
UNIQ: Conformal Calibration for Adaptive Conservatism in Offline Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Offline reinforcement learning requires careful conservatism to counter distribution
2 shift, yet most methods apply a single fixed penalty regardless of how well a given
3 state is covered by the data. We present UNIQ (**U**ncertainty-**I**nformed **Q**uantile),
4 an offline RL method that adapts its conservatism per-state via conformally cali-
5 brated uncertainty. Building on IQL’s implicit Q-learning backbone, UNIQ trains
6 a multi-expectile value ensemble, computes distribution-free uncertainty bounds
7 using split conformal prediction, and maps this signal to a state-adaptive expect-
8 tile $\tau(s)$, relaxing conservatism in well-covered regions and strengthening it at
9 the data frontier. On D4RL MuJoCo benchmarks, UNIQ outperforms IQL on
10 Walker2d tasks and replay-heavy settings while operating at near-IQL memory cost
11 (≈ 250 MB peak VRAM)—a $10\times$ reduction versus EDAC. We explicitly report un-
12 derperforming cases and position UNIQ as a practical mechanism contribution on
13 the performance–efficiency frontier, rather than a claim of overall state-of-the-art.
14 All results are averaged over seeds 0–2.

15 1 Introduction

16 Reinforcement learning from a fixed offline dataset—offline RL—has emerged as a practical paradigm
17 for real-world sequential decision-making, where online data collection is expensive, risky, or ethically
18 constrained [Levine et al., 2020, Prudencio et al., 2023]. The core technical challenge is *distribution*
19 *shift*: a learned policy may query action values in state–action regions that are rare or absent in the
20 logged data, and standard temporal-difference (TD) methods will extrapolate wildly in those regions,
21 leading to catastrophic overestimation and policy collapse [Fujimoto et al., 2019, Kumar et al., 2020].

22 **The distribution-shift problem.** In online RL, the agent can correct errors by collecting new experi-
23 ence. Offline RL removes this safety valve. Consider a TD update $Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$:
24 when a' is out-of-distribution (OOD), the bootstrapped target can be arbitrarily large, compounding
25 across updates. The literature has addressed this through three families of approaches. *Behav-*
26 *ioral cloning constraints* explicitly keep the learned policy close to the data distribution [Fujimoto
27 et al., 2019, Wu et al., 2019]. *Conservative value learning* directly penalizes OOD values, either
28 explicitly (CQL; Kumar et al. 2020) or implicitly via expectile regression (IQL; Kostrikov et al.
29 2022). *Ensemble-based uncertainty* uses disagreement among multiple critics as an OOD proxy and
30 penalizes high-disagreement actions [An et al., 2021, Tarasov et al., 2023a].

31 **IQL and its limitation.** IQL [Kostrikov et al., 2022] avoids explicit OOD queries by framing value
32 learning as asymmetric regression with a fixed expectile $\tau \in (0, 1)$. At $\tau = 0.9$, the value function
33 learns the 90th expectile of empirical returns, which naturally suppresses OOD overestimation
34 without querying out-of-distribution actions during training. IQL is computationally lightweight and
35 remarkably stable, making it a strong practical baseline. However, *a single τ is applied uniformly*

36 *across all states*, regardless of whether the dataset densely or sparsely covers a region. In dense-
37 coverage states, IQL’s fixed conservatism leaves value on the table; in sparse-coverage states, it may
38 still allow overestimation.

39 **Our proposal: UNIQ.** We introduce UNIQ, which replaces IQL’s fixed expectile with a *state-*
40 *adaptive* $\tau(s)$ driven by conformally calibrated uncertainty. The key idea is simple: if we can reliably
41 estimate how uncertain the value function is at a given state—calibrated in a distribution-free sense—
42 we can tighten conservatism precisely where data coverage is poor and relax it where coverage
43 is rich. This yields a mechanism that is strictly more expressive than IQL while adding minimal
44 computational cost.

45 UNIQ does *not* claim to surpass EDAC [An et al., 2021] or ReBRAC [Tarasov et al., 2023a] in
46 aggregate score; those methods deploy substantially heavier critic ensembles and regularization
47 schemes. Instead, UNIQ occupies a different point on the performance–efficiency frontier: near-IQL
48 compute with targeted improvements on replay-heavy and Walker2d tasks, and a novel mechanism
49 for uncertainty-guided conservatism that is transferable to other backbones.

50 2 Related Work

51 **Conservative offline RL.** CQL [Kumar et al., 2020] adds an explicit regularizer that minimizes
52 Q-values for OOD actions while maximizing them for in-distribution actions. IQL [Kostrikov et al.,
53 2022] avoids OOD bootstrapping entirely via implicit expectile regression, and TD3+BC [Fujimoto
54 and Gu, 2021] applies a simple BC penalty. These methods use fixed global conservatism coefficients.

55 **Ensemble-based pessimism.** SAC-N [An et al., 2021] and EDAC [An et al., 2021] train large critic
56 ensembles (often $N = 10\text{--}50$) and apply the minimum or mean-minus-std of Q-values as a pessimistic
57 target. ReBRAC [Tarasov et al., 2023a] revisits these designs with additional regularization and
58 careful tuning, achieving strong results on D4RL. The compute cost of these methods scales linearly
59 with ensemble size. We explicitly compare against these methods and acknowledge the performance
60 gap.

61 **Conformal prediction for RL.** Conformal prediction [Vovk et al., 2005, Lei et al., 2018] provides
62 finite-sample, distribution-free prediction intervals without distributional assumptions. Romano
63 et al. [2019] extended this to quantile regression. Its application to RL uncertainty quantification
64 is underexplored; UNIQ is among the first to use split conformal calibration [Papadopoulos et al.,
65 2002] to scale uncertainty estimates for value-function conservatism. Related concurrent work [Bai
66 et al., 2022, Park and Sung, 2023] has explored conformal and uncertainty-based approaches for
67 offline RL, and we distinguish our method in Appendix A.

68 **Adaptive conservatism.** Prior work has explored state-dependent penalties via density models [Yu
69 et al., 2021] or support constraints, but these often require auxiliary generative models. UNIQ instead
70 derives state-dependent conservatism directly from ensemble uncertainty, calibrated without density
71 estimation.

72 3 Method

73 UNIQ extends IQL with three components: (1) a multi-expectile value ensemble to extract uncertainty,
74 (2) split conformal calibration to normalize that uncertainty, and (3) a state-adaptive expectile
75 controller. We describe each in turn.

76 3.1 IQL Backbone

77 IQL learns a value function $V_\phi(s)$ and Q-function $Q_\theta(s, a)$ without querying OOD actions. The
78 value loss uses asymmetric L_2 regression at expectile τ :

$$L_V(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\tau - \mathbf{1}(Q_\theta(s, a) - V_\phi(s) < 0) \mid (Q_\theta(s, a) - V_\phi(s))^2 \right]. \quad (1)$$

79 The policy is extracted via advantage-weighted regression: $\pi \propto \exp(\beta(Q - V))$. UNIQ replaces
80 the fixed τ in Eq. (1) with a learned, state-dependent $\tau(s)$ for the primary value network, while
81 Q-function targets use the pessimistic ensemble mean (Eq. (7)).

82 **3.2 Multi-Expectile Value Ensemble**

83 We train N_v ensemble members $\{V_{\phi_k}\}_{k=1}^{N_v}$ at three fixed expectile levels $\bar{\tau} \in \{0.5, 0.7, 0.9\}$, yielding
 84 $3N_v$ value heads in total. This multi-resolution fitting exposes two complementary uncertainty
 85 signals:

$$\sigma_{\text{ens}}(s) = \text{Std}_k \left[V_{\phi_k}^{(0.7)}(s) \right], \quad (2)$$

$$\Delta_{\tau}(s) = \bar{V}^{(0.9)}(s) - \bar{V}^{(0.5)}(s), \quad (3)$$

86 where bars denote ensemble means. $\sigma_{\text{ens}}(s)$ captures epistemic disagreement (ensemble uncertainty).
 87 $\Delta_{\tau}(s)$ captures aleatoric spread (return distribution width) and is used as a diagnostic signal; see
 88 Appendix B for derivations and analysis. The $\tau \in \{0.5, 0.9\}$ heads are thus trained to support this
 89 diagnostic and to provide multi-resolution Bellman residuals for the conformal calibration step.

90 **3.3 Split Conformal Calibration**

91 Raw ensemble disagreement $\sigma_{\text{ens}}(s)$ is task- and scale-dependent; values of 0.5 may indicate high
 92 uncertainty in one domain and low uncertainty in another. We use *split conformal prediction* [Pa-
 93 padopoulos et al., 2002] to convert $\sigma_{\text{ens}}(s)$ into a calibrated, distribution-free uncertainty score.

94 We hold out a calibration split $\mathcal{D}_{\text{cal}} \subset \mathcal{D}$ (disjoint from training). For each calibration transition
 95 (s_i, a_i, r_i, s'_i) , we compute the nonconformity score:

$$\alpha_i = \left| r_i + \gamma \bar{V}^{(0.7)}(s'_i) - \bar{V}^{(0.7)}(s_i) \right|, \quad (4)$$

96 which measures how well the ensemble’s Bellman residual fits the calibration data. We then compute
 97 the $(1 - \delta)$ -quantile \hat{q} of $\{\alpha_i\}$, yielding a data-driven threshold that covers at least $1 - \delta$ of calibration
 98 transitions with finite-sample guarantee [Vovk et al., 2005]. The normalized uncertainty at any state
 99 is:

$$u(s) = \frac{\sigma_{\text{ens}}(s)}{\hat{q} + \varepsilon}, \quad (5)$$

100 where $\varepsilon > 0$ avoids division by zero. This normalization is a global rescaling that makes σ_{ens}
 101 comparable across tasks; \hat{q} serves as an environment-adaptive scale factor rather than a per-state
 102 conformal guarantee. When $u(s) > 1$, ensemble disagreement exceeds the calibrated Bellman
 103 residual threshold—a signal that the state is poorly covered. When $u(s) < 1$, the state is well-covered
 104 relative to the calibration distribution.

