
Beyond One-Size-Fits-All: Diagnosis-Driven Online Reinforcement Learning with Offline Priors

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Online reinforcement learning (RL) agents increasingly depend on knowledge
2 acquired offline to achieve practical efficiency. Originally studied in offline-to-
3 online RL, this paradigm now spans foundation model post-training and embodied
4 intelligence, with prior types expanding from offline datasets and pre-trained
5 policies to increasingly diverse knowledge sources such as multimodal foundation
6 models and generative world models. Offline priors have become central to how
7 deep RL is developed and deployed. However, this reliance introduces a challenge
8 that the prevailing benchmark-driven paradigm cannot resolve: because prior
9 validity varies across deployments and shifts during training, no single approach to
10 managing it is universally optimal, and benchmark rankings offer limited guidance
11 for real-world deployments. Rather than pursuing universal solutions, we argue that
12 the field should shift to diagnosis-driven tension management, in which deployment-
13 specific evidence guides how the learner relates to its priors throughout training,
14 enabling both flexible and adaptive deployment. We support this position with a
15 framework characterizing how priors reshape online optimization through three
16 functional roles, controlled experiments demonstrating help-or-hurt reversals, cross-
17 domain evidence from foundation model post-training to embodied intelligence,
18 and engagement with five substantive counterarguments.

19 1 Introduction

20 Online reinforcement learning (RL) agents increasingly rely on knowledge acquired offline to achieve
21 practical performance. In foundation model post-training, large-scale pre-trained models provide the
22 base capability that online RL refines for alignment, reasoning, and agentic applications [1, 2]. In
23 embodied intelligence, simulators, demonstration datasets, world models, and pre-trained policies
24 supply the prior knowledge that makes online learning on physical platforms feasible [3, 4, 5, 6].
25 Despite differences in domain, objective, and prior type, offline priors now provide the foundation
26 from which agents learn through online experience [7, 8, 9, 10].

27 Despite these benefits, how to use offline priors remains a persistent challenge: the same reliance
28 decision that helps in one setting often hurts in another. In offline-to-online RL, whether to preserve the
29 pre-trained policy, retain the offline dataset, or maintain conservative value estimates each produces
30 different outcomes depending on the quality of the offline sources and their relationship to the
31 deployment task [11, 8, 7, 12]. In foundation model post-training, the role of the reference constraint
32 varies with the fidelity of the reward signal: when rewards are verifiable, strong regularization
33 restricts the discovery of novel strategies [13], while when rewards come from learned preference
34 models, the same regularization is essential for preventing overoptimization of proxy scores [1, 14].
35 In robotics, whether to train in simulation and transfer or to learn directly on physical hardware
36 produces different outcomes depending on the fidelity of the simulator and the complexity of the

Online RL from Scratch



Online RL with Offline Priors

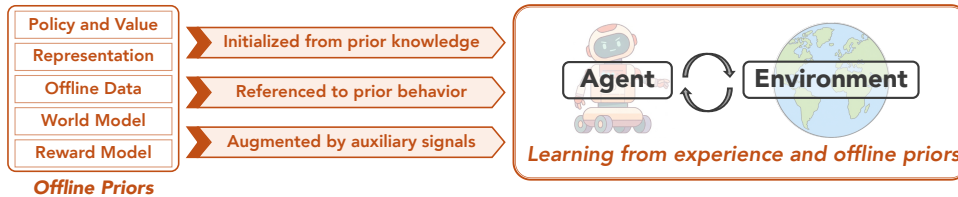


Figure 1: Online RL from scratch (top) versus with offline priors (bottom). The agent-environment interaction loop is identical in both cases. Offline priors add knowledge sources that accelerate learning but also introduce structural tensions analyzed in [Sec. 3](#).

37 contact dynamics [15, 6, 16]. These inconsistencies recur across communities, domains, and prior
38 types, suggesting a shared structural property rather than isolated engineering issues.

