
Can Active Sampling Reduce Causal Confusion in Offline Reinforcement Learning?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Causal confusion is a phenomenon where an agent learns a policy that reflects
2 *imperfect spurious correlations* in the data. The resulting causally confused be-
3 haviors may appear desirable during training but may fail at deployment. This
4 problem gets exacerbated in domains such as robotics with potentially large gaps
5 between open- and closed-loop performance of an agent. In such cases, a causally
6 confused model may appear to perform well according to open-loop metrics but fail
7 catastrophically when deployed in the real world. In this paper, we conduct the first
8 study of causal confusion in offline reinforcement learning and hypothesise that
9 selectively sampling data points that may help disambiguate the underlying causal
10 mechanism of the environment may alleviate causal confusion. To investigate this
11 hypothesis, we consider a set of simulated setups to study causal confusion and
12 the ability of active sampling schemes to reduce its effects. We provide empirical
13 evidence that random and active sampling schemes are able to consistently reduce
14 causal confusion as training progresses and that active sampling is able to do so
15 more efficiently than random sampling.

16 1 Introduction

17 Offline learning offers opportunities to scale reinforcement learning to domains where offline data
18 is plentiful but online interaction with the environment is costly. The fundamental challenge of
19 offline reinforcement learning is to identify cause and effect of actions from a fixed dataset, which is
20 often intractable. In the absence of online interactions, our hope is that the dataset covers a uniform
21 distribution of an exhaustive set of plausible scenarios. This is often not the case in datasets for robotic
22 control, which are long-tailed and often contain only a handful of samples for rare (and informative)
23 events. Causal confusion occurs when agents misinterpret the underlying causal mechanisms of the
24 environment and erroneously associate certain actions or states with a given reward. For example,
25 if an agent happens to simultaneously observe independent events X and Y in its environment
26 whenever it receives a reward R , and R causally depends on Y but not on X , the agent may attribute
27 the reward R to X and Y occurring jointly even though R may be independent of Y . Problematically,
28 if the spurious correlation between Y and R observed at training time ceases to hold at deployment
29 time, a causally-confused model may show a significant deterioration in performance. Often, spurious
30 correlations are not *perfectly* held in offline data, but optimisation schemes like mini-batched gradient
31 descent can still produce models that latch onto them since they help in optimising the training loss.
32 In this paper, we explore whether causal confusion in offline reinforcement learning from datasets
33 exhibiting causal ambiguity can be alleviated by random or active sampling. We provide empirical
34 evidence that random and active sampling schemes are able to consistently reduce causal confusion
35 and that active sampling is able to do so more efficiently than random sampling.

36 2 Related Work

37 **Causal Confusion in Supervised Learning.** Several works in imitation learning have proposed
38 solutions to mitigate causal confusion, which was first defined in [de Haan et al., 2019]. Wen et al.
39 [2020] proposes adversarial training to prune out any *known* sources of spurious correlations from
40 the policy’s representation, for instance, the previous control commands given to a robot; Wen et al.
41 [2021] propose loss-reweighting of datapoints based on the loss of a model trained with just the
42 spurious correlates as the input; OREO [Park et al., 2021] regularises the model’s representation to
43 be invariant to any individual object being dropped out in a scene. Lee et al. [2022] propose training
44 a diversified ensemble in the case when *perfect spurious correlations* exist in the data and later select
45 from these hypotheses based on validation data. Causal Confusion has also recently been studied
46 in reward-learning from preferences [Tien et al., 2022], where spurious correlations can be drawn
47 between a human evaluator’s preferences and certain actions or parts of the state space.

48 **Ensemble Models in RL.** Ensembles have been studied extensively to guide exploration in online
49 RL [Osband et al., 2016] [Lee et al., 2021], and recently to construct adaptive pessimism constraints
50 in offline RL, to disincentivise uncertain actions from having high estimated returns. Recent work
51 [An et al., 2021] showed that increasing the size and diversity of the ensembled critic in Soft-Actor-
52 Critic [Haarnoja et al., 2018] performs competitively with state-of-the-art offline RL algorithms.
53 However prior work hasn’t explored how the uncertainty from ensembles could be used to sample
54 transitions in RL. Prioritised replay [Schaul et al., 2015] is a sampling scheme based on the TD-error
55 of transitions, that was proposed in off-policy RL but hasn’t been studied in offline RL.

