
From Static Policies to Adaptive Priors in Offline Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Offline reinforcement learning (RL) has traditionally focused on learning policies
2 for direct deployment under conservative objectives, where uncertainty outside the
3 offline dataset is treated pessimistically to ensure robustness. We argue that this
4 formulation becomes incomplete when an offline-trained policy is subsequently
5 updated through online interaction, as increasingly occurs in modern intelligent
6 systems through test-time adaptation and online fine-tuning. This position paper
7 argues that, in such settings, the objective of offline RL should extend beyond im-
8 mediate deployment and instead prioritize learning *adaptive policy priors*: policies
9 that preserve the capacity to improve during subsequent interaction through mem-
10 ory, exploration, and self-correction. We formalize this perspective as *adaptive*
11 *offline reinforcement learning* (AORL), distinguish it from offline-to-online RL,
12 and explain why adaptability becomes important under distributional shift, limited
13 dataset coverage, and changing test-time conditions. We further discuss Bayesian
14 offline RL as one principled direction for constructing adaptive policy priors by
15 preserving epistemic uncertainty over plausible environments. Finally, we outline
16 connections, open challenges, and research directions for treating offline RL as
17 preparation for future experience rather than as a static deployment problem.

18 1 Introduction

19 Offline reinforcement learning (RL) studies how to optimize a policy using a static dataset of
20 trajectories, without relying on online interaction during training¹. In its classical formulation (Levine
21 et al., 2020), the objective is to learn a policy for *direct deployment*: the offline-learned policy is treated
22 as a static, final decision rule whose performance is expected to derive entirely from offline data rather
23 than from subsequent interaction. This makes distributional shift particularly challenging (Kumar
24 et al., 2019): out-of-dataset actions compound errors over time with no mechanism for correction.

25 As a result, the dominant design principle of offline RL has been *conservatism*: uncertainty about
26 out-of-dataset actions is treated pessimistically so that learned policies remain safe and robust
27 under limited support. This principle is realized through various objectives and constraints (Kumar
28 et al., 2020; Yu et al., 2020; Fujimoto & Gu, 2021), and has driven much of the progress in offline
29 RL (Reed et al., 2022; Lee et al., 2022a; Kumar et al., 2023; Park et al., 2025b).

30 A fundamental limitation is that this leaves little room for improvement after offline learning, even
31 though modern systems routinely continue learning after deployment through in-context learn-
32 ing (Brown et al., 2020; Wei et al., 2022), planning (Yao et al., 2023; Snell et al., 2025), or online RL
33 fine-tuning (Guo et al., 2025a,b). In the *era of experience* (Silver & Sutton, 2025), agent capability
34 depends not only on static data but on acquiring new information through interaction (Hadsell et al.,

¹We refer to “training” as parameter updates in the policy, in contrast to in-context learning (Brown et al., 2020).

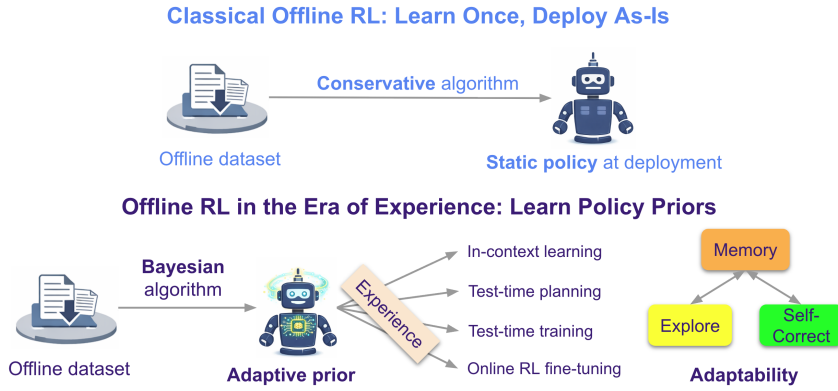


Figure 1: Top: **Classical offline RL** learns a static policy from a static dataset under conservative objectives. Bottom: **Adaptive offline RL** learns a policy prior that remains improvable through subsequent experience, including in-context learning, planning, or online fine-tuning. Here, **adaptability** arises from the interaction between memory, exploration, and self-correction. Bayesian perspective is highlighted as one principled direction for constructing such adaptive priors.

35 2020; Hughes et al., 2024), yet excessive conservatism makes this difficult by suppressing actions
 36 that interaction would reveal to be beneficial.

37 In this position paper, we use the term *adaptive offline reinforcement learning* (AORL) to refer to
 38 settings where offline RL should prioritize adaptability beyond direct deployment. **Our position**
 39 **is that when policies will continue improving through interaction, offline RL should learn**
 40 **adaptive policy priors rather than static policies.** Under this view, offline learning should not
 41 eliminate uncertainty solely for immediate robustness, but preserve sufficient behavioral flexibility
 42 for later adaptation. Such adaptability requires three ingredients: *memory*, so that decisions depend
 43 on online history; *exploration*, so that uncertain but potentially valuable actions remain reachable;
 44 and *self-correction*, so that new evidence can revise early mistakes during interaction. Actions
 45 outside the dataset are therefore not inherently undesirable; rather, they reflect epistemic uncertainty
 46 that can later be resolved through in-context learning, planning, or fine-tuning. Fig. 1 summarizes
 47 this shift from static policies to adaptive policy priors.

48 One principled direction for AORL arises from Bayesian perspectives (Ghosh et al., 2022; Ni et al.,
 49 2025). In this view, offline RL is cast as an epistemic POMDP (Ghosh et al., 2021): limited dataset cov-
 50 erage induces a posterior distribution over plausible MDPs that agree on observed data while differing
 51 beyond dataset support. The resulting optimal policy is naturally *adaptive in context*: by conditioning
 52 on online history, it can first explore uncertain but potentially good actions, then exploit the most
 53 promising ones. Bayesian offline learning therefore offers a concrete interpretation of adaptive policy
 54 priors whose value emerges through test-time interaction rather than as static behavior at deployment.

55 2 Formulation of Adaptive Offline Reinforcement Learning

56 We consider the standard offline RL setting (Levine et al., 2020) for a discrete-time, infinite-horizon,
 57 discounted-reward MDP defined by the tuple $\mathcal{M}^* = (\mathcal{S}, \mathcal{A}, \rho^*, P^*, R^*, \gamma)$. Here, the state space \mathcal{S} ,
 58 the action space \mathcal{A} , and the discount factor $\gamma \in (0, 1)$ are assumed known, while the environment
 59 components, including the initial state distribution $\rho^* \in \Delta(\mathcal{S})$, the transition function $P^* : \mathcal{S} \times \mathcal{A} \rightarrow$
 60 $\Delta(\mathcal{S})$, and reward function $R^* : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathbb{R})$, are *unknown*.