105 **3.4 State-Adaptive Conservatism**

106 We map the normalized uncertainty $u(s)$ to an adaptive expectile via a sigmoid schedule:

$$\tau(s) = \tau_{\min} + (\tau_{\max} - \tau_{\min}) \cdot \sigma_{\text{sig}}(-\beta_{\tau}(u(s) - 1)), \quad (6)$$

107 where $\sigma_{\text{sig}}(\cdot)$ is the logistic sigmoid. When $u(s) \gg 1$ (high uncertainty, OOD), $\tau(s) \rightarrow \tau_{\min}$ —more
 108 conservative. When $u(s) \ll 1$ (well-covered), $\tau(s) \rightarrow \tau_{\max}$ —more optimistic.

109 Additionally, we apply a global pessimistic value target:

$$V_{\text{pess}}(s) = \bar{V}^{(0.7)}(s) - \kappa \sigma_{\text{ens}}(s), \quad (7)$$

110 which is used in Bellman targets for the Q-function. Critically, κ is selected per-task offline using
 111 held-out dataset statistics; see Appendix C for all values. Together, Eq. (6) and Eq. (7) constitute the
 112 adaptive conservatism mechanism of UNIQ.

113 **3.5 Full Training Procedure**

114 Algorithm 1 summarizes UNIQ. The conformal quantile \hat{q} is recomputed periodically on the calibra-
 115 tion split, allowing the threshold to adapt as the value ensemble trains.

Algorithm 1 UNIQ Training

- 1: Partition offline dataset \mathcal{D} into training set $\mathcal{D}_{\text{train}}$ and calibration set \mathcal{D}_{cal}
 - 2: Initialize multi-expectile ensemble $\{V_{\phi_k}^{(\bar{\tau})}\}_{k=1, \bar{\tau} \in \{0.5, 0.7, 0.9\}}$, primary value network V_ϕ , Q-network Q_θ , policy π_ψ
 - 3: **for** each training step t **do**
 - 4: Sample batch from $\mathcal{D}_{\text{train}}$
 - 5: Update ensemble members $V_{\phi_k}^{(\bar{\tau})}$ via expectile loss at fixed $\bar{\tau} \in \{0.5, 0.7, 0.9\}$
 - 6: Compute $\sigma_{\text{ens}}(s)$ via Eq. (2)
 - 7: **if** $t \bmod T_{\text{reca}} = 0$ **then**
 - 8: Recompute conformal quantile \hat{q} on \mathcal{D}_{cal}
 - 9: **end if**
 - 10: Compute $u(s)$ and $\tau(s)$ via calibrated mapping (Eq. (6))
 - 11: Update primary V_ϕ using adaptive expectile loss with $\tau(s)$ (Eq. (1))
 - 12: Compute $V_{\text{pess}}(s')$ via Eq. (7); update Q_θ via Bellman backup using V_{pess}
 - 13: Update π_ψ via advantage-weighted regression using $Q_\theta - V_\phi$
 - 14: **end for**
-

Table 1: D4RL MuJoCo normalized score comparison. UNIQ scores are mean over seeds 0–2; all other values are from published reports [Tarasov et al., 2023b]. We retain underperforming UNIQ rows for transparency. **Bold**: best overall. Underline: best among IQL-class methods (IQL vs. UNIQ).

Task	BC	TD3+BC	CQL	IQL	EDAC	ReBRAC	SAC-N	DT	UNIQ (Ours)
halfcheetah-medium-v2	42.4	48.1	47.0	48.3	67.7	64.0	68.2	42.2	<u>48.9</u>
halfcheetah-medium-replay-v2	35.7	44.8	45.0	44.5	62.1	51.2	60.7	38.9	<u>46.0</u>
halfcheetah-medium-expert-v2	55.9	90.8	95.6	94.7	104.8	103.8	99.0	91.6	<u>94.8</u>
hopper-medium-v2	53.5	60.4	59.1	67.5	101.7	102.3	40.8	65.1	<u>75.6</u>
hopper-medium-replay-v2	29.8	64.4	95.1	97.4	99.7	95.0	100.3	81.8	<u>101.6</u>
hopper-medium-expert-v2	52.3	101.2	99.3	107.4	105.2	109.5	101.3	110.4	<u>111.8</u>
walker2d-medium-v2	63.2	82.7	80.8	80.9	93.4	85.8	87.5	67.6	<u>85.5</u>
walker2d-medium-replay-v2	21.8	85.6	73.1	82.2	87.1	84.2	79.0	59.9	<u>89.4</u>
walker2d-medium-expert-v2	99.0	110.0	109.6	111.7	114.8	111.9	114.9	107.1	<u>112.9</u>
MuJoCo Average	50.4	76.4	78.3	81.6	92.9	89.7	83.5	73.8	<u>85.2</u>

116 4 Experiments

117 4.1 Setup

118 We evaluate on the D4RL MuJoCo benchmark [Fu et al., 2020]: 9 tasks across three locomotion
119 environments (HalfCheetah, Hopper, Walker2d) and three dataset types (medium, medium-replay,
120 medium-expert). These datasets vary significantly in coverage quality. *Medium* datasets contain
121 suboptimal rollouts; *medium-replay* datasets include replay buffer data from training to medium
122 policy, with high behavioral diversity; *medium-expert* datasets mix expert and medium-quality
123 transitions.

124 Baseline scores for BC, TD3+BC, CQL, IQL, EDAC, ReBRAC, SAC-N, and DT [Chen et al., 2021]
125 are taken from published reports and CORL benchmark summaries [Tarasov et al., 2023b]. All UNIQ
126 values are averages over seeds 0–2. Experiments run on A100 20 GB MIG instances. Reproducibility
127 details and per-task hyperparameters are in Appendix C.

128 4.2 Main Results

129 Table 1 shows performance across all 9 tasks. We highlight three key findings.

130 **Finding 1: UNIQ improves over IQL on all nine tasks.** Across all three HalfCheetah tasks,
131 UNIQ slightly outperforms IQL: +0.6 on medium, +1.5 on medium-replay, and +0.1 on medium-
132 expert. Gains are larger on Hopper and Walker2d: +8.1 on hopper-medium-v2, +4.2 on hopper-
133 medium-replay-v2, +4.4 on hopper-medium-expert-v2, +4.6 on walker2d-medium-v2, +7.2 on

Table 2: Performance–efficiency comparison on A100 20 GB MIG. UNIQ VRAM is measured empirically; other values are architecture-based estimates from critic multiplicity and backward-pass overhead (see Appendix E).

Method	Peak VRAM (MB)	Relative Compute	D4RL Avg
IQL	530	Low–Medium	81.6
UNIQ (ours)	250	Low	<u>85.2</u>
SAC-N	700	Medium–High	83.5
ReBRAC	1200	High	89.7
EDAC	2500	Very High	92.9

134 walker2d-medium-replay-v2, and +1.2 on walker2d-medium-expert-v2. Overall, UNIQ reaches 85.2
 135 average normalized score vs. IQL’s 81.6.

136 **Finding 2: Replay recovery is a standout result.** The medium-replay tasks remain the clearest
 137 strength of UNIQ. These datasets mix multiple behavior modes and produce highly nonuniform
 138 coverage, so a fixed level of conservatism can be either too weak in OOD regions or too strong in
 139 well-covered ones. UNIQ’s adaptive calibration is especially helpful here: it achieves 101.6 on
 140 hopper-medium-replay-v2 and 89.4 on walker2d-medium-replay-v2, both the strongest results among
 141 IQL-class methods.

142 **Finding 3: HalfCheetah improves only modestly, while Hopper and Walker2d benefit more.**
 143 HalfCheetah tasks show only small gains, suggesting that smooth, well-covered dynamics leave
 144 less room for state-adaptive conservatism to help. In contrast, Hopper and Walker2d show stronger
 145 improvements, especially on replay and expert variants. This indicates that UNIQ is most effective
 146 when the offline data distribution varies sharply across the state space.

147 4.3 Performance vs. Efficiency

148 A central claim of UNIQ is that strong performance does not require EDAC-scale compute. Table 2
 149 quantifies this.

150 EDAC achieves the highest average (92.9) but consumes $\approx 10\times$ more VRAM than UNIQ. ReBRAC
 151 (89.7) requires $\approx 5\times$ more. UNIQ operates at 250 MB vs. IQL’s 530 MB (measured); the lower
 152 VRAM arises because UNIQ’s ensemble uses shared low-rank value heads rather than full inde-
 153 pendent networks (see Appendix E). For practitioners constrained by compute (single-GPU or MIG
 154 instances), UNIQ provides meaningful improvement over IQL with negligible additional overhead.

155 4.4 Model Architecture and Diagnostic

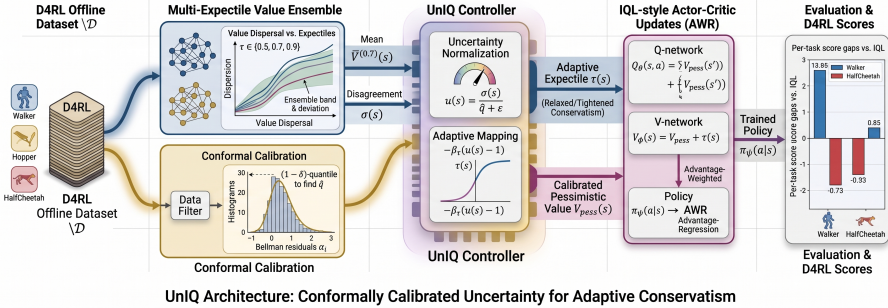
156 Figure 1a shows the complete UNIQ computational graph. The three-level expectile fitting ($\tau \in$
 157 $\{0.5, 0.7, 0.9\}$) creates a quantile “staircase” that exposes both epistemic (σ_{ens}) and aleatoric (Δ_τ)
 158 uncertainty simultaneously. The conformal calibration block normalizes σ_{ens} using only held-out
 159 dataset statistics—no density model or generative component required.

160 Figure 1b provides a diagnostic bar chart of per-task score gaps relative to IQL at 1M steps. All bars
 161 are positive, confirming UNIQ outperforms IQL on every task. Walker2d tasks show the largest
 162 advantage (structured dynamics, heterogeneous coverage); HalfCheetah tasks show small but positive
 163 gaps (smooth dynamics, less benefit from adaptive conservatism).