56 **AI Alignment.** AI alignment seeks to align the behavior of agents with the intentions of their
57 creators by investigating the incentives behind demonstrated tasks. Recent work on *Goal Misgen-*
58 *eralisation* [Langosco et al., 2022] explores how online RL agents in Procgen [Cobbe et al., 2019]
59 can get confused about the goal they’re pursuing since those goals co-occur with irrelevant artifacts
60 in the environment most of the time. In this case the specification is correct, but the agent still
61 pursues an unintended objective (as opposed to poor reward definitions that predictably lead to reward
62 hacking). We build upon an environment introduced in this work to collect data for reproducing the
63 phenomenon of causal confusion in offline RL.

64 3 Alleviating Causal Confusing in Offline RL via Active Sampling

65 3.1 Offline Reinforcement Learning

66 Offline RL algorithms aim to learn an optimal policy along with estimates of the value (or Q -value)
67 function from a dataset of transitions $\mathcal{D} = \{(s, \mathbf{a}, r, s')\}$ collected by a behaviour policy π_β .

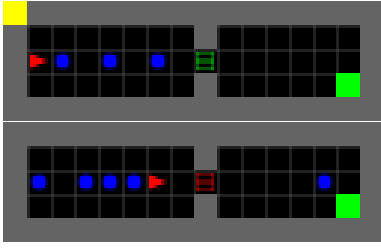
68 **Conservative Q -Learning.** For our experimnts we choose CQL [Kumar et al., 2020] for it’s
69 simplicity and competitive performance. The CQL objective, which combines the standard TD-error
70 of Q -learning with a penalty constraining deviations from the behaviour policy, is defined as:

$$\mathcal{L}_{\text{critic}}^{\text{CQL}}(\theta) = \frac{1}{2} \mathbb{E}_{(s, a, s') \sim \mathcal{D}} \left[(Q_\theta - \mathcal{B}^\pi Q_{\bar{\theta}})^2 \right] + \alpha_0 \mathbb{E}_{s \sim \mathcal{D}} \left[\log \sum_a \exp Q(s, a) - \mathbb{E}_{a \sim \pi_\beta} [Q(s, a)] \right], \quad (1)$$

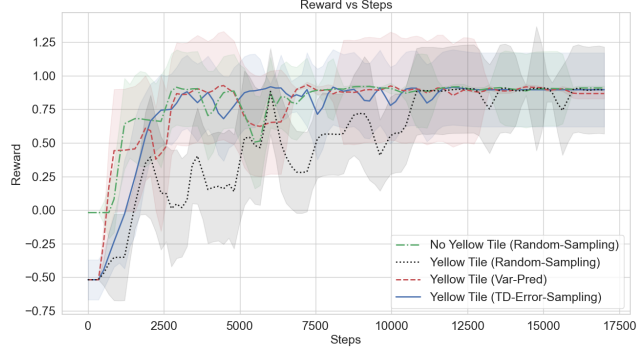
71 where the Bellman operator $\mathcal{B}^\pi Q = r + \gamma P^\pi Q$, and P^π is the transition matrix coupled with the
72 policy π . We model the uncertainty of the learned Q -function parameterised by θ by ensembling the
73 model and training on identical transitions across the ensemble members, with their own corresponding
74 targets ($\bar{\theta}$) as proposed in Ghasemipour et al. [2022]

75 3.2 Active Sampling

76 The focus of this work is on experimenting with data-sampling strategies without making any
77 modifications to the objective. Algorithm 1 describes the CQL setup with active-sampling of
78 transitions, where the modifications from vanilla random-sampling are highlighted in blue. We study
79 the following uncertainty-based and loss-based data acquisition schemes:



(a) **Top:** The leading vehicle is static and the top-left tile flashes yellow since the leading vehicle is static. **Bottom:** The agent is in front of a red light, and the top-left tile isn't yellow since the leading vehicle isn't static or blocked.