61 Instead of interacting with \mathcal{M}^* during training, the offline learner receives a static offline dataset
 62 $\mathcal{D} = \{\tau^i\}_{i=1}^N$ of trajectories collected by an unknown behavior policy β . Each trajectory
 63 $\tau = (s_0, a_0, r_1, s_1, a_1, r_2, \dots)$ is generated by $s_0 \sim \rho^*$, $a_t \sim \beta(h_t)$, $s_{t+1} \sim P^*(s_t, a_t)$, $r_{t+1} \sim$
 64 $R^*(s_t, a_t)$, $\forall t \geq 0$, where the behavior policy may depend on the interaction history. We define
 65 the *history* at time t as $h_t := (s_0, a_0, r_1, s_1, \dots, a_{t-1}, r_t, s_t)$, with $h_t \in \mathcal{H}_t$ and \mathcal{H}_t denoting the
 66 corresponding history space. The *ideal goal* for offline RL is to find a policy that maximizes expected
 67 discounted return under the true environment: $\max_{\pi} J(\pi; \mathcal{M}^*) := \mathbb{E}_{\tau} [\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid \pi, \mathcal{M}^*]$. In
 68 the most general form, this policy may depend on interaction history, denoted as $\pi : \mathcal{H}_t \rightarrow \Delta(\mathcal{A})$.

69 Because \mathcal{M}^* is inaccessible during offline optimization, this objective cannot be optimized directly.
 70 The key difficulty is epistemic uncertainty in state-action regions poorly covered by \mathcal{D} . We define the
 71 **behavioral support** of a state s within \mathcal{D} as $\text{supp}_{\mathcal{D}}(s) = \{a \mid \Pr_{\mathcal{D}}(a \mid s) > \epsilon\} \subseteq \mathcal{A}$, where
 72 $\Pr_{\mathcal{D}}(\cdot \mid s)$ denotes the empirical action distribution in \mathcal{D} and $\epsilon > 0$ is a small threshold. For states

73 s not observed in the dataset, $\text{supp}_{\mathcal{D}}(s) = \emptyset$. Epistemic uncertainty therefore remains substantial
 74 over actions in $\mathcal{A} \setminus \text{supp}_{\mathcal{D}}(s)$ for each $s \in \mathcal{S}$. Different treatments of this uncertainty give rise to
 75 different algorithmic principles, including conservative offline RL (Levine et al., 2020), Bayesian
 76 offline RL (Ghosh et al., 2022), and optimistic offline RL (Agarwal et al., 2020).

77 2.1 Core Components of Adaptive Policy Priors

Table 1: Comparison between classical offline RL and adaptive offline RL.

Property	Classical Offline RL	Adaptive Offline RL
Decision rule	Markovian	History-dependent
Uncertainty treatment	Within offline support	Preserve flexibility
Design focus	Offline stage	Offline + test-time stages
Core capabilities	Safety, robustness	Memory, exploration, self-correction

78 **History dependence in decision making (memory).** Classical offline RL typically assumes a
 79 Markovian decision rule, where actions are selected solely from the current state through a policy
 80 $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ (Levine et al., 2020). The offline-learned Markovian policy remains fixed during
 81 deployment. In adaptive offline RL (AORL), by contrast, effective decision making should depend
 82 on online interaction history, because limited offline coverage leaves uncertainty about environment
 83 dynamics. This dependence can be implemented *explicitly* through a history-dependent policy $\pi : \mathcal{H}_t \rightarrow \Delta(\mathcal{A})$ (Chen et al., 2021b), or *implicitly* through online adaptation of latent variables (Ghosh
 84 et al., 2022; Liu et al., 2023), plans (Janner et al., 2022), or even policy parameters (Xu et al., 2025).

86 **Behavioral flexibility (exploration and self-correction).** Classical offline RL typically restricts
 87 learned actions to the behavioral support $\text{supp}_{\mathcal{D}}(s)$ or heavily penalizes actions outside it (Levine
 88 et al., 2020). In AORL, by contrast, out-of-support actions are not treated as intrinsically undesirable:
 89 some may be potentially useful but remain unresolved under offline data alone because of epistemic
 90 uncertainty rather than evidence of poor value. Preserving non-negligible probability on such actions
 91 allows the policy to explore uncertain behaviors during interaction. Equally importantly, because
 92 these actions may fail, the policy should be able to self-correct by using online feedback to revise
 93 future decisions rather than committing to early mistakes (Wang et al., 2024; Kumar et al., 2025).

94 These components define the basic requirements for an **adaptive policy prior**: the policy should
 95 retain memory from interaction history while preserving sufficient flexibility to explore uncertain
 96 behaviors and correct them online. Table 1 summarizes the key differences.

97 **Forms of adaptation.** Adaptation mechanisms differ in whether policy parameters are updated and
 98 how much online data they require. Common forms are summarized in Table 2.

Table 2: Common forms of adaptation after offline optimization. See Sec. A for details.

Form of Adaptation	Online Data	Parameter Update?	Mechanism
Test-time in-context learning	Limited	No	Implicit inference
Test-time planning	Limited	No	Explicit inference
Test-time training	Limited	Yes	Temporary update
Online RL fine-tuning	Extended	Yes	Persistent update

99 3 Why Offline RL Needs Adaptive Policy Priors?

100 The previous section defined adaptive policy priors through memory and behavioral flexibility. We
 101 now ask why these properties are needed beyond the classical offline RL objective. The key reason is
 102 that uncertainty often persists beyond offline optimization, appearing in online environments through
 103 distributional shift, limited data coverage, exploration bottlenecks, or changing conditions.

104 **Inevitable distributional shift.** Classical offline RL reduces risk by keeping actions within
 105 behavioral support, yet out-of-distribution (OOD) states may still arise during deployment because
 106 small errors accumulate over time. This is well known in imitation learning, where compounding
 107 errors induce *covariate shift* (Ross & Bagnell, 2010), and remains a practical bottleneck in offline
 108 RL (Park et al., 2024a). Since such OOD states are rarely trained explicitly under conservative
 109 objectives, the learned policy may become brittle when uncertainty matters most. Adaptability
 110 helps address this by preparing the offline-learned policy to use memory and online feedback to
 111 **self-correct** once OOD states are encountered.

112 **Limited high-quality coverage.** Classical offline RL effectively bounds policy improvement by
 113 the quality of behaviors represented in the dataset (Jin et al., 2021; Uehara & Sun, 2022). This

114 limitation becomes pronounced when high-quality actions are absent or underrepresented (Ni et al.,
115 2025). In practice, collecting high-quality trajectories is expensive, and in open-ended domains
116 even carefully curated datasets may remain incomplete. Adaptability mitigates this limitation by
117 preserving behavioral flexibility, allowing subsequent interaction to **explore** uncertain actions and
118 reinforce those that prove beneficial.