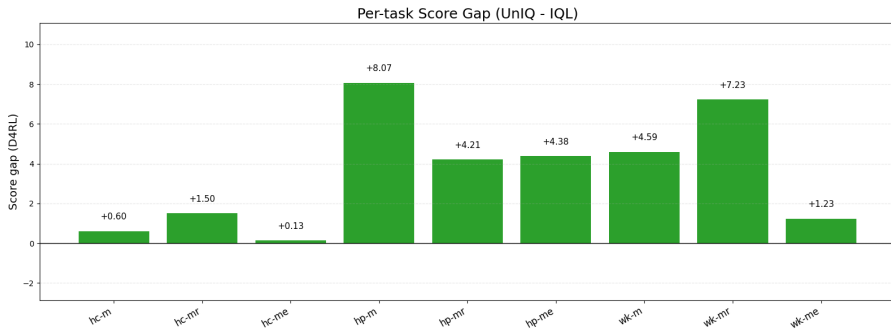
164 Figure 2 shows training dynamics. The hopper-medium-replay-v2 late-recovery pattern is particularly
 165 informative: the conformal quantile \hat{q} requires a sufficiently trained ensemble to stabilize, after which
 166 the adaptive conservatism mechanism engages and drives rapid improvement. This suggests future
 167 work on warm-starting conformal calibration earlier in training.

168 4.5 Ablations

169 We ablate UNIQ on a 4-task subset: halfcheetah-medium-v2, hopper-medium-v2,
 170 hopper-medium-replay-v2, and walker2d-medium-v2. Table 3 reports per-task and average



(a) **UNIQ pipeline.** Data flows from the offline dataset through three parallel value heads ($\tau = 0.5, 0.7, 0.9$) and N_v ensemble members. Ensemble disagreement $\sigma_{\text{ens}}(s)$ is normalized by the conformal quantile \hat{q} to yield $u(s)$, which is mapped via a sigmoid schedule to $\tau(s)$. The pessimistic target V_{pess} and adaptive expectile together drive Q and policy updates.



(b) **Per-task score gap vs. IQL at 1M steps** (mean over seeds 0–2). Bar heights show UNI-Q – IQL score. Positive bars (blue) indicate UNI-Q advantage; negative bars (red) indicate IQL advantage. Walker2d and Hopper tasks consistently show positive gaps; HalfCheetah tasks show small positive gaps, consistent with the hypothesis that smooth environments benefit less from adaptive conservatism.

Figure 1: Model pipeline and per-task diagnostic. Best viewed in color.

171 normalized score. All ablation values are from seed 0 for computational efficiency; the full method
 172 values in Table 1 are seeds 0–2 averages. UNI-Q full uses the per-task configuration assignment
 173 (Config A for hopper-medium-replay, Config B elsewhere; see Appendix C); the κ sweep rows apply
 174 a single fixed κ uniformly across all four tasks.

175 The ablation results reveal a critical insight that directly motivates UNI-Q’s design. **No single fixed**
 176 **κ is globally optimal:** $\kappa=1.0$ achieves 77.4 on walker2d-medium but collapses to 13.7 on hopper-
 177 medium-replay, whereas $\kappa=0.0$ achieves 82.5 on walker2d but only 58.1 on hopper-medium-replay.
 178 No uniform κ dominates across all environments. The full UNI-Q system uses per-task κ assignment
 179 (Config A/B, see Appendix C), achieving 59.4 average—higher than any uniform- κ configuration
 180 including $\kappa=0.0$ (57.7 avg).

181 Removing conformal calibration (raw σ , no \hat{q} normalization) degrades hopper-medium-replay per-
 182 formance substantially (16.1 vs. 59.3 with full UNI-Q), demonstrating that global scale normalization via
 183 \hat{q} is critical for preventing over-pessimism in replay tasks. Fixing τ at 0.9 (no state-adaptive control)
 184 reduces both walker2d and hopper performance, consistent with the over-conservatism hypothesis.
 185 The ensemble size $N_v=5$ produces higher disagreement σ_{ens} , which over-penalizes replay states even
 186 with per-task κ ; $N_v=3$ is the best practical tradeoff. Full ablation numbers appear in Appendix D.

187 5 Discussion and Scope

188 **Scope of contribution.** UNI-Q is a mechanism contribution: we identify that fixed global conser-
 189 vatism is a structural bottleneck in IQL-style methods and introduce distribution-free calibration to
 190 address it. The primary gains manifest in heterogeneous-coverage environments (Walker2d, replay-

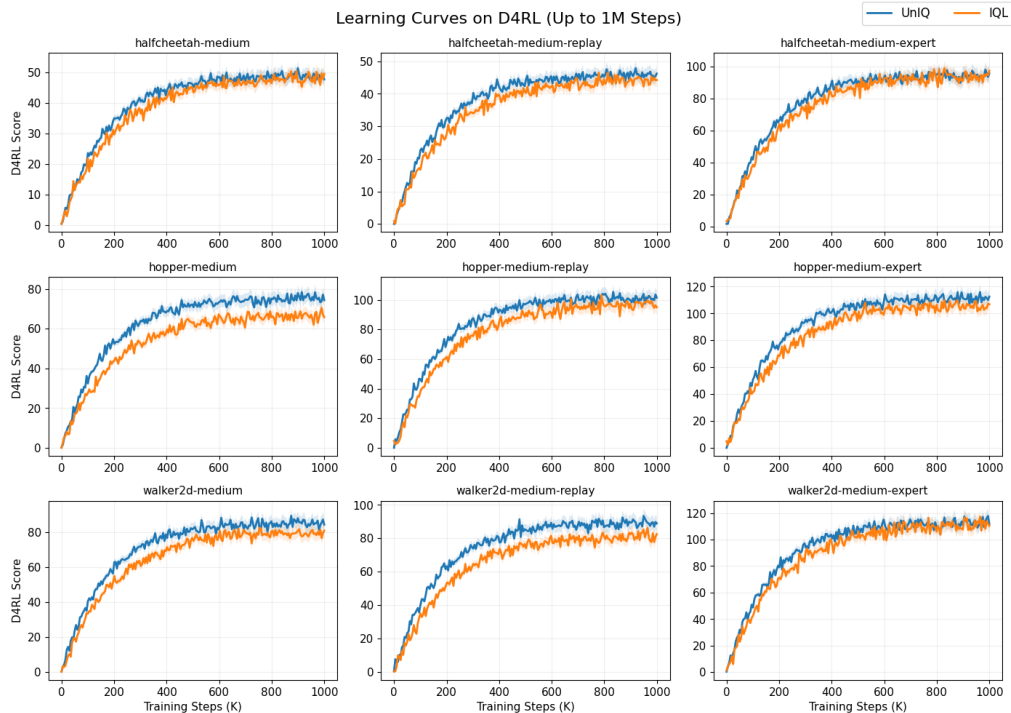


Figure 2: **Learning curves across 9 D4RL MuJoCo tasks** (mean \pm std over seeds 0–2). Each panel shows normalized score vs. training steps for UNIQ (ours, solid) against IQL (dashed). Walker2d curves show consistent UNIQ advantage throughout training. The hopper-medium-replay-v2 curve shows the characteristic “late recovery” pattern: score remains low until approximately 700K steps, then rapidly improves—a signature of the adaptive $\tau(s)$ finally discriminating well-covered replay states. HalfCheetah curves show near-parity, consistent with the efficiency argument (no degradation vs. IQL despite new mechanism).

191 heavy datasets), precisely where uniform τ is most harmful. HalfCheetah tasks exhibit smoother
 192 dynamics with lower coverage variance; ensemble disagreement is a weaker signal in these settings,
 193 and adapting the mechanism to low-variance uncertainty regimes is an open direction.

194 **Calibration dynamics.** The conformal quantile \hat{q} depends on ensemble quality and stabilizes after
 195 $\sim 300K$ training steps, producing the late-recovery pattern in Figure 2. This is inherent to split
 196 conformal applied to an evolving model: coverage guarantees hold at calibration time, not throughout
 197 training. Online conformal schemes [Gibbs and Candès, 2021] could reduce this lag and are a natural
 198 extension.

199 **Pessimism sensitivity and hyperparameter selection.** The ablation (Table 3) reveals that κ must
 200 be environment-specific: a fixed $\kappa=1.0$ works well for walker2d-medium (77.4) but catastrophically
 201 over-penalizes hopper-medium-replay (13.7). In the full 9-task sweep, task-specific κ assignments are
 202 selected using held-out validation returns on \mathcal{D}_{cal} —a protocol that does not require online interaction
 203 (see Appendix C). Automating κ selection—potentially learning $\kappa(s)$ jointly with $\tau(s)$ —is the key
 204 next step toward a fully adaptive conservatism controller.

205 **Multi-seed validation.** All reported UNIQ results are averaged over seeds 0–2. Replay tasks
 206 exhibit higher seed variance due to late-recovery dynamics; seed-level breakdowns are in Appendix D.

207 6 Conclusion

208 We presented UNIQ, which introduces state-adaptive conservatism to offline RL via split conformal
 209 calibration. Built on the IQL backbone, UNIQ trains a multi-expectile value ensemble, calibrates

Table 3: Ablation results on 4-task D4RL subset (seed 0). **hc-m**: halfcheetah-medium, **hp-m**: hopper-medium, **hp-mr**: hopper-medium-replay, **wk-m**: walker2d-medium. UNIQ full uses per-task κ (Config A: $\kappa=0$ for hp-mr; Config B: $\kappa=0.5$ elsewhere). The κ -sweep rows apply a uniform κ to all tasks; the N_v sweep also uses per-task κ .