(b) Random sampling takes 4x gradient steps to recover the correct solution compared to active sampling (Variance and TD-Error-based) when both are trained on data with the spurious yellow tile.

Figure 1: Traffic-world environment.

80 **Uncertainty about the greedy action (variance-based):** The Q -values of different ensemble
 81 members could have arbitrary numerical offsets but still be equivalent, due to bootstrapping. Instead,
 82 we estimate the uncertainty of actions by computing the variance of their advantage over the ensemble,
 83 where the advantage of an action a^* for a Q -learner can be written as follows:

$$A^\pi(s, a^*) = Q^\pi(s, a^*) - V^\pi(s) \approx Q^\pi(s, a^*) - \sum_a \left[Q(s, a) \cdot \frac{e^{Q(s, a)}}{\sum_{a'} e^{Q(s, a')}} \right] \quad (2)$$

84 **TD-Error (loss-based):** Based on the Temporal Difference error similar to Prioritised Experience
 85 Replay [Schaul et al., 2015]

86 In practice, computing the acquisition scores over all the transitions in the dataset can be very
 87 expensive and redundant since high-error or high-uncertainty points will likely stay informative for
 88 a short window of subsequent gradient steps. We thus recompute the scores after every n gradient
 89 steps, and vary n as a hyper-parameter in our experiments.

90 4 Experiments

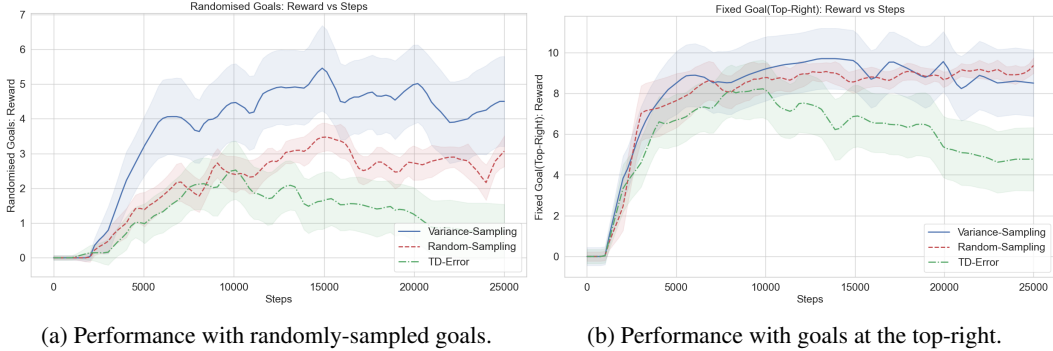
91 We investigate the following questions: (1) Can causal confusion be consistently observed in
 92 CQL when sampling transitions randomly from a long-tailed demonstration dataset? (2) Does
 93 active-sampling based on the (implicit) policy's uncertainty or loss help? (3) What is the compu-
 94 tation time involved with each of these acquisition schemes?
 95
 96
 97
 98

99 4.1 Illustrative Example: Traffic-World

100 The autonomous driving literature cites many ex-
 101 amples where models training on large datasets
 102 are very performant but exhibit causal confusion
 103 on the tail cases of their operational domain, for instance: (1) models stopping at pedestrian crossings
 104 regardless of whether a pedestrian is present or not since the two often co-occur; (2) self-driving
 105 agents that simply try to *cruise* if they know their current speed since expert driving datasets contain
 106 cruising behaviour in a large fraction of each trajectory. We build on the environment proposed in
 107 [Anonymous, 2021] to construct a gridworld (shown in Figure 1a), where an agent (red triangle)
 108 starts at the leftmost point in a row behind leading vehicles (blue circles), and needs to cross a traffic
 109 light to reach a goal location (green square) on the right side of the grid. We collect data such that the
 110 probability of the traffic light turning red becomes lower as the agent approaches it, and so the data
 111 distribution contains (1) mostly episodes where the light is green, (2) some episodes where the traffic