119 **Exploration bottleneck in online RL fine-tuning.** A standard way to overcome the limited adaptabil-
120 ity of offline RL is to continue learning through online RL fine-tuning (Nair et al., 2020). However,
121 in classical offline-to-online RL, online improvement is often slow and incremental (Luo et al., 2023;
122 Zhao et al., 2023). In foundation model post-training, recent evidence further suggests that RL often
123 mainly reweights behaviors already present in the policy prior rather than reliably discovering new
124 ones (Yue et al., 2025; Zhao et al., 2025). As a result, behavioral support may contract further during
125 fine-tuning, as illustrated in Fig. 2, making underrepresented but potentially good actions increasingly
126 difficult to recover. Although prolonged fine-tuning can alleviate this effect (Liu et al., 2025), a
127 complementary and often cheaper solution is to prepare a **behaviorally flexible** policy prior that
128 already preserves diverse modes beyond well-supported actions in the offline dataset.

129 **Test-time condition change.** Standard offline RL typically assumes that the deployment environment
130 matches the one that generated the offline dataset (i.e., \mathcal{M}^*). In realistic open-ended domains,
131 dynamics, rewards, or task structure may shift after offline data collection, so actions that were optimal
132 under the offline dataset may no longer remain optimal at test time. Prior work has shown that offline-
133 learned policies can become brittle under such mismatches, including settings where variation is
134 specified through environment parameters (Liang et al., 2023), external data sources (Lyu et al., 2024),
135 natural-language instructions (Karthikeyan & Pant, 2025), or left unspecified (Mediratta et al., 2024).

136 4 A Bayesian Perspective for Adaptive Policy Priors

137 A natural question is whether adaptive policy priors admit a principled computational interpretation.
138 One promising direction arises from Bayesian principles in offline RL (Ghosh et al., 2022), building on
139 Bayes-adaptive MDPs (Duff, 2002), whose practical potential has recently been demonstrated (Chen
140 et al., 2021b; Choi et al., 2024; Ni et al., 2025). Rather than collapsing uncertainty conservatively or
141 treating unseen actions optimistically, Bayesian offline RL maintains multiple plausible environment
142 hypotheses consistent with the offline dataset, as illustrated in Fig. 3.

143 **Bayesian formulation of offline RL.** In Bayesian model-based offline RL, the agent maintains
144 epistemic uncertainty over environments consistent with the offline dataset \mathcal{D} . A prior distribution
145 $\Pr(\mathcal{M})$ over MDPs induces a posterior $\Pr(\mathcal{M} \mid \mathcal{D})$ after observing the dataset, where each
146 MDP \mathcal{M} is specified by environment parameters (ρ, P, R) . The offline objective becomes
147 $\max_{\pi} \mathbb{E}_{\mathcal{M} \sim \Pr(\mathcal{M} \mid \mathcal{D})} [J(\pi; \mathcal{M})]$. Because each sampled MDP \mathcal{M} is unknown to the agent, decision
148 making becomes a partially observable problem over environment hypotheses, often called
149 an *epistemic POMDP* (Ghosh et al., 2021, 2022). The optimal policy is therefore naturally
150 history-dependent: it maintains Bellman consistency under the posterior over plausible MDPs. By
151 contrast, model-free offline RL may assign unreliable values to out-of-support actions without an
152 explicit mechanism for Bellman consistency.

153 **Connection to adaptive policy priors.** This formulation directly explains the ingredients of adaptive
154 policy priors. Online history h_t provides *memory* about the environment; uncertain actions induce
155 *exploration*; and newly observed evidence supports *self-correction*. Formally, online interaction
156 further updates the environment posterior: $\Pr(\mathcal{M} \mid \mathcal{D}, h_t) \propto \Pr(\mathcal{M} \mid \mathcal{D}) \Pr(h_t \mid \mathcal{M}, \mathcal{D})$, where
157 the offline posterior $\Pr(\mathcal{M} \mid \mathcal{D})$ serves as the effective *prior* for online interaction. Adaptability is
158 thus a direct consequence of preserving epistemic uncertainty during offline optimization. Moreover,
159 test-time condition shift is naturally handled by reweighting the posterior toward the test-time MDP
160 \mathcal{M}' , provided that \mathcal{M}' falls within the support of $\Pr(\mathcal{M} \mid \mathcal{D})$.

161 5 Conclusion

162 We argue that offline RL should prioritize learning *adaptive policy priors* that preserve memory,
163 exploration, and self-correction for subsequent interaction, rather than static policies optimized
164 for immediate deployment. Bayesian offline RL provides one principled direction, though many
165 algorithmic and theoretical questions remain open (see Sec. B). Alternative views are discussed in
166 Sec. C and connections to related RL areas are in Sec. D. More broadly, offline RL should move
167 beyond a supervised-learning view shaped by data scaling laws (Kaplan et al., 2020), toward a
168 continual-learning view aligned with the emerging era of experience (Silver & Sutton, 2025).

References

- 169
170 Zaheer Abbas, Rosie Zhao, Joseph Modayil, Adam White, and Marlos C Machado. Loss of plasticity in continual
171 deep reinforcement learning. In *Conference on Lifelong Learning Agents*, 2023. 13
- 172 Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforce-
173 ment learning. In *International Conference on Machine Learning*, pp. 104–114. PMLR, 2020. 3
- 174 Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B. Tenenbaum, Tommi S. Jaakkola, and Pulkit Agrawal. Is
175 conditional generative modeling all you need for decision making? In *International Conference on Learning
176 Representations*, 2023. 10
- 177 Arthur Argenson and Gabriel Dulac-Arnold. Model-based offline planning. In *International Conference on
178 Learning Representations*, 2020. 10, 12
- 179 Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning with
180 offline data. In *International Conference on Machine Learning*, pp. 1577–1594. PMLR, 2023. 12
- 181 Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon Whiteson.
182 A tutorial on meta-reinforcement learning. *Foundations and Trends in Machine Learning*, 18(2-3):224–384,
183 2025. 13
- 184 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
185 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
186 *Advances in Neural Information Processing Systems*, 2020. 1, 10
- 187 Howard Chen, Noam Razin, Karthik Narasimhan, and Danqi Chen. Retaining by doing: The role of on-policy
188 data in mitigating forgetting. *arXiv preprint arXiv:2510.18874*, 2025. 12
- 189 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind
190 Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In
191 *Advances in Neural Information Processing Systems*, 2021a. 11, 12
- 192 Xiong-Hui Chen, Yang Yu, Qingyang Li, Fan-Ming Luo, Zhiwei Qin, Wenjie Shang, and Jieping Ye. Offline
193 model-based adaptable policy learning. *Advances in Neural Information Processing Systems*, 34:8432–8443,
194 2021b. 3, 4
- 195 Ching-An Cheng, Andrey Kolobov, Dipendra Misra, Allen Nie, and Adith Swaminathan. LLF-bench: Benchmark
196 for interactive learning from language feedback. In *ICLR 2024 Workshop on Large Language Model (LLM)
197 Agents*, 2024. 10
- 198 Egor Cherepanov, Nikita Kachaev, Artem Zhohus, Alexey Kovalev, and Aleksandr Panov. Unraveling the
199 complexity of memory in RL agents: an approach for classification and evaluation. In *International Conference
200 on Learning Representations*, 2026. 11
- 201 Yunseon Choi, Li Zhao, Chuheng Zhang, Lei Song, Jiang Bian, and Kee-Eung Kim. Diversification of adaptive
202 policy for effective offline reinforcement learning. In *International Joint Conference on Artificial Intelligence*,
203 pp. 3863–3871, 2024. 4
- 204 Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. RL2: Fast reinforcement
205 learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016. 11, 13
- 206 Michael O’Gordon Duff. *Optimal Learning: Computational procedures for Bayes-adaptive Markov decision
207 processes*. University of Massachusetts Amherst, 2002. 4
- 208 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep
209 networks. In *International Conference on Machine Learning*, 2017. 10
- 210 Kevin Frans, Seohong Park, Pieter Abbeel, and Sergey Levine. Unsupervised zero-shot reinforcement learning
211 via functional reward encodings. In *International Conference on Machine Learning*, 2024. 13
- 212 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven
213 reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020. 11, 12
- 214 Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *Advances in
215 neural information processing systems*, 34:20132–20145, 2021. 1
- 216 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration.
217 In *International conference on machine learning*, pp. 2052–2062. PMLR, 2019. 10, 11, 12