39 We argue that these inconsistencies reflect an inherent conflict: offline priors carry knowledge that is
40 inevitably bounded, while online RL exists precisely to push beyond those bounds. • Specifically,
41 whether the prior takes the form of a pre-trained policy, a simulator, or an offline dataset, it encodes
42 knowledge acquired under conditions that may differ from those of deployment. As a result, the
43 extent to which this knowledge remains useful is uncertain and can only be revealed through online
44 interaction. We call this the *bounded commitment* of offline priors: valuable knowledge with an
45 inherently uncertain scope of validity. • On the other hand, online RL drives the agent toward
46 optimal performance through interaction with the deployment environment. This requires the agent
47 to eventually surpass what the prior covers, while relying on its knowledge to get there efficiently.
48 As long as the agent operates within the scope of the prior’s validity, reliance is purely beneficial.
49 Once learning pushes beyond that scope, a genuine tension emerges: stronger reliance constrains
50 adaptation, while freer adaptation risks discarding knowledge that still helps elsewhere. Crucially,
51 the agent cannot know where this boundary lies, and the boundary itself shifts as online experience
52 accumulates. Unlike a static trade-off that can be settled at design time, this tension persists and
53 evolves throughout learning. Hence, tension management is a central challenge of this paradigm.

54 One consequence is that tension management has no universal optimum: the right reliance on each
55 prior depends on the deployment and shifts as learning progresses. Despite this, the field continues to
56 evaluate progress by comparing methods on fixed benchmarks, implicitly assuming that the resulting
57 rankings reflect universal truths rather than deployment-specific matches. The field accumulates
58 condition-specific performance rankings rather than transferable understanding. The way forward
59 requires a shift in research perspective. Rather than asking which method is best, the field should ask
60 what determines when each design choice helps or hurts.

We argue that the field of online RL with offline priors should move beyond one-size-fits-all methods toward diagnosis-driven tension management. Since no universal optimum exists for how agents should rely on their priors, effective deployment requires diagnostic infrastructure that can assess the prior-deployment match and monitor how tensions evolve during learning, enabling both flexible and adaptive deployment.

61

62 We develop this position as follows. [Sec. 2](#) defines the paradigm and introduces a taxonomy of
63 offline prior types. [Sec. 3](#) analyzes why offline priors sharpen the core tensions of online RL and
64 introduces the concept of bounded commitments. [Sec. 4](#) presents evidence that tension management
65 has no universal optimum and argues for a shift from benchmark-driven to diagnosis-driven tension
66 management. [Sec. 5](#) engages with several potential objections and counterarguments to our core
67 position. [Sec. 6](#) concludes with research opportunities and a broader perspective.

68 2 The Paradigm of Online RL with Offline Priors

69 Instead of *learning from scratch*¹, online RL agents increasingly learn not only from their own online
 70 interaction with the environment but also from knowledge acquired offline: setting initial parameters,
 71 constraining how far the learner may deviate from prior behavior, or supplying supplementary data or
 72 predictions alongside real interaction. Though studied under different names across offline-to-online
 73 RL [18, 19], LLM post-training [1, 2], sim-to-real transfer [20, 15], model-based RL [4, 21], and
 74 vision-language-action model fine-tuning [22, 23], these mechanisms define a common paradigm
 75 that we call *online RL with offline priors*. As in any RL system, the interaction loop has two sides:
 76 an agent that selects actions and an environment that produces states and rewards. Offline knowledge
 77 can concern either side, and we organize the resulting priors accordingly.

Table 1: Taxonomy of offline priors in online RL, classified by *what knowledge they encode* and *how they function* during online learning. Agent-side priors encode knowledge about how to act or evaluate; environment-side priors encode knowledge about how the world behaves.

Knowledge Side	Functional Role	Representative Prior Type		
Agent	Initialization	Policy π_0	Value Q_0 or V_0	Representation ϕ_0
	Reference	Policy π_{ref}	Value Q_{ref} or V_{ref}	[†] Offline Data \mathcal{D}
Environment	Auxiliary	* World Model \hat{M}	Reward Model \hat{R}	[†] Offline Data \mathcal{D}

[†] Offline data \mathcal{D} records how the environment transitions and how the collecting policy acts; its functional role depends on which aspect the online algorithm extracts [7, 24].

* World model \hat{M} spans a broad spectrum: compact latent dynamics models [4, 21], engineered simulators used to supplement real-world interaction [20, 15], and world foundation models [25, 26, 27, 28].