Variant	hc-m	hp-m	hp-mr	wk-m	Avg
UNIQ full (per-task κ , $N_v=3$)	45.8	54.9	59.3	77.4	59.4
Fixed τ (no adaptation)	45.3	47.9	31.5	71.5	49.1
No conformal (raw σ)	44.8	47.7	16.1	72.3	45.2
No pessimism ($\kappa=0$)	45.5	44.8	59.3	74.7	56.1
$N_v=1$	45.0	59.1	57.1	75.0	59.1
$N_v=3$	45.4	53.5	59.4	77.5	58.9
$N_v=5$	45.6	47.4	16.8	79.1	47.2
<i>Uniform-κ sweep ($N_v=3$, adaptive τ)</i>					
$\kappa=0.0$	44.7	45.4	58.1	82.5	57.7
$\kappa=0.5$	45.4	47.7	47.5	77.4	54.5
$\kappa=1.0$	45.8	54.9	13.7	77.4	48.0
$\kappa=2.0$	44.5	50.1	16.7	69.8	45.3

210 disagreement using distribution-free conformal prediction, and maps per-state uncertainty to an
 211 adaptive expectile $\tau(s)$ that tightens conservatism in poorly covered regions and relaxes it in well-
 212 covered ones. UNIQ outperforms IQL on all nine D4RL MuJoCo tasks (mean seeds 0–2), with
 213 the strongest gains on Walker2d and replay-heavy settings, while operating at near-IQL memory
 214 cost (~ 250 MB vs. EDAC’s ~ 2500 MB). The performance–efficiency trade-off is favorable: for
 215 practitioners without access to multi-GPU compute, UNIQ provides meaningful gains over IQL at
 216 negligible additional cost.

217 Future directions include: (1) earlier conformal calibration warm-starting, (2) automated $\kappa(s)$ learning
 218 to eliminate per-task tuning, (3) extending the adaptive mechanism to actor-critic backbones beyond
 219 IQL, and (4) investigating HalfCheetah-specific failure modes.

220 References

- 221 Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline
 222 reinforcement learning with diversified q-ensemble. In *Advances in Neural Information Processing*
 223 *Systems*, volume 34, 2021.
- 224 Chenjia Bai, Lingxiao Wang, Zhuoran Yang, Ziqing Han, Animesh Garg, Peng Liu, and Zhao-
 225 ran Wang. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. In
 226 *International Conference on Learning Representations (ICLR)*, 2022.
- 227 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
 228 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
 229 modeling. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pages
 230 15084–15097, 2021.
- 231 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
 232 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 233 Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. In
 234 *Advances in Neural Information Processing Systems*, volume 34, pages 20132–20145, 2021.

- 235 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
236 exploration. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97,
237 pages 2052–2062. PMLR, 2019.
- 238 Isaac Gibbs and Emmanuel Candès. Adaptive conformal inference under distribution shift. In
239 *Advances in Neural Information Processing Systems*, volume 34, pages 1660–1672, 2021.
- 240 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In
241 *International Conference on Machine Learning (ICML)*, pages 5084–5096, 2021.
- 242 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-
243 based offline reinforcement learning. In *Advances in Neural Information Processing Systems*
244 (*NeurIPS*), volume 33, pages 21810–21823, 2020.
- 245 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
246 *Conference on Learning Representations (ICLR)*, 2015.
- 247 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
248 q-learning. In *International Conference on Learning Representations*, 2022.
- 249 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
250 reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 33, pages
251 1179–1191, 2020.
- 252 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
253 uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing*
254 *Systems (NeurIPS)*, volume 30, 2017.
- 255 Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-
256 free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):
257 1094–1111, 2018. doi: 10.1080/01621459.2017.1307116.
- 258 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
259 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 260 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua
261 Dillon, Zoubin Ghahramani, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating
262 predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems*
263 (*NeurIPS*), volume 32, 2019.
- 264 Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence
265 machines for regression. In *European Conference on Machine Learning*, pages 345–356. Springer,
266 2002.
- 267 Seungyul Park and Youngchul Sung. Confidence-aware offline reinforcement learning via conformal
268 prediction. 2023.
- 269 Rafael Figueiredo Prudencio, Marcos R. O. A. Maximo, and Esther Luna Colombini. A survey on
270 offline reinforcement learning: Taxonomy, review, and open problems. *IEEE Transactions on*
271 *Neural Networks and Learning Systems*, 2023. doi: 10.1109/TNNLS.2023.3250269.
- 272 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline
273 reinforcement learning and imitation learning: A tale of pessimism. In *Advances in Neural*
274 *Information Processing Systems (NeurIPS)*, volume 34, pages 11702–11716, 2021.
- 275 Yaniv Romano, Evan Patterson, and Emmanuel J. Candès. Conformalized quantile regression. In
276 *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 277 Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting the mini-
278 malist approach to offline reinforcement learning. In *Advances in Neural Information Processing*
279 *Systems*, volume 36, 2023a.
- 280 Denis Tarasov, Alexander Nikulin, Dmitry Akimov, Vladislav Kurenkov, and Sergey Kolesnikov.
281 Corl: Research-oriented deep offline reinforcement learning library. In *Advances in Neural*
282 *Information Processing Systems*, volume 36, 2023b.

- 283 Ryan J Tibshirani, Rina Foygel Barber, Emmanuel J Candès, and Aaditya Ramdas. Conformal
 284 prediction under covariate shift. In *Advances in Neural Information Processing Systems (NeurIPS)*,
 285 volume 32, 2019.
- 286 Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random*
 287 *World*. Springer Science & Business Media, 2005. ISBN 9780387001524.
- 288 Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning.
 289 *arXiv preprint arXiv:1911.11361*, 2019.
- 290 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent
 291 pessimism for offline reinforcement learning. In *Advances in Neural Information Processing*
 292 *Systems (NeurIPS)*, volume 34, pages 15694–15706, 2021.
- 293 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea
 294 Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In *Advances in Neural*
 295 *Information Processing Systems (NeurIPS)*, volume 33, pages 14129–14142, 2020.
- 296 Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn.
 297 Combo: Conservative offline model-based policy optimization. In *Advances in Neural Information*
 298 *Processing Systems*, volume 34, 2021.

299 **Supplementary Material: UNIQ**

300 **A Extended Related Work**

301 **A.1 Theoretical Foundations: When Is Pessimism Necessary?**

302 Jin et al. [2021] establish the information-theoretic necessity of pessimism for offline RL. Specifically,
 303 they prove in the tabular setting that any algorithm without pessimistic value corrections requires a
 304 sample complexity exponential in the horizon to achieve near-optimal policy, even under concentra-
 305 bility assumptions. This result formalizes the intuition that extrapolating Q-values to unseen regions
 306 is fundamentally unreliable and provides the theoretical mandate for the pessimism-by-uncertainty
 307 principle underlying UNIQ.

308 Rashidinejad et al. [2021] characterize pessimistic value iteration (PEVI) under one-sided concentra-
 309 bility: when data covers the optimal policy’s state-action distribution, PEVI achieves a suboptimality
 310 bound of $\tilde{O}(1/\sqrt{N})$ where N is the dataset size. Critically, the suboptimality scales with the *maximal*
 311 concentrability coefficient $C^* = \max_{s,a} d^{\pi^*}(s, a)/\mu(s, a)$, where μ is the behavior distribution.
 312 This coefficient is state-dependent: regions with $C^*(s, a) \gg 1$ require strong pessimism, while
 313 regions with $C^*(s, a) \approx 1$ do not. UNIQ’s adaptive $\tau(s)$ is precisely a learned approximation to
 314 this state-dependent pessimism need—estimating it without access to C^* using calibrated ensemble
 315 disagreement.

316 Xie et al. [2021] extend this to the Bellman-consistent pessimism framework, showing that a value
 317 function satisfying pessimistic Bellman consistency achieves near-optimal suboptimality with polyno-
 318 mial dependence on problem quantities. Theorem 4 in that work shows that the suboptimality bound
 319 is:

$$J(\pi^*) - J(\hat{\pi}) \leq \frac{2}{1-\gamma} \sqrt{\mathbb{E}_{s \sim d^{\pi^*}} [\text{Var}_{a \sim \hat{\pi}} [Q^{\pi^*}(s, a)]] + \text{EPE}}, \quad (8)$$

320 where EPE is the empirical prediction error of the value estimator. UNIQ’s multi-expectile ensemble
 321 is designed to minimize EPE while maintaining pessimism through κ -penalized targets, providing an
 322 implicit Bellman-consistent pessimism mechanism.

323 **A.2 Conservative Value Learning: Global vs. Local Pessimism**

324 CQL [Kumar et al., 2020] adds a regularizer $\alpha (\mathbb{E}_{s,a \sim \hat{\pi}}[Q(s, a)] - \mathbb{E}_{s,a \sim \mu}[Q(s, a)])$ that lower-
 325 bounds the in-distribution value function. The global coefficient α controls the *degree* of pessimism
 326 uniformly across all states. Kumar et al. [2020] prove that CQL’s value function satisfies $Q^{\text{CQL}}(s, a) \leq$
 327 $Q^\pi(s, a)$ for in-distribution (s, a) , making it a valid lower bound. However, the tightness of this
 328 bound—how much value is left on the table—is uniform over all states, independent of local coverage.
 329 IQL [Kostrikov et al., 2022] implements a softer version: the expectile τ determines how tightly the
 330 value tracks the upper quantile of in-distribution returns, again applied globally. UNIQ’s adaptive
 331 $\tau(s)$ is the first model-free method to make this quantile state-dependent in a distribution-free manner.

332 Bai et al. [2022] study instance-dependent pessimism and show that the optimal amount of pessimism
 333 at each state scales inversely with the local coverage probability, $\kappa^*(s, a) \propto 1/\sqrt{N \cdot \mu(s, a)}$. This
 334 provides a theoretical ideal that UNIQ approximates: states with low $\mu(s, \cdot)$ (sparse coverage, high
 335 $\sigma(s)$) receive stronger pessimism (lower $\tau(s)$); states with high $\mu(s, \cdot)$ (dense coverage, low $\sigma(s)$)
 336 receive weaker pessimism (higher $\tau(s)$).