- 218 Dibya Ghosh, Jad Rahme, Aviral Kumar, Amy Zhang, Ryan P. Adams, and Sergey Levine. Why generalization
219 in RL is difficult: Epistemic pomdps and implicit partial observability. In *Advances in Neural Information*
220 *Processing Systems*, 2021. 2, 4
- 221 Dibya Ghosh, Anurag Ajay, Pulkit Agrawal, and Sergey Levine. Offline rl policies should be trained to be
222 adaptive. In *International Conference on Machine Learning*, pp. 7513–7530. PMLR, 2022. 2, 3, 4, 11
- 223 Jake Grigsby, Linxi Fan, and Yuke Zhu. Amago: Scalable in-context reinforcement learning for adaptive agents.
224 In *International Conference on Learning Representations*, 2024. 11
- 225 Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Thomas Paine, Sergio Gómez, Konrad Zolna, Rishabh
226 Agarwal, Josh S Merel, Daniel J Mankowitz, Cosmin Paduraru, et al. Rl unplugged: A suite of benchmarks
227 for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:7248–7259, 2020.
228 11
- 229 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang,
230 Shirong Ma, Xiao Bi, et al. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning.
231 *arXiv preprint arXiv:2501.12948*, 2025a. 1, 11
- 232 Yanjiang Guo, Jianke Zhang, Xiaoyu Chen, Xiang Ji, Yen-Jen Wang, Yucheng Hu, and Jianyu Chen. Improving
233 vision-language-action model with online reinforcement learning. In *International Conference on Robotics*
234 *and Automation*, 2025b. 1
- 235 Raia Hadsell, Dushyant Rao, Andrei A Rusu, and Razvan Pascanu. Embracing change: Continual learning in
236 deep neural networks. *Trends in Cognitive Sciences*, 24(12):1028–1040, 2020. 1
- 237 Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv preprint*
238 *arXiv:1502.02259*, 2015. 13
- 239 Wenchong He, Zhe Jiang, Tingsong Xiao, Zelin Xu, and Yukun Li. A survey on uncertainty quantification
240 methods for deep learning. *ACM Computing Surveys*, 58(7):1–35, 2026. 11
- 241 Hao Hu, Yiqin Yang, Jianing Ye, Chengjie Wu, Ziqing Mai, Yujing Hu, Tangjie Lv, Changjie Fan, Qianchuan
242 Zhao, and Chongjie Zhang. Bayesian design principles for offline-to-online reinforcement learning. In
243 *International Conference on Machine Learning*, pp. 19491–19515. PMLR, 2024. 12
- 244 Edward Hughes, Michael D Dennis, Jack Parker-Holder, Feryal Behbahani, Aditi Mavalankar, Yuge Shi, Tom
245 Schaul, and Tim Rocktäschel. Position: Open-endedness is essential for artificial superhuman intelligence. In
246 *International Conference on Machine Learning*, 2024. 2
- 247 Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior
248 synthesis. In *International Conference on Machine Learning*, pp. 9902–9915. PMLR, 2022. 3, 10, 12
- 249 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In *International*
250 *Conference on Machine Learning*, pp. 5084–5096. PMLR, 2021. 3
- 251 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray,
252 Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint*
253 *arXiv:2001.08361*, 2020. 4
- 254 Akash Karthikeyan and Yash Vardhan Pant. Genplan: Generative sequence models as adaptive planners. In
255 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025. 4, 11
- 256 Timo Klein, Lukas Mikloutz, Kevin Sidak, Claudia Plant, and Sebastian Tschiatschek. Plasticity loss in deep
257 reinforcement learning: A survey. *arXiv preprint arXiv:2411.04832*, 2024. 13
- 258 Martin Klissarov, Jonathan Cook, Diego Antognini, Hao Sun, Jingling Li, Natasha Jaques, Claudiu Musat,
259 and Edward Grefenstette. Improving interactive in-context learning from natural language feedback. *arXiv*
260 *preprint arXiv:2602.16066*, 2026. 10, 13
- 261 Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via
262 bootstrapping error reduction. *Advances in Neural Information Processing Systems*, 32, 2019. 1, 11, 12
- 263 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement
264 learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020. 1
- 265 Aviral Kumar, Rishabh Agarwal, Xinyang Geng, George Tucker, and Sergey Levine. Offline q-learning on
266 diverse multi-task data both scales and generalizes. In *International Conference on Learning Representations*,
267 2023. 1