78 Offline priors can shape online optimization through exactly three channels: by determining where
 79 optimization begins, by constraining what objective it pursues, or by providing experience beyond
 80 direct interaction. • Initialization priors set the starting point for online learning: an offline-trained
 81 policy in offline-to-online RL [19], a supervised fine-tuned model in LLM post-training [1], or a
 82 simulation-trained controller in sim-to-real transfer [3] each play this role, providing initial com-
 83 petence and reducing the exploration burden that dominates learning from scratch. • Reference
 84 priors modify the learning objective by anchoring updates to prior behavior: the KL penalty to a
 85 reference policy in RLHF [1], the conservative value penalty in Cal-QL [19], and behavioral cloning
 86 regularization [24, 29] are different mechanisms serving the same structural function. In each case,
 87 the prior defines a trust region that the online learner is penalized from leaving. • Auxiliary priors
 88 provide additional information outside the online loop, whether through retained offline data [7],
 89 model-based rollouts [4], or learned reward signals [1]. Unlike reference priors, these do not alter
 90 what the optimizer aims to achieve, but expand the evidence it can draw on.

91 Offline data occupies a unique position in this taxonomy. It may originate from behavior-policy
 92 rollouts, expert demonstrations, human preference comparisons, or in-the-wild recordings [30, 5],
 93 and it inherently records both how the environment transitions and how the collecting policy acts.
 94 Which aspect the online algorithm extracts determines the functional role: replay for value estimation
 95 treats data as auxiliary information [7], while behavioral regularization treats the same data as a
 96 reference [24, 31]. This dependence on algorithmic use extends beyond data: a pessimistic value
 97 function always initializes the learner [19, 8], but additionally serves as an ongoing reference when
 98 its conservative penalty is maintained during fine-tuning [32]. In general, functional role is not an
 99 inherent property of any prior; it is determined by the algorithm that deploys it.

100 These priors fundamentally reshape the online learning process. The next section examines why these
 101 changes, despite their well-documented benefits, introduce structural tensions into online learning.

102 3 Why Offline Priors Sharpen Tensions

103 In principle, two fundamental tensions govern the online RL learning process: exploration must be
 104 balanced against exploitation, and the plasticity to incorporate new experience must be balanced
 105 against the stability of what has already been learned. However, in the challenging tasks that

¹Also referred to as *tabula rasa* RL [17].

106 organize mainstream deep RL research the balance tilts sharply to one side. Similarly, sparse rewards,
 107 high-dimensional action spaces, and long horizons make useful discoveries so rare that the need
 108 for exploration overwhelms any concern about premature exploitation [33, 34]. Bootstrapping
 109 from the agent’s own shifting value estimates induces optimization pathologies that progressively
 110 degrade the network’s ability to incorporate new experience [35, 36, 37, 38]. In these settings,
 111 failure does not stems from unbalanced tensions, but typically from one-side insufficiency: agents
 112 explore too little to discover useful behavior and remain too rigid to learn effectively from new
 113 evidence. The research priorities of the past decade reflect this asymmetry: exploration methods
 114 overwhelmingly aim to increase coverage [39, 34], and plasticity interventions overwhelmingly aim
 115 to restore adaptability [40, 41, 42].

116 **How Priors Sharpen the Tensions.** Over the past several years, diverse research communities
 117 have effectively addressed these bottlenecks by equipping agents with offline priors before online
 118 interaction begins. Pre-trained policies and value functions reduce the exploration burden by providing
 119 informed starting behavior rather than random search [18, 19, 1]. Offline data and world models
 120 ground the learning process in prior experience, mitigating the cold-start pathologies that degrade
 121 network capacity from the earliest updates [7, 4]. These gains are substantial and well documented.
 122 However, as priors grow stronger, the previously negligible side of each tension becomes increasingly
 123 consequential, and the full two-sided character of both oppositions re-emerges.

124 This shift is visible along both axes. • On the *stability-plasticity* axis, aggressive online updates risk
 125 catastrophic forgetting at the offline-to-online transition [8, 43, 44], while distributional mismatch
 126 between offline data and online rollouts can destabilize value estimation [7]. Meanwhile, plasticity
 127 failure is no longer only a matter of network capacity, such as dormant neurons or rank collapse [41,
 128 37], but also one of optimization bias, where strong initialization shapes the loss landscape in ways
 129 that bias subsequent training toward the prior [11, 45]. • On the *exploration-exploitation* axis,
 130 strong priors can narrow the agent’s policy toward pre-trained behavioral modes [46, 47], while
 131 optimization against learned reward models can drive the agent to exploit proxy scores rather than
 132 explore genuinely better behavior [14]. Meanwhile, exploration failure shifts from an inability to
 133 reach informative states [33] to a difficulty in moving beyond the signals carried by the prior [46]. In
 134 both cases, priors do not simply solve the bottlenecks; they restore and sharpen the full two-sided
 135 character of a tension that task difficulty had compressed into a one-sided bottleneck.