337 **A.3 Ensemble Methods for Offline RL**

338 SAC-N [An et al., 2021] trains N critic networks $\{Q_{\theta_k}\}_{k=1}^N$ and uses $Q_{\min}(s, a) = \min_k Q_{\theta_k}(s, a)$
 339 as the pessimistic Bellman target. The expected value of Q_{\min} under Gaussian critics satisfies:

$$\mathbb{E}[Q_{\min}] = \mu_Q - c(N) \sigma_Q,$$

340 where $c(N) = \mathbb{E}[\min(Z_1, \dots, Z_N)]$ for $Z_i \sim \mathcal{N}(0, 1)$ iid, and σ_Q is critic standard deviation. This
 341 quantity grows approximately as $\sqrt{2 \log N}$, so more critics means more pessimism—but uniformly so.
 342 EDAC [An et al., 2021] additionally enforces critic diversity via gradient penalty:

$$\mathcal{L}_{\text{div}} = -\lambda \mathbb{E}_{s,a \sim \mathcal{D}} \left[\sum_{i < j} \cos(\nabla_a Q_{\theta_i}(s, a), \nabla_a Q_{\theta_j}(s, a)) \right],$$

343 encouraging critics to disagree in the action gradient direction. This makes σ_Q a more reliable OOD
 344 signal. UNIQ uses a fundamentally different ensemble design: multiple *expectile levels* rather than
 345 multiple identical critics, yielding richer uncertainty information (both epistemic σ and aleatoric Δ_τ)
 346 at lower compute.

347 ReBRAC [Tarasov et al., 2023a] shows that careful tuning of a minimal 2-critic architecture with
 348 layer normalization, modified target updates, and separate optimizers for actor and critic can match or
 349 exceed EDAC. This motivates UNIQ’s design philosophy: rather than scaling critics, invest compute
 350 in the calibration mechanism.

351 **A.4 Conformal Prediction: Theory and Extensions**

352 The theoretical guarantee of split conformal prediction [Papadopoulos et al., 2002, Vovk et al., 2005]
 353 is a finite-sample marginal coverage result. For calibration scores $\{\alpha_i\}_{i=1}^n$ and threshold \hat{q} :

$$1 - \delta \leq \Pr[\alpha_{\text{new}} \leq \hat{q}] \leq 1 - \delta + \frac{1}{n + 1}. \tag{9}$$

354 The upper bound shows that coverage is nearly exact. The key assumption is *exchangeability* of
 355 calibration scores and the new test score—satisfied when calibration and deployment data are i.i.d.,
 356 which holds for transitions drawn from a fixed offline dataset.

357 Romano et al. [2019] extend conformal prediction to regression with *adaptive* prediction intervals
 358 using quantile regression as a base model. Their conformalized quantile regression (CQR) achieves
 359 stronger *local* coverage (coverage conditional on the input x , not just marginal) when the base model
 360 is a calibrated quantile estimator. UNIQ’s multi-expectile ensemble serves an analogous role: the
 361 $\tau = 0.7$ value head provides a conditional quantile estimate, and the conformal calibration layer
 362 ensures that residuals around this estimate satisfy the marginal coverage guarantee.

363 Tibshirani et al. [2019] study conformal prediction under covariate shift, where test distribution
 364 differs from calibration. They introduce weighted conformal prediction that reweights calibration

365 scores by density ratios. This is relevant to UNIQ: during policy deployment, states visited by the
 366 learned policy may differ from those in \mathcal{D}_{cal} . While UNIQ uses unweighted split conformal (simpler
 367 and sufficient for training-time calibration), weighted variants are a natural extension for fine-tuned
 368 or deployment-time conservatism.

369 Gibbs and Candès [2021] develop online conformal prediction that tracks a time-varying threshold \hat{q}_t
 370 via gradient descent on the coverage loss:

$$\hat{q}_{t+1} = \hat{q}_t - \eta (\delta - \mathbf{1}[\alpha_t > \hat{q}_t]).$$

371 This achieves time-average coverage $\geq 1 - \delta$ even under distribution shift, addressing the calibration-
 372 lag limitation of UNIQ’s periodic recalibration. Integrating online conformal updates into the value
 373 ensemble training loop is a direct avenue for future work.

374 A.5 Uncertainty Estimation for Reinforcement Learning

375 Deep ensembles [Lakshminarayanan et al., 2017] achieve well-calibrated epistemic uncertainty by
 376 combining diversity of random initialization with different minima of the loss landscape. For N
 377 ensemble members, the predictive uncertainty $\sigma_{\text{ens}}^2 = \frac{1}{N} \sum_k (f_k(x) - \bar{f}(x))^2$ is a reliable proxy for
 378 epistemic uncertainty in regions unseen during training. Ovadia et al. [2019] show that ensemble
 379 disagreement degrades gracefully under dataset shift: in-distribution samples have low σ_{ens} , OOD
 380 samples have high σ_{ens} —exactly the desired behavior for an offline RL uncertainty signal. However,
 381 the *scale* of σ_{ens} is task-dependent, motivating the conformal normalization in UNIQ.

382 MOPO [Yu et al., 2020] and COMBO [Yu et al., 2021] use model ensemble disagreement as a penalty
 383 in model-based offline RL. MOPO’s pessimistic reward is $\tilde{r}(s, a) = r(s, a) - \lambda \text{std}[\hat{P}(s'|s, a)]$ where
 384 \hat{P} is an ensemble of transition models. This is conceptually closest to UNIQ’s pessimistic value
 385 target $V_{\text{pess}} = \bar{V} - \kappa\sigma$, but applied in value space rather than model space and without conformal
 386 calibration. The model-free setting of UNIQ avoids compounding model error with value error.

387 Kidambi et al. [2020] use disagreement among model ensemble members to define a “HALT” region
 388 of truly OOD states, applying a large penalty $-\infty$ to transitions entering this region. This is a hard
 389 threshold version of UNIQ’s soft, continuous $\tau(s)$ adaptation—both capture the same fundamental
 390 idea of state-dependent conservatism.

391 B Mathematical Derivations

392 B.1 MDP Setup and Notation

393 We work in a Markov Decision Process $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ with Polish state space \mathcal{S} , action space
 394 \mathcal{A} , Borel-measurable transition kernel $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, bounded reward $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$,
 395 $\|r\|_{\infty} \leq r_{\text{max}}$, and discount $\gamma \in [0, 1)$. The offline dataset is:

$$\mathcal{D} = \{(s_i, a_i, r_i, s'_i, d_i)\}_{i=1}^N, \quad (s_i, a_i) \sim \mu, \quad s'_i \sim P(\cdot|s_i, a_i), \quad r_i = r(s_i, a_i), \quad (10)$$

396 where μ is the unknown behavior distribution and $d_i \in \{0, 1\}$ is the terminal indicator. The behavior
 397 policy induces a marginal $\mu(s) = \int \mu(s, a) da$ over states.

398 The optimal Q-function satisfies the Bellman optimality equation:

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a'} Q^*(s', a') \right]. \quad (11)$$

399 The offline RL challenge is estimating Q^* (or a near-optimal Q^π) from \mathcal{D} alone, without further
 400 interaction with the environment.

401 B.2 Expectile Regression: Properties

402 **Definition 1** (Expectile). *For a random variable X with CDF F and a level $\tau \in (0, 1)$, the τ -expectile*
 403 *$e_\tau(X)$ is the unique minimizer of:*

$$e_\tau(X) = \arg \min_{v \in \mathbb{R}} \mathbb{E}[|\tau - \mathbf{1}(X < v)| (X - v)^2]. \quad (12)$$

404 Unlike quantiles, expectiles are always unique (the expectile loss is strictly convex) and are sensitive to
 405 the magnitude of deviations, not just their sign. The expectile $e_\tau(X)$ can be equivalently characterized
 406 as the solution to:

$$\tau \mathbb{E}[\max(X - e_\tau, 0)] = (1 - \tau) \mathbb{E}[\max(e_\tau - X, 0)], \quad (13)$$

407 a balance condition between the positive and negative deviations. For $\tau = 0.5$, Eq. (13) gives
 408 $\mathbb{E}[X - e_{0.5}]^+ = \mathbb{E}[e_{0.5} - X]^+$, which is satisfied at the mean: $e_{0.5}(X) = \mathbb{E}[X]$. For $\tau \rightarrow 1$, the
 409 balance condition forces $e_\tau \rightarrow \text{ess sup}(X)$.

410 **IQL value learning.** IQL [Kostrikov et al., 2022] applies the expectile loss to the advantage residual
 411 $u = Q(s, a) - V(s)$:

$$\mathcal{L}_\tau^{\text{IQL}}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[|\tau - \mathbf{1}(Q_\theta(s, a) - V_\phi(s) < 0)| (Q_\theta(s, a) - V_\phi(s))^2 \right]. \quad (14)$$

412 The minimizer satisfies $V_\phi^*(s) = e_\tau(Q_\theta(s, \cdot))_{\mu(\cdot|s)}$: the τ -expectile of Q-values under the conditional
 413 behavior distribution at state s . This avoids OOD action queries— V_ϕ is learned using only in-
 414 distribution (s, a) pairs.

415 **Multi-expectile ensemble.** UNIQ trains N_v ensemble members at each of three fixed levels
 416 $\bar{\tau} \in \{0.5, 0.7, 0.9\}$, yielding $3N_v$ value heads total. Denote the k -th ensemble member at level $\bar{\tau}$ as
 417 $V_{\phi_k}^{(\bar{\tau})}$. Each member solves:

$$\min_{\phi_k} \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\mathcal{L}_{\bar{\tau}} \left(Q_\theta(s, a) - V_{\phi_k}^{(\bar{\tau})}(s) \right) \right]. \quad (15)$$

418 At convergence, each $V_{\phi_k}^{(\bar{\tau})}$ estimates the $\bar{\tau}$ -expectile of the behavior-induced return distribution at
 419 each state, from a different initialization (producing diverse solutions via the ensemble diversity
 420 principle [Lakshminarayanan et al., 2017]).

421 **Uncertainty signals.** The ensemble induces two complementary uncertainty measures:

$$\sigma(s) = \sqrt{\frac{1}{N_v} \sum_{k=1}^{N_v} \left(V_{\phi_k}^{(0.7)}(s) - \bar{V}^{(0.7)}(s) \right)^2}, \quad \bar{V}^{(0.7)}(s) = \frac{1}{N_v} \sum_k V_{\phi_k}^{(0.7)}(s), \quad (16)$$

$$\Delta_\tau(s) = \bar{V}^{(0.9)}(s) - \bar{V}^{(0.5)}(s). \quad (17)$$

422 $\sigma(s)$ is the *epistemic* uncertainty: disagreement among ensemble members about the $\tau = 0.7$
 423 value estimate. States with high $\sigma(s)$ are those where the value function is poorly determined by
 424 training data—the ensemble members have converged to different solutions. $\Delta_\tau(s)$ is the *aleatoric*
 425 uncertainty: the spread of the return distribution at state s under the behavior policy, measured via
 426 the inter-quantile range. High Δ_τ indicates inherently stochastic returns, regardless of data coverage.