- 268 Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli, Shariq
269 Iqbal, Colton Bishop, Rebecca Roelofs, Lei M Zhang, Kay McKinney, Disha Shrivastava, Cosmin Paduraru,
270 George Tucker, Doina Precup, Feryal Behbahani, and Aleksandra Faust. Training language models to
271 self-correct via reinforcement learning. In *International Conference on Learning Representations*, 2025. 3
- 272 Gaspard Lambrechts, Adrien Bolland, and Damien Ernst. Informed POMDP: Leveraging additional information
273 in model-based RL. In *Reinforcement Learning Conference*, 2024. 11
- 274 Kuang-Huei Lee, Ofir Nachum, Mengjiao Yang, Lisa Lee, Daniel Freeman, Sergio Guadarrama, Ian Fischer,
275 Winnie Xu, Eric Jang, Henryk Michalewski, and Igor Mordatch. Multi-game decision transformers. In
276 *Advances in Neural Information Processing Systems*, 2022a. 1
- 277 Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online reinforcement
278 learning via balanced replay and pessimistic q-ensemble. In *Conference on Robot Learning*, pp. 1702–1712.
279 PMLR, 2022b. 11
- 280 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review,
281 and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020. 1, 2, 3
- 282 Chenhao Li, Andreas Krause, and Marco Hutter. Uncertainty-aware robotic world model makes offline model-
283 based reinforcement learning work on real robots. *arXiv preprint arXiv:2504.16680*, 2025a. 11
- 284 Lu Li, Tianwei Ni, Yihao Sun, and Pierre-Luc Bacon. The three regimes of offline-to-online reinforcement
285 learning. *arXiv preprint arXiv:2510.01460*, 2025b. 13
- 286 Zhixuan Liang, Yao Mu, Mingyu Ding, Fei Ni, Masayoshi Tomizuka, and Ping Luo. AdaptDiffuser: diffusion
287 models as adaptive self-evolving planners. In *International Conference on Machine Learning*, 2023. 4, 11
- 288 Haoxin Lin, Siyuan Xiao, Yi-Chen Li, Zhilong Zhang, Yihao Sun, Chengxing Jia, and Yang Yu. ADM-v2:
289 Pursuing full-horizon roll-out in dynamics models for offline policy learning and evaluation. In *International
290 Conference on Learning Representations*, 2026. 11
- 291 Jinxin Liu, Hongyin Zhang, Zifeng Zhuang, Yachen Kang, Donglin Wang, and Bin Wang. Design from policies:
292 Conservative test-time adaptation for offline policy optimization. *Advances in Neural Information Processing
293 Systems*, 2023. 3
- 294 Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. ProRL:
295 Prolonged reinforcement learning expands reasoning boundaries in large language models. In *Conference on
296 Neural Information Processing Systems*, 2025. 4
- 297 Zuxin Liu, Zijian Guo, Haohong Lin, Yihang Yao, Jiacheng Zhu, Zhepeng Cen, Hanjiang Hu, Wenhao Yu,
298 Tingnan Zhang, Jie Tan, and Ding Zhao. Datasets and benchmarks for offline safe reinforcement learning.
299 *Journal of Data-centric Machine Learning Research*, 2024. 11
- 300 Cong Lu, Philip J. Ball, Tim G. J. Rudner, Jack Parker-Holder, Michael A Osborne, and Yee Whye Teh.
301 Challenges and opportunities in offline reinforcement learning from visual observations. *Transactions on
302 Machine Learning Research*, 2023. 11
- 303 Fan-Ming Luo, Zuolin Tu, Zefang Huang, and Yang Yu. Efficient recurrent off-policy rl requires a context-
304 encoder-specific learning rate. *Advances in Neural Information Processing Systems*, 37:48484–48518, 2024.
305 11
- 306 Yicheng Luo, Jackie Kay, Edward Grefenstette, and Marc Peter Deisenroth. Finetuning from offline rein-
307 forcement learning: Challenges, trade-offs and practical solutions. *arXiv preprint arXiv:2303.17396*, 2023.
308 4
- 309 Jiafei Lyu, Kang Xu, Jiacheng Xu, Jing-Wen Yang, Zongzhang Zhang, Chenjia Bai, Zongqing Lu, Xiu Li, et al.
310 Odrl: A benchmark for off-dynamics reinforcement learning. *Advances in Neural Information Processing
311 Systems*, 2024. 4, 11
- 312 Ishita Mediratta, Qingfei You, Minqi Jiang, and Roberta Raileanu. The generalization gap in offline reinforcement
313 learning. In *International Conference on Learning Representations*, 2024. 4, 11
- 314 Amir Moeini, Jiuqi Wang, Jacob Beck, Ethan Blaser, Shimon Whiteson, Rohan Chandra, and Shangdong Zhang.
315 A survey of in-context reinforcement learning. *arXiv preprint arXiv:2502.07978*, 2025. 10
- 316 Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online reinforcement
317 learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020. 4, 11, 13

- 318 Tianwei Ni, Michel Ma, Benjamin Eysenbach, and Pierre-Luc Bacon. When do transformers shine in RL?
319 decoupling memory from credit assignment. In *Advances in Neural Information Processing Systems*, 2023.
320 [11](#)
- 321 Tianwei Ni, Benjamin Eysenbach, Erfan SeyedSalehi, Michel Ma, Clement Gehring, Aditya Mahajan, and Pierre-
322 Luc Bacon. Bridging state and history representations: Understanding self-predictive rl. In *International*
323 *Conference on Learning Representations*, 2024. [11](#)
- 324 Tianwei Ni, Esther Derman, Vineet Jain, Vincent Taboga, Siamak Ravanbakhsh, and Pierre-Luc Bacon. Long-
325 horizon model-based offline reinforcement learning without conservatism. *arXiv preprint arXiv:2512.04341*,
326 2025. [2](#), [4](#), [11](#), [12](#)
- 327 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang,
328 Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,
329 Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training
330 language models to follow instructions with human feedback. In *Advances in Neural Information Processing*
331 *Systems*, 2022. [12](#)
- 332 Charles Packer, Katelyn Gao, Jernej Kos, Philipp Krähenbühl, Vladlen Koltun, and Dawn Song. Assessing
333 generalization in deep reinforcement learning. *arXiv preprint arXiv:1810.12282*, 2018. [10](#)
- 334 Seohong Park, Kevin Frans, Sergey Levine, and Aviral Kumar. Is value learning really the main bottleneck in
335 offline rl? *Advances in Neural Information Processing Systems*, 2024a. [3](#), [10](#)
- 336 Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation policies with hilbert representations. In
337 *International Conference on Machine Learning*, 2024b. [13](#)
- 338 Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking offline
339 goal-conditioned rl. In *International Conference on Learning Representations*, 2025a. [11](#)
- 340 Seohong Park, Kevin Frans, Deepinder Mann, Benjamin Eysenbach, Aviral Kumar, and Sergey Levine. Horizon
341 reduction makes RL scalable. In *Conference on Neural Information Processing Systems*, 2025b. [1](#)
- 342 Rong-Jun Qin, Xingyuan Zhang, Songyi Gao, Xiong-Hui Chen, Zewen Li, Weinan Zhang, and Yang Yu. Neorl:
343 A near real-world benchmark for offline reinforcement learning. *Advances in Neural Information Processing*
344 *Systems*, 35:24753–24765, 2022. [11](#)
- 345 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov, Gabriel Barth-
346 maron, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi,
347 Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando
348 de Freitas. A generalist agent. *Transactions on Machine Learning Research*, 2022. Featured Certification,
349 Outstanding Certification. [1](#)
- 350 Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *Proceedings of the thirteenth in-*
351 *ternational conference on artificial intelligence and statistics*, pp. 661–668. JMLR Workshop and Conference
352 Proceedings, 2010. [3](#)
- 353 Flore Sentenac, Ilbin Lee, and Csaba Szepesvari. Balancing optimism and pessimism in offline-to-online
354 learning. *arXiv preprint arXiv:2502.08259*, 2025. [11](#)
- 355 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language
356 agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 2023. [10](#)
- 357 David Silver and Richard S Sutton. Welcome to the era of experience. *preprint*, 2025. [1](#), [4](#)
- 358 Anya Sims, Cong Lu, Jakob N Foerster, and Yee W Teh. The edge-of-reach problem in offline model-based
359 reinforcement learning. *Advances in Neural Information Processing Systems*, 37:63029–63056, 2024. [11](#)
- 360 Amit Sinha and Aditya Mahajan. Agent-state based policies in pomdps: Beyond belief-state mdps. In *Conference*
361 *on Decision and Control*. IEEE, 2024. [11](#)
- 362 Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally
363 can be more effective than scaling parameters for reasoning. In *International Conference on Learning*
364 *Representations*, 2025. [1](#), [12](#)
- 365 Yuda Song, Yifei Zhou, Ayush Sekhari, Drew Bagnell, Akshay Krishnamurthy, and Wen Sun. Hybrid RL: Using
366 both offline and online data can make RL efficient. In *International Conference on Learning Representations*,
367 2023. [12](#)