Algorithm 1 Conservative Q -Learning (+ active-sampling)

- 1: Initialise ensemble Q -function Q_θ , n_{ep} =epochs, d_{sz} =dataset size, b_{sz} =batch size, T =steps-per-epoch.
 - 2: **for** epoch e in $\{1, \dots, n_{ep}\}$ **do**
 - 3: **for** step t in $\{1, \dots, T\}$ **do**
 - 4: compute scores acq_i over $\mathcal{D}_{train} = [s_i, a_i]_{i=1}^{d_{sz}}$ according to the acquisition function
 - 5: $acq_i = \frac{acq_i}{\sum_{j=1}^{d_{sz}} acq_j}$ (normalise scores)
 - 6: sample batch $B = [s_i, a_i, s'_i, r_i]_{i=1}^{b_{sz}}$ from $\mathcal{D}_{train} \sim \text{multinomial}(acq)$
 - 7: Train the Q -function on \mathcal{D}_{train} using objective from Equation (1)
 - 8: **end for**
 - 9: **end for**
-



(a) Performance with randomly-sampled goals. (b) Performance with goals at the top-right.

Figure 2: Agents trained on a dataset containing 6000 episodes with a fixed goal and 200 episodes with a randomly sampled goal in Maze. We see that random-sampling and active-sampling perform similarly on the fixed goal evaluation environment (right), but the active-sampling variants achieve higher reward in the environments with randomly sampled goals. This verifies that the model is not just performing well in one of the two kinds of environments, is not constrained by capacity, and the reason behind the lower performance of random-sampling in this case is causal confusion.

122 light is red and the agent is waiting behind the vehicle in front (referenced here onward as the leading
 123 vehicle), and (3) only a couple of episodes where the light turns red with the agent at the front of
 124 the traffic queue. In this setup, the agent could just learn to follow the leading vehicle, instead of
 125 learning traffic light rules. To test causal confusion explicitly here, we introduce a related spurious
 126 correlate: a flashing yellow tile at the top left of the grid, that is yellow whenever the leading vehicle
 127 is stopped or blocked, and grey otherwise. The agent could follow this as an indicator of whether to
 128 stop or go ahead, and this policy would be correct for 98% of the data points. Figure 1b shows the
 129 training curves of CQL agents trained with randomly-sampled data, with and without the yellow tile
 130 present in images in the dataset - we see that the performance of the former agent degrades and it
 131 takes 4x the number of gradient steps to converge to the solution of the latter agent which is trained
 132 without the spurious correlate present. Also shown are the active sampling variants trained with the
 133 spurious yellow tile, which perform very similarly to random-sampling when the spurious correlate is
 134 not present.

125 **4.2 Generalization in Offline RL: Progen**

126 The Maze environment in Progen [Cobbe et al., 2019] defines a navigation task where the agent
 127 starts at the bottom left in the maze and receives a reward of +10 upon successfully reaching the goal
 128 which is sampled at any valid location in the maze. [Langosco et al., 2022] recently showed that
 129 an agent trained on a series of environments with the goal always at the top-right will be causally
 130 confused about the source of the reward. It will still navigate to the top-right even when the goal is
 131 sampled elsewhere. We generate a skewed *mixture* dataset containing mostly episodes where the goal
 132 is sampled at the top-right, and a few episodes where the goal is sampled randomly. Further details
 133 about the setup described in the Appendix. Figure 2b shows the evaluation performance of random
 134 and active-sampling agents trained on the *mixture* dataset, when goals are sampled randomly in the
 135 evaluation environment. We plot the computation time for the random and active-sampling variants
 136 in Figure 3 in the Appendix. Qualitative evaluations show that agents which achieve lower reward
 137 still successfully navigate to the top-right corner of the maze.