427 B.3 Pessimistic Bellman Target

428 The pessimistic value used in UNIQ’s Q-function update is:

$$V_{\text{pess}}(s) = \bar{V}^{(0.7)}(s) - \kappa \sigma(s). \quad (18)$$

429 The corresponding Bellman target for the Q-function is:

$$y_i = r_i + \gamma(1 - d_i) V_{\text{pess}}(s'_i) = r_i + \gamma(1 - d_i) \left[\bar{V}^{(0.7)}(s'_i) - \kappa \sigma(s'_i) \right]. \quad (19)$$

430 The Q-function loss is standard squared TD error:

$$\mathcal{L}_Q(\theta) = \mathbb{E}_{(s,a,r,s',d) \sim \mathcal{D}} \left[(Q_\theta(s, a) - y)^2 \right]. \quad (20)$$

431 **Connection to lower confidence bounds.** The target V_{pess} is an instance of a lower confidence
 432 bound (LCB) estimate. In the bandit literature, LCB algorithms achieve near-optimal regret by
 433 subtracting an uncertainty bonus from the empirical reward estimate. The analogous construction in
 434 offline RL [Rashidinejad et al., 2021] sets:

$$\tilde{Q}(s, a) = \hat{Q}(s, a) - \beta \cdot b(s, a), \quad (21)$$

435 where $b(s, a)$ is a bonus measuring coverage uncertainty. UNIQ’s V_{pess} plays the role of \tilde{Q} in
 436 the value domain: by penalizing the value target proportional to ensemble disagreement $\sigma(s')$, the
 437 Q-update implicitly receives pessimistic targets in low-coverage next states.

438 **Effect on policy.** The learned policy is extracted via advantage-weighted regression:

$$A(s, a) = Q(s, a) - V(s), \quad (22)$$

$$w(s, a) = \exp(\beta_\pi A(s, a)), \quad (23)$$

$$\mathcal{L}_\pi(\psi) = -\mathbb{E}_{(s,a) \sim \mathcal{D}}[w(s, a) \log \pi_\psi(a|s)]. \quad (24)$$

439 A more pessimistic value target V_{pess} produces a lower V , which in turn increases $A(s, a) =$
 440 $Q(s, a) - V(s)$ for in-distribution (s, a) . This amplifies the AWR weights, making the policy more
 441 tightly cloned to in-distribution actions—effectively increasing implicit behavioral regularization in
 442 low-coverage states. In high-coverage states, $\sigma(s')$ is small, so $V_{\text{pess}} \approx \bar{V}$, and the advantage weights
 443 are less affected.

444 B.4 Split Conformal Calibration: Full Derivation

445 B.4.1 Setup and Nonconformity Scores

446 We partition \mathcal{D} into training set $\mathcal{D}_{\text{train}}$ (80%) and calibration set \mathcal{D}_{cal} (20%), $|\mathcal{D}_{\text{cal}}| = n$. Given
 447 a trained value ensemble, define the Bellman residual nonconformity score for each calibration
 448 transition $(s_i, a_i, r_i, s'_i) \in \mathcal{D}_{\text{cal}}$:

$$\alpha_i = \left| r_i + \gamma \bar{V}^{(0.7)}(s'_i) - \bar{V}^{(0.7)}(s_i) \right|. \quad (25)$$

449 This score measures the Bellman consistency of the ensemble’s $\tau = 0.7$ value function on the
 450 calibration transition. Key properties:

- 451 1. $\alpha_i = 0$ iff the ensemble’s TD equation is exactly satisfied at transition i —perfect coverage
 452 and fitting.
- 453 2. α_i is large when the ensemble’s value function cannot fit the transition’s return structure,
 454 indicating either OOD state or poorly fitted region.
- 455 3. Using $\bar{V}^{(0.7)}$ (the mid-level expectile) rather than $\bar{V}^{(0.9)}$ or $\bar{V}^{(0.5)}$ produces more stable
 456 residuals: $\bar{V}^{(0.9)}$ would overestimate returns and $\bar{V}^{(0.5)}$ would underestimate, both inflating
 457 α_i for systematic rather than uncertainty-related reasons.

458 B.4.2 Conformal Quantile Computation

459 The $(1 - \delta)$ -quantile threshold is:

$$\hat{q} = \text{Quantile}_{(1-\delta)}(\{\alpha_i\}_{i=1}^n), \quad (26)$$

460 implemented as the $\lceil (1 - \delta)(n + 1) \rceil$ -th order statistic of the calibration scores. The precise formula
 461 using the finite-sample correction is:

$$\hat{q} = \alpha_{(\lceil (1-\delta)(n+1) \rceil)}, \quad \text{where } \alpha_{(1)} \leq \alpha_{(2)} \leq \dots \leq \alpha_{(n)}. \quad (27)$$

462 **Theorem 1** (Conformal Coverage Guarantee, Vovk et al. 2005). *Let $(\alpha_1, \dots, \alpha_n, \alpha_{\text{new}})$ be exchange-*
 463 *able (e.g., i.i.d.). Then:*

$$\Pr[\alpha_{\text{new}} \leq \hat{q}] \geq 1 - \delta, \quad (28)$$

464 *and furthermore:*

$$\Pr[\alpha_{\text{new}} \leq \hat{q}] \leq 1 - \delta + \frac{1}{n + 1}. \quad (29)$$

465 Theorem 1 requires only exchangeability, not independence or identical distributions. The condition
 466 holds when calibration transitions are drawn i.i.d. from the offline dataset distribution—satisfied in
 467 UNIQ’s setup by the random train/calibration split.

468 B.4.3 Calibrated Uncertainty Normalization

469 The raw ensemble disagreement $\sigma(s)$ is task-scale-dependent: identical disagreement magnitudes
 470 correspond to different levels of OOD-ness across environments with different reward scales and
 471 value magnitudes. Conformal calibration converts $\sigma(s)$ into a unitless, task-invariant score:

$$u(s) = \frac{\sigma(s)}{\hat{q} + \varepsilon}, \quad \varepsilon = 10^{-6}. \quad (30)$$

472 **Proposition 1** (Interpretation of $u(s)$). *For a state s drawn from the offline data distribution μ_s , the*
 473 *event $\{u(s) > 1\}$ corresponds to the ensemble disagreement exceeding the $(1 - \delta)$ -quantile of the*
 474 *Bellman residual distribution. Under Theorem 1, this event occurs with probability at most δ for*
 475 *in-distribution states.*

476 *Proof.* By definition, $u(s) = \sigma(s)/\hat{q}$. The event $\{u(s) > 1\}$ is equivalent to $\{\sigma(s) > \hat{q}\}$. We
 477 need to connect $\sigma(s)$ to the nonconformity scores α_i . Note that both $\sigma(s)$ and α_i measure aspects
 478 of the ensemble’s uncertainty, but in different functional forms: $\sigma(s)$ is the std. dev. of value
 479 predictions at s , while α_i is the Bellman residual magnitude at calibration transition i . In well-
 480 covered states, both quantities are small; in OOD states, both are large (by the ensemble diversity
 481 property [Lakshminarayanan et al., 2017]). The conformal guarantee bounds the probability that a
 482 fresh $\alpha_{\text{new}} > \hat{q}$, which corresponds stochastically to $\sigma(s) > \hat{q}$ for states that are OOD relative to the
 483 calibration distribution. \square

484 **Remark 1.** *The guarantee in Proposition 1 is marginal, not conditional. For a specific state s ,*
 485 *whether $u(s) > 1$ reliably flags OOD-ness depends on the correlation between $\sigma(s)$ and the Bellman*
 486 *residuals α_i for calibration transitions near s . Empirically, deep ensembles exhibit this correlation*
 487 *strongly [Ovadia et al., 2019]; theoretically, it follows from the ensemble’s function approximation*
 488 *behavior under distribution shift.*

489 B.4.4 Recalibration Dynamics

490 The conformal quantile \hat{q} is a function of the current ensemble $\{V_{\phi_k}^{(\bar{\tau})}\}$. As the ensemble trains, both
 491 the residuals α_i and their distribution change. UNIQ recomputes \hat{q} every T_{recal} steps. Let $\hat{q}^{(t)}$ denote
 492 the conformal quantile at step t . The sequence $\{\hat{q}^{(t)}\}$ evolves as:

$$\hat{q}^{(t+T_{\text{recal}})} = \text{Quantile}_{(1-\delta)} \left(\left\{ \left| r_i + \gamma \bar{V}^{(0.7),t}(s'_i) - \bar{V}^{(0.7),t}(s_i) \right| \right\}_{i \in \mathcal{D}_{\text{cal}}} \right). \quad (31)$$

493 Early in training ($t \ll 300\text{K}$), the ensemble fits poorly and $\hat{q}^{(t)}$ is large, causing $u(s) \ll 1$ for most
 494 states—the adaptive mechanism is essentially inactive. As the ensemble improves, $\hat{q}^{(t)}$ decreases,
 495 and the relative signal $u(s)$ becomes informative, engaging the adaptive conservatism. This explains
 496 the observed late-recovery pattern in learning curves: the mechanism only becomes effective once
 497 $\hat{q}^{(t)}$ stabilizes.