- 368 Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with
369 self-supervision for generalization under distribution shifts. In *International Conference on Machine Learning*,
370 2020. 10
- 371 Phillip Swazinna, Steffen Udluft, and Thomas Runkler. User-interactive offline reinforcement learning. In
372 *International Conference on Learning Representations*, 2023. 12
- 373 Ahmed Touati, Jérémy Rapin, and Yann Ollivier. Does zero-shot reinforcement learning exist? In *International*
374 *Conference on Learning Representations*, 2023. 12, 13
- 375 Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under partial coverage.
376 In *International Conference on Learning Representations*, 2022. 3
- 377 Andrew Wagenmaker, Perry Dong, Raymond Tsao, Chelsea Finn, and Sergey Levine. Posterior behavioral
378 cloning: Pretraining bc policies for efficient rl finetuning. *arXiv preprint arXiv:2512.16911*, 2025. 12
- 379 Jane Wang, Zeb Kurth-Nelson, Hubert Soyer, Joel Z. Leibo, Dhruva Tirumala, Rémi Munos, Charles Blundell,
380 Dharshan Kumaran, and Matt M. Botvinick. Learning to reinforcement learn. In *Annual Meeting of the*
381 *Cognitive Science Society*, 2017. 13
- 382 Yifei Wang, Yuyang Wu, Zeming Wei, Stefanie Jegelka, and Yisen Wang. A theoretical understanding of
383 self-correction through in-context alignment. *Advances in Neural Information Processing Systems*, 2024. 3
- 384 Hui Wei, Zihao Zhang, Shenghua He, Tian Xia, Shijia Pan, and Fei Liu. Plangenllms: A modern survey of llm
385 planning capabilities. In *Annual Meeting of the Association for Computational Linguistics*, pp. 19497–19521,
386 2025. 10, 12
- 387 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al.
388 Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information*
389 *Processing Systems*, 35:24824–24837, 2022. 1
- 390 Manuel Wendl, Yarden As, Manish Prajapat, Anton Pollak, Stelian Coros, and Andreas Krause. Safe exploration
391 via policy priors. In *International Conference on Learning Representations*, 2026. 12
- 392 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as
393 implicit bayesian inference. In *International Conference on Learning Representations*, 2022. 10
- 394 Shoukai Xu, Mingkui Tan, Liu Liu, Zhong Zhang, Peilin Zhao, et al. Test-time adapted reinforcement learning
395 with action entropy regularization. In *International Conference on Machine Learning*, 2025. 3, 10
- 396 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree
397 of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information*
398 *Processing Systems*, 2023. 1
- 399 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu
400 Ma. Mopo: Model-based offline policy optimization. *Advances in Neural Information Processing Systems*,
401 33:14129–14142, 2020. 1, 11
- 402 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. Does
403 reinforcement learning really incentivize reasoning capacity in LLMs beyond the base model? In *Conference*
404 *on Neural Information Processing Systems*, 2025. 4, 10, 12
- 405 Dylan Zhang, Yufeng Xu, Haojin Wang, Qingzhi Chen, and Hao Peng. Good sft optimizes for sft, better sft
406 prepares for reinforcement learning. *arXiv preprint arXiv:2602.01058*, 2026a. 12
- 407 Shenao Zhang, Yaqing Wang, Yinxiao Liu, Tianqi Liu, Peter Grabowski, Eugene Ie, Zhaoran Wang, and Yunxuan
408 Li. Beyond markovian: Reflective exploration via bayes-adaptive RL for LLM reasoning. In *International*
409 *Conference on Learning Representations*, 2026b. 12
- 410 Kai Zhao, Jianye Hao, Yi Ma, Jinyi Liu, Yan Zheng, and Zhaopeng Meng. Enoto: Improving offline-to-online
411 reinforcement learning with q-ensembles. *arXiv preprint arXiv:2306.06871*, 2023. 4
- 412 Rosie Zhao, Alexandru Meterez, Sham M. Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach. Echo
413 chamber: RL post-training amplifies behaviors learned in pretraining. In *Conference on Language Modeling*,
414 2025. 4
- 415 Guangyao Zhou, Sivaramakrishnan Swaminathan, Rajkumar Vasudeva Raju, J Swaroop Guntupalli, Wolfgang
416 Lehrach, Joseph Ortiz, Antoine Dedieu, Miguel Lázaro-Gredilla, and Kevin Murphy. Diffusion model
417 predictive control. *Transactions on Machine Learning Research*, 2025. 10, 11

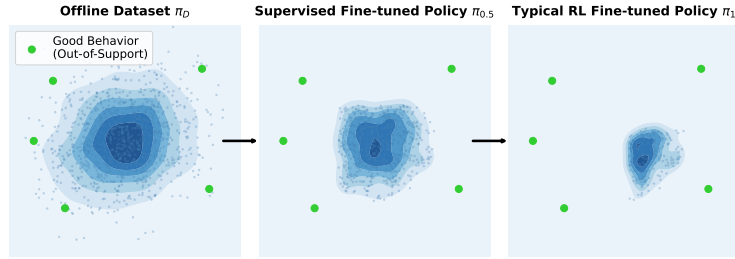


Figure 2: Behavioral support may shrink across offline (here, SFT) and online fine-tuning stages, making underrepresented but possibly good actions increasingly difficult to recover (Yue et al., 2025).

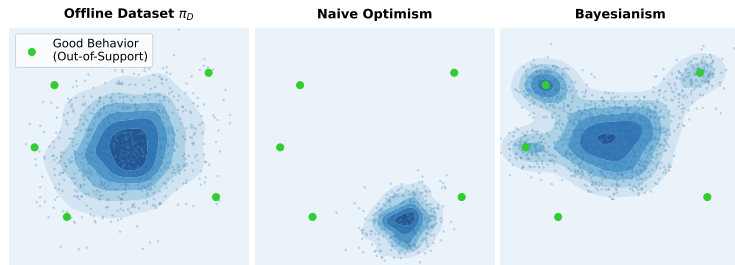


Figure 3: Bayesian offline RL treats unseen actions as uncertain rather than inherently bad or reliably good. In contrast, naive optimism in model-free off-policy RL may assign unreliable values to out-of-support actions due to extrapolation error (Fujimoto et al., 2019).