138 **5 Conclusions**

139 In this paper we designed preliminary setups to study causal confusion in offline RL, which occurs
 140 when a policy is learnt with random sampling of data from a skewed offline dataset. We designed
 141 uncertainty-based and loss-based data sampling baselines to selectively sample transitions for training,
 142 and saw promising evidence that active sampling can recover a less causally-confused model in
 143 significantly fewer training steps as compared to random-sampling. An interesting line of future
 144 work would be to scale this up to larger benchmarks, and extend this analysis to the case when
 145 acquisition scores for active sampling aren't computed for all the data at once, instead maintaining an
 146 approximation through running scores as is done in [Schaul et al., 2015].

147 **References**

- 148 Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline
149 reinforcement learning with diversified q-ensemble. In A. Beygelzimer, Y. Dauphin, P. Liang, and
150 J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL
151 <https://openreview.net/forum?id=ZUvaSo1QZh3>.
- 152 Anonymous. Resolving causal confusion in reinforcement learning via robust exploration. In
153 *Self-Supervision for Reinforcement Learning Workshop - ICLR 2021*, 2021. URL [https://](https://openreview.net/forum?id=DKCXncD4Xtq)
154 openreview.net/forum?id=DKCXncD4Xtq.
- 155 Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation
156 to benchmark reinforcement learning. *arXiv preprint arXiv:1912.01588*, 2019.
- 157 Pim de Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learn-
158 ing. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Gar-
159 nett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran
160 Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/file/](https://proceedings.neurips.cc/paper/2019/file/947018640bf36a2bb609d3557a285329-Paper.pdf)
161 [947018640bf36a2bb609d3557a285329-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/947018640bf36a2bb609d3557a285329-Paper.pdf).
- 162 Seyed Kamyar Seyed Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic?
163 estimating uncertainties for offline RL through ensembles, and why their independence matters.,
164 2022. URL <https://openreview.net/forum?id=wQ7RCayXUS1>.
- 165 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
166 maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy and Andreas
167 Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80
168 of *Proceedings of Machine Learning Research*, pages 1861–1870. PMLR, 10–15 Jul 2018. URL
169 <https://proceedings.mlr.press/v80/haarnoja18b.html>.
- 170 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for
171 offline reinforcement learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and
172 H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1179–1191.
173 Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/0d2b2061826a5df3221116a5085a6052-Paper.pdf)
174 [0d2b2061826a5df3221116a5085a6052-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/0d2b2061826a5df3221116a5085a6052-Paper.pdf).
- 175 Lauro Langosco Di Langosco, Jack Koch, Lee D Sharkey, Jacob Pfau, and David Krueger. Goal mis-
176 generalization in deep reinforcement learning. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,
177 Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International*
178 *Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,
179 pages 12004–12019. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/v162/](https://proceedings.mlr.press/v162/langosco22a.html)
180 [langosco22a.html](https://proceedings.mlr.press/v162/langosco22a.html).
- 181 Kimin Lee, Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Sunrise: A simple unified
182 framework for ensemble learning in deep reinforcement learning. In Marina Meila and Tong
183 Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume
184 139 of *Proceedings of Machine Learning Research*, pages 6131–6141. PMLR, 18–24 Jul 2021.
185 URL <https://proceedings.mlr.press/v139/lee21g.html>.
- 186 Yoonho Lee, Huaxiu Yao, and Chelsea Finn. Diversify and disambiguate: Learning from underspeci-
187 fied data. 2022.
- 188 Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep explo-
189 ration via bootstrapped dqn. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and
190 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Cur-
191 ran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper/2016/file/](https://proceedings.neurips.cc/paper/2016/file/8d8818c8e140c64c743113f563cf750f-Paper.pdf)
192 [8d8818c8e140c64c743113f563cf750f-Paper.pdf](https://proceedings.neurips.cc/paper/2016/file/8d8818c8e140c64c743113f563cf750f-Paper.pdf).
- 193 Jongjin Park, Younggyo Seo, Chang Liu, Li Zhao, Tao Qin, Jinwoo Shin, and Tie-Yan
194 Liu. Object-aware regularization for addressing causal confusion in imitation learning. In
195 M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors,
196 *Advances in Neural Information Processing Systems*, volume 34, pages 3029–3042. Cur-
197 ran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper/2021/file/](https://proceedings.neurips.cc/paper/2021/file/17a3120e4e5fbdc3cb5b5f946809b06a-Paper.pdf)
198 [17a3120e4e5fbdc3cb5b5f946809b06a-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/17a3120e4e5fbdc3cb5b5f946809b06a-Paper.pdf).