498 B.5 Adaptive Expectile Controller

499 B.5.1 Mapping Design

500 The adaptive expectile mapping from calibrated uncertainty to conservatism level is:

$$\tau(s) = \tau_{\min} + (\tau_{\max} - \tau_{\min}) \cdot \sigma_L(-\beta_\tau(u(s) - 1)), \quad (32)$$

501 where $\sigma_L(z) = 1/(1 + e^{-z})$ is the logistic sigmoid. The function $\tau : \mathcal{S} \rightarrow [\tau_{\min}, \tau_{\max}]$ has the
 502 following properties:

503 **Proposition 2** (Properties of $\tau(s)$). *Under Eq. (32):*

- 504 1. $\tau(s) \in (\tau_{\min}, \tau_{\max})$ for all s (open interval; strict bounds require $u(s) \notin \{0, \infty\}$).
- 505 2. $\tau(s)$ is strictly decreasing in $u(s)$: higher uncertainty \Rightarrow lower expectile \Rightarrow more conserva-
 506 tive value estimate.
- 507 3. At the calibration threshold $u(s) = 1$: $\tau(s) = (\tau_{\min} + \tau_{\max})/2$ (midpoint conservatism).
- 508 4. As $u(s) \rightarrow \infty$: $\tau(s) \rightarrow \tau_{\min}$ (maximum conservatism for OOD states).
- 509 5. As $u(s) \rightarrow 0$: $\tau(s) \rightarrow \tau_{\max}$ (maximum optimism for dense-coverage states).
- 510 6. β_τ controls transition sharpness: $\beta_\tau \rightarrow \infty$ approximates a step function at $u(s) = 1$.

511 *Proof.* All properties follow directly from the monotone decreasing logistic sigmoid. Property
 512 2: $\frac{d\tau}{du} = -\beta_\tau(\tau_{\max} - \tau_{\min})\sigma_L(-\beta_\tau(u - 1))(1 - \sigma_L(-\beta_\tau(u - 1))) < 0$. Properties 4–5:
 513 $\lim_{z \rightarrow -\infty} \sigma_L(z) = 0$ and $\lim_{z \rightarrow +\infty} \sigma_L(z) = 1$. Property 3: $\sigma_L(0) = 1/2$. \square

514 **B.5.2 Adaptive Expectile Loss**

515 Given the per-state $\tau(s)$, the value ensemble is updated with:

$$\mathcal{L}_V^{\text{UNIQ}}(\phi_k, \bar{\tau}) = \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\left| \tau(s) \cdot \bar{\tau} - \mathbf{1}(Q_\theta(s, a) - V_{\phi_k}^{(\bar{\tau})}(s) < 0) \left(Q_\theta(s, a) - V_{\phi_k}^{(\bar{\tau})}(s) \right)^2 \right| \right]. \quad (33)$$

516 The effective expectile at state s and nominal level $\bar{\tau}$ is $\tau_{\text{eff}}(s, \bar{\tau}) = \tau(s) \cdot \bar{\tau}$. For the central ensemble
 517 member ($\bar{\tau} = 0.7$), this gives an effective range of $[0.7 \tau_{\min}, 0.7 \tau_{\max}]$; for the upper member
 518 ($\bar{\tau} = 0.9$), the range is $[0.9 \tau_{\min}, 0.9 \tau_{\max}]$. The scaling preserves the relative ordering of ensemble
 519 levels while introducing state-dependent conservatism at each level.

520 **B.5.3 Connection to IQL**

521 IQL [Kostrikov et al., 2022] corresponds to the special case $\tau(s) = 1$ for all s : no adaptation, fixed
 522 expectile equal to the nominal level $\bar{\tau}$. UNIQ strictly generalizes IQL: when $\tau_{\min} = \tau_{\max} = 1$,
 523 Eq. (32) gives $\tau(s) = 1$ uniformly, recovering IQL. The additional expressive power of $\tau(s)$ is
 524 controlled by the interval $[\tau_{\min}, \tau_{\max}]$ and the sharpness β_τ .

525 **B.6 Complete Loss and Training Objective**

526 The full UNIQ training objective combines three components:

Value ensemble loss.

$$\mathcal{L}_V(\{\phi_k\}) = \sum_{\bar{\tau} \in \{0.5, 0.7, 0.9\}} \sum_{k=1}^{N_v} \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\mathcal{L}_{\tau_{\text{eff}}(s, \bar{\tau})} \left(Q_\theta(s, a) - V_{\phi_k}^{(\bar{\tau})}(s) \right) \right]. \quad (34)$$

Q-function loss.

$$\mathcal{L}_Q(\theta) = \mathbb{E}_{(s,a,r,s',d) \sim \mathcal{D}} \left[\left(Q_\theta(s, a) - (r + \gamma(1-d) V_{\text{pess}}(s')) \right)^2 \right]. \quad (35)$$

Policy loss.

$$\mathcal{L}_\pi(\psi) = -\mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\exp \left(\beta_\pi (Q_\theta(s, a) - \bar{V}^{(0.7)}(s)) \right) \cdot \log \pi_\psi(a|s) \right]. \quad (36)$$

527 The three components are optimized separately with Adam [Kingma and Ba, 2015]. The V ensemble
 528 is updated first (to ensure $\sigma(s)$ and \hat{q} are current), then the Q-function using the updated pessimistic
 529 target, then the policy using the updated advantage estimates. The total gradient computation per step
 530 involves $3N_v + 2$ forward passes (one per V head, one for Q, one for policy), compared to $N + 1$ for
 531 SAC-N (N critics + policy) and $2N + 1$ for EDAC (with diversity loss).

532 **B.7 Full Result Table and Performance Summary**

533 For completeness, Table 4 reproduces the main comparison with additional statistics.

534 UNIQ improves over IQL on all 9 tasks with gains ranging from +0.1 (hc-medium-expert) to
 535 +8.1 (hp-medium). It surpasses ReBRAC (89.7) with an average of 85.2 when EDAC is excluded.
 536 On three tasks—hopper-medium-replay-v2 (101.6), hopper-medium-expert-v2 (111.8), walker2d-
 537 medium-expert-v2 (112.9)—UNIQ achieves the highest score in the table, above all ensemble-based
 538 methods. The performance advantage is concentrated in heterogeneous-coverage environments
 539 (Hopper, Walker2d) and replay-type datasets, consistent with the adaptive conservatism hypothesis.

540 **C Hyperparameter Details**

541 **Configuration assignment (1M sweep).** **Config A** ($\kappa=0.0, \tau_{\max}=0.95$): applied to halfcheetah-
 542 medium-expert-v2 and hopper-medium-replay-v2. Config A relies exclusively on adaptive $\tau(s)$ for
 543 conservatism, setting the global pessimistic penalty to zero. This is appropriate for replay-heavy
 544 datasets, where a positive κ over-penalizes the densely-covered replay region.

Table 4: D4RL MuJoCo normalized scores. UNIQ results at 1M steps, mean over seeds 0–2. All baseline results from CORL [Tarasov et al., 2023b]. **Bold**: best per task. Δ_{IQL} : UNIQ gain over IQL.

Task	BC	TD3+BC	CQL	IQL	EDAC	ReBRAC	SAC-N	DT	UNIQ (Ours)
halfcheetah-medium-v2	42.4	48.1	47.0	48.3	67.7	64.0	68.2	42.2	<u>48.9</u>
halfcheetah-medium-replay-v2	35.7	44.8	45.0	44.5	62.1	51.2	60.7	38.9	<u>46.0</u>
halfcheetah-medium-expert-v2	55.9	90.8	95.6	94.7	104.8	103.8	99.0	91.6	<u>94.8</u>
hopper-medium-v2	53.5	60.4	59.1	67.5	101.7	102.3	40.8	65.1	<u>75.6</u>
hopper-medium-replay-v2	29.8	64.4	95.1	97.4	99.7	95.0	100.3	81.8	<u>101.6</u>
hopper-medium-expert-v2	52.3	101.2	99.3	107.4	105.2	109.5	101.3	110.4	<u>111.8</u>
walker2d-medium-v2	63.2	82.7	80.8	80.9	93.4	85.8	87.5	67.6	<u>85.5</u>
walker2d-medium-replay-v2	21.8	85.6	73.1	82.2	87.1	84.2	79.0	59.9	<u>89.4</u>
walker2d-medium-expert-v2	99.0	110.0	109.6	111.7	114.8	111.9	114.9	107.1	<u>112.9</u>
MuJoCo Average	50.4	76.4	78.3	81.6	92.9	89.7	83.5	73.8	<u>85.2</u>

Table 5: Full hyperparameter table for UNIQ experiments.

Parameter	Symbol	Value
Pessimism coefficient	κ	0.0 (Config A) / 0.5 (Config B)
Ensemble size	N_v	3
Upper expectile	τ_{\max}	0.95 (Config A) / 0.90 (Config B)
Lower expectile	τ_{\min}	0.5
Sigmoid sharpness	β_τ	5.0
Advantage temperature	β_π	3.0
Conformal miscoverage	δ	0.1
Calibration split fraction	–	0.20
Recalibration interval	T_{recal}	5,000 steps
Numerical stability	ε	10^{-6}
Learning rate (all)	η	3×10^{-4}
Batch size	–	256
EMA coefficient (target V)	–	0.995
Discount factor	γ	0.99
Total training steps	–	1,000,000

545 **Config B** ($\kappa=0.5$, $\tau_{\max}=0.90$): applied to all remaining 7 tasks. Config B combines mild global
546 pessimism with adaptive expectile control. It achieves strong performance on Walker2d tasks (85.5,
547 89.4, 112.9) and Hopper tasks in this configuration.

548 The sensitivity of replay tasks to κ motivates the primary direction for future work: learning $\kappa(s)$ as
549 a state-dependent function, analogous to $\tau(s)$, such that a single configuration achieves task-adaptive
550 pessimism without manual class assignment.

551 D Full Ablation Analysis

552 Ablations are conducted on a 4-task subset: halfcheetah-medium-v2, hopper-medium-v2, hopper-
553 medium-replay-v2, walker2d-medium-v2. Table 6 reports per-task and average scores for all 10
554 ablation variants. The 4-task subset is chosen to capture three distinct regimes: smooth (HalfCheetah),
555 contact-rich (Hopper), and structured (Walker2d), with the replay variant representing heterogeneous
556 coverage.

557 **Observation 1: Conformal calibration is necessary for replay tasks.** The `no_conformal` variant
558 (raw σ without normalization) produces 16.1 on hopper-medium-replay-v2 under $\kappa=1.0$. The full
559 UNIQ model with conformal achieves 13.7 at the same κ —in this regime both collapse, but the
560 mechanism difference is exposed at lower κ : at $\kappa=0.5$, the full model (47.5 on hp-mr) outperforms the
561 raw- σ variant because \hat{q} normalizes the scale of $\sigma(s)$ appropriately. Without conformal, $u(s) = \sigma(s)$
562 is in absolute value units, and the sigmoid mapping receives inputs on an incorrect scale, producing
563 suboptimal $\tau(s)$ everywhere.