419 A Forms of Adaptation in Details

420 **Test-time in-context learning.** The simplest form of adaptation occurs when decision making is
 421 history-dependent, so that behavior changes directly through interaction history without modifying
 422 policy parameters or performing additional optimization. This mechanism is commonly referred to
 423 as *in-context learning* (Brown et al., 2020; Moeini et al., 2025), and has often been interpreted as
 424 implicit Bayesian inference over latent tasks (Xie et al., 2022). It differs from *in-weight learning*,
 425 where adaptation occurs through parameter updates. A special case of in-context learning is *self-*
 426 *improvement* (Shinn et al., 2023), where explicit reward signals may be unavailable at test time,
 427 yet the agent can still improve through informative observations that reduce uncertainty about the
 428 environment. Self-improvement includes settings where the uncertainty lies in transition dynamics
 429 rather than reward function (Packer et al., 2018), as well as settings where feedback is conveyed
 430 through language observations (Cheng et al., 2024; Klissarov et al., 2026).

431 **Test-time planning.** In-context learning relies primarily on the generalization ability of the offline-
 432 learned policy, which may be insufficient when the true environment \mathcal{M}^* differs substantially from
 433 what is specified by \mathcal{D} . Test-time planning addresses this limitation by improving policy decisions
 434 with online search over future trajectories, often using learned world models (Argenson & Dulac-
 435 Arnold, 2020; Zhou et al., 2025; Wei et al., 2025). Recent planning approaches further enable
 436 trajectory optimization through probabilistic inference (e.g., diffusion-based denoising) without
 437 requiring an explicit world model (Janner et al., 2022; Ajay et al., 2023).

438 **Test-time training.** Adaptation may also occur through parameter updates using a limited amount of
 439 online data collected at test time (Finn et al., 2017; Sun et al., 2020). Unlike in-context learning or
 440 planning, test-time training modifies the policy itself, often through maximizing Q-value on test-time
 441 states without reward signals (Park et al., 2024a; Xu et al., 2025). Although this lies outside the
 442 classical offline RL formulation, it is relevant under AORL because the offline-learned policy is
 443 viewed as a prior that should support efficient improvement from limited online experience.

444 **Online RL fine-tuning.** A more extended form of adaptation is online RL fine-tuning, where the
 445 offline-learned policy serves as an initialization for continued RL with substantial online interaction.

446 This includes offline-to-online RL (Nair et al., 2020; Lee et al., 2022b) and RL post-training of
447 foundation models (Guo et al., 2025a). Compared with test-time training, the focus is not rapid local
448 adaptation, but sustained policy improvement through accumulating experience.

449 B Open Challenges and Research Directions

450 **Benchmarks and evaluation for adaptability.** Existing offline RL benchmarks evaluate many
451 important settings, including high-dimensional observations (Gulcehre et al., 2020; Lu et al., 2023),
452 sparse rewards and stitching (Fu et al., 2020; Park et al., 2025a), stochasticity (Qin et al., 2022), and
453 safety (Liu et al., 2024). However, they rarely evaluate whether an offline-learned policy can *improve*
454 *through subsequent interaction*. Prior work on adaptability has so far focused either on bandit-like
455 settings (Ghosh et al., 2022; Ni et al., 2025) or illustrative task modifications (Liang et al., 2023;
456 Zhou et al., 2025; Karthikeyan & Pant, 2025). Two notable exceptions are the off-dynamics RL
457 benchmark (Lyu et al., 2024) and the generalization benchmark (Mediratta et al., 2024). In Lyu et al.
458 (2024), test-time condition changes are specified through externally provided datasets, whereas in
459 Mediratta et al. (2024), changes are conveyed either through text instructions in WebShop or remain
460 unspecified in Procgen. If AORL is to become a meaningful research direction, progress will require
461 more benchmarks that explicitly measure the capability for test-time adaptation.

462 On the offline side, dataset design should include at least two cases: *underrepresented but valuable*
463 *actions*, where useful behaviors appear only sparsely in the data, and *unseen but plausible actions*,
464 where beneficial behaviors are absent altogether and must be discovered through later interaction. On
465 the online side, evaluation should go beyond a single fixed test environment and instead consider *a*
466 *distribution of environments*: (1) environments fully aligned with the offline data, and (2) environments
467 with test-time condition changes under varying degrees of specification.

468 More broadly, AORL benchmarks should treat *multi-episode interaction* as the basic unit of training
469 and evaluation, akin to meta-RL setups (Duan et al., 2016), in order to capture the improvable capacity
470 of adaptive agents. Evaluation should measure not only final returns, but also *how* improvement
471 unfolds through interaction, including adaptation speed, robustness to early mistakes, recovery after
472 failed exploration, and sensitivity to test-time horizon length (Sentenac et al., 2025).

473 **Theory of adaptive offline RL.** Classical offline RL theory typically relies on coverage assumptions
474 over good actions and Markovian policies. Extending these guarantees to history-dependent policies
475 that can discover unseen but valuable actions through online interaction remains largely open. Key
476 questions include which class of offline datasets favors adaptive priors over static policies, when
477 in-context learning suffices versus when parameter updates are necessary, and how guarantees should
478 change under test-time condition shift.

479 **Overcoming value overestimation.** A major challenge in AORL is that classical *value overestimation*
480 re-emerges once explicit conservatism is relaxed (Fujimoto et al., 2019; Kumar et al., 2019; Sims
481 et al., 2024). Because adaptive priors must preserve uncertain actions for later exploration, avoiding
482 overestimation *without* collapsing back to conservative behavior becomes a central algorithmic
483 difficulty. Recent evidence from Bayesian offline RL suggests that sufficiently long-horizon rollouts
484 can mitigate this problem (Ni et al., 2025). This points to a broader direction: scaling long-horizon
485 planning in model-based offline RL, potentially informed by recent progress in long-horizon world
486 modeling (Li et al., 2025a; Lin et al., 2026).

487 **Scalable Bayesian inference.** A practical bottleneck of the Bayesian direction in Sec. 4 is that exact
488 posterior inference over MDPs is typically intractable. In practice, uncertainty is often approximated
489 through model ensembles, where disagreement serves as a proxy for epistemic uncertainty (Yu et al.,
490 2020). Improving the scalability of such approximations, while achieving more reliable uncertainty
491 quantification (He et al., 2026), remains an important direction for AORL.

492 **Memory-based RL.** Memory remains underexplored in offline RL, with the most notable progress
493 such as the Decision Transformer family (Chen et al., 2021a), which still operate under conservative
494 objectives. We believe that advances in understanding memory (Ni et al., 2023; Cherepanov et al.,
495 2026), learning history representations (Lambrechts et al., 2024; Ni et al., 2024; Sinha & Mahajan,
496 2024), and scaling memory-based RL (Grigsby et al., 2024; Luo et al., 2024) in online POMDPs
497 can help AORL by understanding how in-context learning emerges from offline-trained policies.

498 **Beyond Bayesian model-based direction.** Although we emphasize the Bayesian model-based
499 direction in Sec. 4 because it offers a principled interpretation of adaptive policy priors, simpler
500 alternatives may also achieve adaptive behavior. These include model-free approaches (Hu et al., 2024;
501 Wagenmaker et al., 2025) as well as non-Bayesian approaches (Touati et al., 2023). Understanding
502 when such methods recover similar adaptive mechanisms, and when they fundamentally differ from
503 model-based formulations, remains an important open direction.