- 199 Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized Experience Replay. *arXiv*
 200 *e-prints*, art. arXiv:1511.05952, November 2015.
- 201 Jeremy Tien, Jerry Zhi-Yang He, Zackory Erickson, Anca Dragan, and Daniel S. Brown. A study
 202 of causal confusion in preference-based reward learning. In *ICML 2022: Workshop on Spurious*
 203 *Correlations, Invariance and Stability*, 2022. URL <https://openreview.net/forum?id=WaZZ0Sw9fwf>.
- 205 Chuan Wen, Jierui Lin, Trevor Darrell, Dinesh Jayaraman, and Yang Gao. Fighting copycat agents in
 206 behavioral cloning from observation histories. In *NeurIPS*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1b113258af3968aaf3969ca67e744ff8-Abstract.html>.
- 208 Chuan Wen, Jierui Lin, Jianing Qian, Yang Gao, and Dinesh Jayaraman. Keyframe-focused visual
 209 imitation learning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International*
 210 *Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,
 211 pages 11123–11133. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/wen21d.html>.

213 A Appendix

214 A.1 Implementation

215 All our environments use a discrete action space. Therefore we build our method on top of the
 216 double-DQN implementation similar to the original CQL paper. As stated in section 1, we use
 217 ensembles of Q networks and at evaluation time, we average the Q-value outputs of the ensemble,
 218 and select the action with the maximum Q-value. In other place where we need to do inference
 219 (for instance: to compute Q -values for the conservative loss) we similarly take the mean across the
 220 ensemble.

221 A.2 Code and Data

222 We will release our code, data and pretrained models once the work is uploaded online. The code
 223 repository will also contain code to reproduce all the figures in this work.

224 A.3 Data Collection

- 225 1. Traffic-World: To collect data for Offline RL, we trained a PPO agent on a slightly modified
 226 version of the Traffic-world environment, with reward shaping on the environment, to
 227 incentivise the agent to reach the goal since this could is a hard exploration environment
 228 (there is the potential to receive many negative rewards before receiving a positive reward,
 229 and without reward shaping the PPO agent just learns to toggle in-place till the episode ends
 230 to avoid negative penalties).
- 231 2. Maze: We use the Impala-based PPO agent trained in [Langosco et al., 2022] for 200M
 232 steps to collect the expert trajectories on 6000 episodes of epsides with randomised goals
 233 and 200 episodes of epsides with fixed goals.

234 A.4 Hyper-parameters

235 **CQL:** We conduct a grid search over the learning rate and conservative penalty coefficient (α). We
 236 use gradient clipping with the norm varied between 3,5,7.

237 **Active Sampling:** We kept all the hyper-parameters the same as random sampling (batch size,
 238 learning rate, α). The parameters related to active sampling are

- 239 1. n : the number of gradient steps with stale scores we take before we recompute acquisition
 240 scores on the data.
- 241 2. the ensemble size which we set to 3, and keep constant across the active and random-
 242 sampling variants for a fair comparison.

243 **A.5 Computational Cost**

244 Figure 3 shows a scatterplot for the wallclock times to achieve highest reward across different active
245 and random baselines. It also plots the time needed for active sampling variants to achieve the best
246 reward that random sampling achieves (denoted as Variance-par-Random and TD-par-Random in the
247



Figure 3: Timing Comparison for different sampling schemes on the Procgen-Maze benchmark plotted as reward achieved versus wallclock time in minutes.