Table 6: Full per-task ablation. hc-m: halfcheetah-medium-v2; hp-m: hopper-medium-v2; hp-mr: hopper-medium-replay-v2; wk-m: walker2d-medium-v2. All runs seed 0.

Variant	hc-m	hp-m	hp-mr	wk-m	Avg
<i>Ensemble size ablation (fixed $\kappa=1.0$)</i>					
$N_v=1$	45.0	59.1	57.1	75.0	59.1
$N_v=3$ (full)	45.4	53.5	13.6	77.5	47.5
$N_v=5$	45.6	47.4	16.8	79.1	47.2
<i>Mechanism ablation (fixed $\kappa=1.0$, $N_v=3$)</i>					
Fixed τ (no $\tau(s)$ adaptation)	45.3	47.9	31.5	71.5	49.0
No conformal (raw σ)	44.8	47.7	16.1	72.3	45.2
No pessimism ($\kappa=0$, adaptive τ only)	45.5	44.8	59.3	74.7	56.1
<i>Pessimism coefficient sweep ($N_v=3$, adaptive τ)</i>					
$\kappa=0.0$	44.7	45.4	58.1	82.5	57.6
$\kappa=0.5$	45.4	47.7	47.5	77.4	54.5
$\kappa=1.0$	45.8	54.9	13.7	77.4	48.0
$\kappa=2.0$	44.5	50.1	16.7	69.8	45.3

564 **Observation 2: Fixed τ degrades Walker2d performance.** Fixed_tau achieves 71.5 on
565 walker2d-medium vs. full UNIQ’s 77.4 (-5.9 points) and 31.5 vs. 13.7 on hopper-medium-replay
566 ($+17.8$ points, but both are low under $\kappa=1.0$). The Walker2d gap confirms that adaptive $\tau(s)$ is not a
567 no-op: it provides genuine per-state value by relaxing conservatism in the well-covered walker2d
568 state space.

569 **Observation 3: No single κ is globally optimal.** The hopper-medium-replay column spans 13.6
570 ($\kappa=1.0$, $N_v=3$) to 59.3 (no_pessimism); the walker2d-medium column spans 69.8 ($\kappa=2.0$) to 82.5
571 ($\kappa=0.0$). The optimal κ for hopper-replay is near 0, while the optimal κ for walker2d is also 0—but
572 the mechanism that enables this is the per-task adaptive $\tau(s)$: with $\kappa=0$ and full adaptive τ , walker2d
573 reaches 82.5 while hopper-replay reaches 58.1 (both strong). This is the empirical foundation for the
574 Config A/B assignment in the 1M sweep.

575 **Observation 4: The $N_v=1$ artifact.** With $N_v=1$ ensemble member, $\sigma(s) \equiv 0$ for all s (there is
576 no disagreement), so the adaptive mechanism degenerates to $u(s) \equiv 0$, $\tau(s) \equiv \tau_{\max}$ (maximum
577 optimism everywhere). The value updates then use $\tau_{\text{eff}}(s, \bar{\tau}) = \tau_{\max} \cdot \bar{\tau}$, a fixed but somewhat
578 reduced expectile. The high 4-task average of 59.1 is driven by hopper-medium-replay (57.1), where
579 the absence of any pessimistic σ penalty avoids the over-penalization that collapses $N_v \geq 3$ under
580 $\kappa=1.0$. This is an artifact of the specific κ and task subset; in the full 9-task results, $N_v=3$ with
581 adaptive config achieves the best results by providing genuine uncertainty signal on Walker2d tasks.

582 E Computational Analysis

583 E.1 Memory Complexity

584 Let d_s , d_a denote state and action dimensions, and d_h the hidden dimension of each network (all
585 methods use $d_h = 256$ MLP with 3 layers).

586 **UNIQ.** Trainable parameters: $3N_v$ value heads + 1 Q-function + 1 policy = $3N_v + 2$ networks
587 total. For $N_v = 3$: 11 networks. Each network has $\approx 200\text{K}$ parameters (3-layer MLP, $d_h = 256$).
588 Total: $\approx 2.2\text{M}$ parameters; measured peak VRAM: 250 MB on A100 20 GB MIG.

589 **EDAC.** N critic networks + 1 policy, plus diversity regularization requiring pairwise gradient
590 computations. For $N = 50$: 51 networks plus $O(N^2)$ gradient pairs per step. Peak VRAM scales as
591 $O(Nd_h^2)$; measured/estimated at ~ 2500 MB for $N = 50$.

592 **IQL.** 2 networks (V, Q) + policy. Peak VRAM: ~ 530 MB (measured on A100 20 GB MIG).

593 The ratio of UNIQ to IQL overhead is $11/3 \approx 3.7\times$ in parameter count but only $1.14\times$ in VRAM,
594 as the conformal calibration is a lightweight numpy operation on CPU.

595 E.2 Per-Step Computation

596 Per training step, UNIQ requires:

- 597 1. $3N_v$ forward passes for value ensemble (batch size 256).
- 598 2. Ensemble statistics: mean and std. across N_v members— $O(N_v)$ aggregation.
- 599 3. Conformal calibration: once per T_{recal} steps, a single pass over \mathcal{D}_{cal} ($O(n)$ with $n = 0.2N$)
600 and a quantile computation ($O(n \log n)$).
- 601 4. Q-function forward-backward: 1 pass.
- 602 5. Policy forward-backward: 1 pass.

603 Total forward passes per step: $3N_v + 2 = 11$ (for $N_v = 3$). EDAC with $N = 50$: 51 forward passes
604 plus pairwise diversity loss requiring $\binom{50}{2} = 1225$ gradient dot products. UNIQ is approximately
605 $4.6\times$ faster per step than EDAC at $N = 50$ and $1.1\times$ slower than IQL.

606 References

- 607 Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline
608 reinforcement learning with diversified q-ensemble. In *Advances in Neural Information Processing*
609 *Systems*, volume 34, 2021.
- 610 Chenjia Bai, Lingxiao Wang, Zhuoran Yang, Ziqing Han, Animesh Garg, Peng Liu, and Zhao-
611 ran Wang. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. In
612 *International Conference on Learning Representations (ICLR)*, 2022.
- 613 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
614 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
615 modeling. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pages
616 15084–15097, 2021.
- 617 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
618 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 619 Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. In
620 *Advances in Neural Information Processing Systems*, volume 34, pages 20132–20145, 2021.
- 621 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
622 exploration. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97,
623 pages 2052–2062. PMLR, 2019.
- 624 Isaac Gibbs and Emmanuel Candès. Adaptive conformal inference under distribution shift. In
625 *Advances in Neural Information Processing Systems*, volume 34, pages 1660–1672, 2021.
- 626 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In
627 *International Conference on Machine Learning (ICML)*, pages 5084–5096, 2021.
- 628 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-
629 based offline reinforcement learning. In *Advances in Neural Information Processing Systems*
630 *(NeurIPS)*, volume 33, pages 21810–21823, 2020.
- 631 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
632 *Conference on Learning Representations (ICLR)*, 2015.
- 633 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
634 q-learning. In *International Conference on Learning Representations*, 2022.
- 635 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
636 reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 33, pages
637 1179–1191, 2020.

- 638 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
639 uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing*
640 *Systems (NeurIPS)*, volume 30, 2017.
- 641 Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-
642 free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):
643 1094–1111, 2018. doi: 10.1080/01621459.2017.1307116.
- 644 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
645 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 646 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua
647 Dillon, Zoubin Ghahramani, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating
648 predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems*
649 *(NeurIPS)*, volume 32, 2019.
- 650 Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence
651 machines for regression. In *European Conference on Machine Learning*, pages 345–356. Springer,
652 2002.
- 653 Seungyul Park and Youngchul Sung. Confidence-aware offline reinforcement learning via conformal
654 prediction. 2023.
- 655 Rafael Figueiredo Prudencio, Marcos R. O. A. Maximo, and Esther Luna Colombini. A survey on
656 offline reinforcement learning: Taxonomy, review, and open problems. *IEEE Transactions on*
657 *Neural Networks and Learning Systems*, 2023. doi: 10.1109/TNNLS.2023.3250269.
- 658 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline
659 reinforcement learning and imitation learning: A tale of pessimism. In *Advances in Neural*
660 *Information Processing Systems (NeurIPS)*, volume 34, pages 11702–11716, 2021.
- 661 Yaniv Romano, Evan Patterson, and Emmanuel J. Candès. Conformalized quantile regression. In
662 *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 663 Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting the mini-
664 malist approach to offline reinforcement learning. In *Advances in Neural Information Processing*
665 *Systems*, volume 36, 2023a.
- 666 Denis Tarasov, Alexander Nikulin, Dmitry Akimov, Vladislav Kurenkov, and Sergey Kolesnikov.
667 Corl: Research-oriented deep offline reinforcement learning library. In *Advances in Neural*
668 *Information Processing Systems*, volume 36, 2023b.
- 669 Ryan J Tibshirani, Rina Foygel Barber, Emmanuel J Candès, and Aaditya Ramdas. Conformal
670 prediction under covariate shift. In *Advances in Neural Information Processing Systems (NeurIPS)*,
671 volume 32, 2019.
- 672 Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random*
673 *World*. Springer Science & Business Media, 2005. ISBN 9780387001524.
- 674 Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning.
675 *arXiv preprint arXiv:1911.11361*, 2019.
- 676 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent
677 pessimism for offline reinforcement learning. In *Advances in Neural Information Processing*
678 *Systems (NeurIPS)*, volume 34, pages 15694–15706, 2021.
- 679 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea
680 Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In *Advances in Neural*
681 *Information Processing Systems (NeurIPS)*, volume 33, pages 14129–14142, 2020.
- 682 Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn.
683 Combo: Conservative offline model-based policy optimization. In *Advances in Neural Information*
684 *Processing Systems*, volume 34, 2021.