504 **Foundation models for AORL.** World foundation models offer an intriguing answer to the *prior*
505 specification problem in Sec. 4: rather than constructing $\Pr(\mathcal{M})$ from scratch, one can specify the
506 MDP prior from a foundation model pretrained on broad dynamics data. Pretrained history-dependent
507 policies, including large language models and vision-language-action models, may further serve as
508 strong initial policies that efficiently explore beyond offline dataset support.

509 **AORL for foundation models.** Recent progress in foundation models already highlights the im-
510 portance of test-time adaptation and continual learning (Wei et al., 2025; Snell et al., 2025). Yet
511 offline post-training remains dominated by supervised fine-tuning (SFT) (Ouyang et al., 2022), which
512 may narrow behavioral support and hinder later test-time adaptation or online RL fine-tuning (Zhang
513 et al., 2026a). Bayesian directions may therefore offer a useful path for AORL in foundation models:
514 model-based formulations can preserve broader support through synthetic on-policy rollouts (Chen
515 et al., 2025), while Bayesian formulations can promote in-context adaptation (Zhang et al., 2026b).

516 C Alternative Views

517 **Conservatism as the first principle.** A dominant alternative view in offline RL is that conservatism
518 should remain the primary design principle, motivated mainly by a technical concern: offline value
519 estimation is prone to extrapolation error and value overestimation (Fujimoto et al., 2019; Kumar
520 et al., 2019). Conservative objectives therefore restrict policy improvement toward in-distribution
521 behavior to stabilize learning; in high-stakes deployment settings, this also aligns naturally with safety
522 requirements. Our position differs when an offline-learned policy is expected to continue improving
523 online. In this setting, uncertainty need not be fully eliminated offline, and excessive conservatism
524 may suppress behaviors needed for later adaptation. Recent evidence suggests that value overestima-
525 tion can also be controlled through alternative mechanisms such as long-horizon planning (Ni et al.,
526 2025), weakening the case for conservatism as the default first principle. Importantly, *adaptability*
527 *does not rule out moderate conservatism*: some degree of conservatism may still be necessary (Wendl
528 et al., 2026), and can even be adjusted at test time (Swazinna et al., 2023). What matters is preserving
529 enough memory and behavioral flexibility for the policy to revise decisions through interaction. The
530 key distinction is whether offline learning aims to produce a static, final policy or an improvable prior.

531 **Adaptation should be left to online RL, not offline RL.** Another alternative view is that offline
532 RL need not itself produce an adaptive prior; its role is simply to provide a stable initialization
533 for subsequent online RL, while adaptation is delegated to downstream fine-tuning. Under this
534 perspective, one may even blur the distinction between offline and online learning by directly training
535 online RL from scratch with offline data (Song et al., 2023; Ball et al., 2023). This is reasonable in
536 settings where online interaction is abundant, but incomplete when online experience is limited or
537 costly. If offline training collapses the policy too narrowly—for example through behavioral cloning
538 or strongly conservative objectives—online RL may inherit a restricted search space and struggle to
539 recover underrepresented but valuable behaviors (Yue et al., 2025). Finally, when offline learning
540 is intended to produce *reusable* priors rather than task-specific warm starts, its objective should
541 therefore extend beyond fast initialization to preserving behavioral flexibility for later adaptation.

542 **Adaptability already appears in many offline RL methods.** Another alternative view is that
543 offline RL already contains mechanisms associated with adaptation, making an explicit shift in
544 objective unnecessary. Common history-dependent policies with sequence modeling (Chen et al.,
545 2021a) use past context during action generation, yet typically remain constrained by patterns
546 learned from the offline dataset. Test-time planning methods (Argenson & Dulac-Arnold, 2020;
547 Janner et al., 2022) improve decisions through online search, but the candidate trajectories they
548 consider are usually still constrained by offline support. Similarly, *stitching* (Fu et al., 2020) enables
549 compositional generalization from offline data, but this remains a passive recombination of in-
550 distribution behaviors rather than adaptation driven by new online evidence. These examples suggest

551 that adaptive ingredients are increasingly present in offline RL, but common formulations still do not
552 explicitly preserve exploration and self-correction under later interaction.

553 **D Connections Beyond Offline RL**

554 Adaptive offline RL connects naturally to neighboring research areas that study adaptation and
555 generalization under different assumptions.

556 **Offline-to-online RL.** Although offline RL is defined by the absence of online data collection during
557 training, the learned policy is still evaluated through interaction with the true environment \mathcal{M}^* at test
558 (deployment) time. This creates a temporal separation between the offline stage, where parameters
559 are optimized from a fixed dataset, and the online stage, where decisions are generated under the
560 true environment. This distinction also separates offline RL (including AORL) from offline-to-online
561 RL (Nair et al., 2020), where the policy (persistently) updates its parameters in the online stage.

562 **Meta-RL.** The Bayesian perspective in [Sec. 4](#) is closely related to contextual MDPs (Hallak et al.,
563 2015) and meta-RL (Duan et al., 2016; Wang et al., 2017; Beck et al., 2025), where policies are
564 trained across a distribution of tasks so that adaptation emerges at test time. The key difference is that
565 in meta-RL the task distribution is pre-specified before training, whereas in Bayesian offline RL the
566 distribution over environments must be self-constructed from offline data. Nevertheless, both settings
567 give rise to similar adaptive behavior at test time, including online exploration and self-correction
568 under uncertainty.

569 **Plasticity in continual RL.** The ability to improve from subsequent experience is closely related
570 to *plasticity* in continual RL (Abbas et al., 2023; Klein et al., 2024) and offline-to-online RL (Li
571 et al., 2025b). AORL differs in that adaptation often occurs without parameter updates at deployment
572 time, also known as *in-context plasticity* (Klissarov et al., 2026). Still, the central concern is similar:
573 whether prior learning preserves the capacity to improve when new experience becomes available.

574 **Behavior foundation models.** A recent line of work trains behavior foundation models (BFMs)
575 from *reward-free* offline datasets, aiming to produce policies that generalize zero-shot to arbitrary
576 downstream reward functions (Touati et al., 2023; Park et al., 2024b; Frans et al., 2024). The key idea
577 is to learn representations that decouple dynamics from reward, so that a new task is solved by simple
578 computation on a few externally-provided reward-labeled samples without further training. Once the
579 task embedding is inferred, the policy is static and memoryless; dynamics are typically assumed fixed.

580 BFMs are therefore complementary to, rather than a substitute for, AORL. Their reward generaliz-
581 ability could be incorporated into adaptive policy priors to handle the *reward specification* problem at
582 test time. This also sharpens what “adaptation” means in AORL: unlike the passive task inference
583 of BFMs from externally provided data, AORL requires *active* generalization in which the agent
584 generates its own experience to resolve uncertainty.