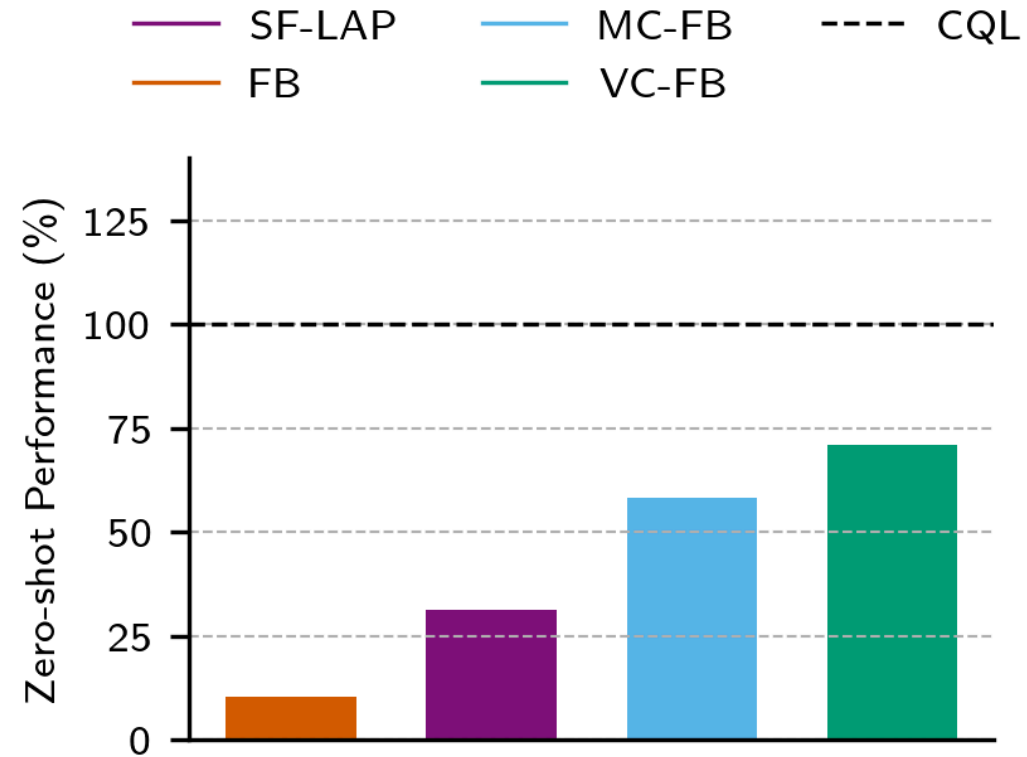
[5] Ahmed Touati, Jérémy Rapin, and Yann Ollivier. *Does zero-shot reinforcement learning exist?* ICLR 2023

[6] Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. *Hiql: Offline goalconditioned rl with latent states as actions.* NeurIPS 2023.

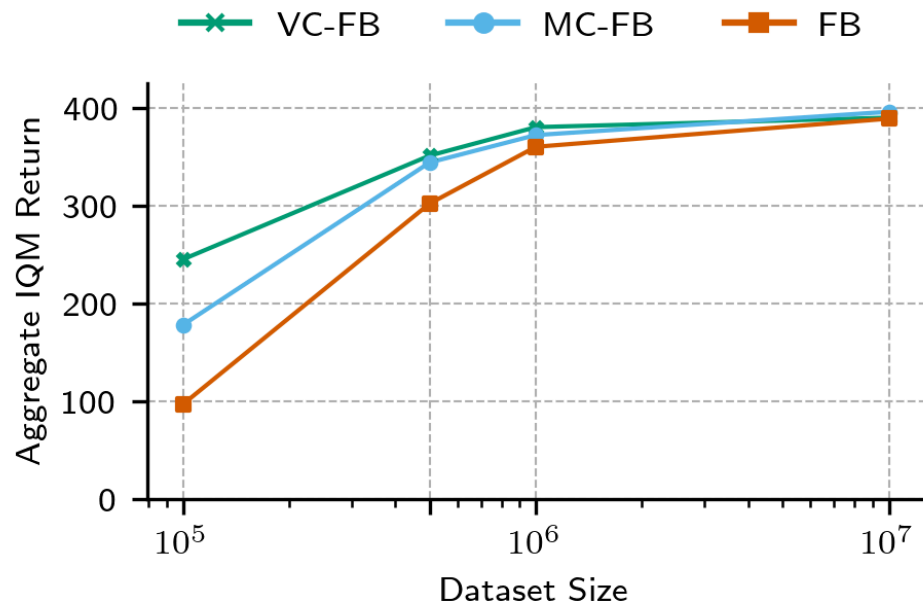
# ExORL Results



# D4RL Results



# Performance on Idealised Datasets is Unaffected



Dataset	Domain	Task	FB	VC-FB	MC-FB
RND	all	all	389	390	396
DIAYN	all	all	269	280	283
RANDOM	all	all	111	131	133
ALL	all	all	256	267	<b>271</b>

# Conclusions

- Like standard offline RL methods, BFM<sub>s</sub> suffer from the *distribution shift*
- As a resolution, we introduce *Conservative* BFM<sub>s</sub>
- *Conservative* BFM<sub>s</sub> considerably outperform standard BFM<sub>s</sub> on low-quality datasets
- *Conservative* BFM<sub>s</sub> do not compromise performance on idealised datasets

Project page:



Twitter/X: [@enjeeneer](https://twitter.com/enjeeneer)

Website: <https://enjeeneer.io>