136 **Priors as Bounded Commitments.** The sharpened tensions described above share a common
 137 structural root. Every offline prior encodes knowledge from a source setting that may differ from the
 138 deployment environment. The agent cannot fully determine where this knowledge remains valid and
 139 where it does not, yet it must *rely on* the prior to learn efficiently and *adapt beyond* it where the prior
 140 falls short. We call this epistemic status a **bounded commitment**: the prior is valuable but its scope of
 141 validity is bounded, and the agent must commit to using it without knowing those bounds precisely.
 142 The insight that source knowledge has limited validity in new contexts is well established across
 143 Bayesian RL [48, 49], transfer learning [50], and adaptive offline RL [51]. The concept of ‘bounded
 144 commitments’ is precisely a name given to this fundamental insight, as it applies to all previous types
 145 of online RL. Table 2 makes this concrete: each functional role from Sec. 2 introduces a reliance
 146 parameter (μ, λ, β) that governs how strongly the learner commits to its priors.

Table 2: Structural changes that offline priors introduce to online RL optimization. The reliance parameters μ, λ, β each govern how strongly the learner relies on the prior; setting any to zero recovers the from-scratch case. $\mathcal{M}_{\text{prior}}$ denotes the effective (possibly approximate) environment implied by auxiliary sources $(\mathcal{D}, \hat{M}, \hat{R})$, and $J(\theta)$ is shorthand for $J(\theta; \mathcal{M}_{\text{deploy}})$.

	Online RL from Scratch	Online RL with Offline Priors	Reliance
(a) Initialization	$\theta_0 \sim \text{random}$	$\theta_0 = \mu \theta_{\text{prior}} + (1 - \mu) \theta_{\text{rand}}$	$\mu \in [0, 1]$
(b) Reference	$\max_{\theta} J(\theta)$	$\max_{\theta} J(\theta) - \lambda L_{\text{ref}}(\theta)$	$\lambda \geq 0$
(c) Auxiliary	$\max_{\theta} J(\theta; \mathcal{M}_{\text{deploy}})$	$\max_{\theta} (1 - \beta) J(\theta; \mathcal{M}_{\text{deploy}}) + \beta J(\theta; \mathcal{M}_{\text{prior}})$	$\beta \in [0, 1]$

147 Because the validity boundary of each prior is never fully knowable, every reliance configuration is
 148 necessarily a bet. Although Table 2 expresses reliance through continuous parameters for analytical
 149 clarity, $\mu, \lambda,$ and β abstract over method-level design decisions such as whether to use pre-trained
 150 weights, whether to constrain the objective, and whether to supplement the replay buffer. Too much

151 reliance risks trapping the agent in knowledge that does not hold; too little wastes knowledge that
152 could have accelerated learning. This difficulty is structural, not algorithmic, and is compounded
153 by coupling: the KL penalty λ in RLHF, for instance, simultaneously controls policy stability and
154 exploration freedom [46], so that adjusting one tension inevitably affects the other. Because the right
155 reliance depends on how well each prior matches the deployment environment, Sec. 4 examines how
156 this challenge manifests empirically and what it implies for how the field evaluates progress.

Takeaway: Offline priors are fundamentally bounded commitments: valuable knowledge whose
scope of validity the agent can never fully determine. This is why offline priors, despite their
substantial benefits, sharpen rather than resolve the core tensions of online RL.

157

158 4 From Pursuing Universality to Diagnosis-Driven Flexibility and Adaptivity

159 Despite the structural tensions identified in Sec. 3, the current dominant research paradigm continues
160 to evaluate progress by comparing methods on fixed benchmarks and seeking algorithms that perform
161 well across the board. This approach implicitly assumes that a single reliance configuration can be
162 universally optimal. This section presents evidence against that assumption (Sec. 4.1) and argues that
163 reliable deployment requires diagnostic infrastructure rather than universal methods (Sec. 4.2).

164 4.1 Tension Management Has No Universal Optimum

165 How well each prior matches the deployment environment determines the optimal reliance config-
166 uration: how closely the offline data covers the online distribution, how accurately a world model
167 reflects real dynamics, and how faithfully a pre-trained policy captures the behavior the task requires.
168 Since this match varies across tasks and conditions, the optimum varies with it.

169 **Illustrative Experiments.** We demonstrate this non-universality through controlled experiments in
170 the offline-to-online RL setting. We choose this setting because it abstracts away domain-specific
171 engineering details present in LLM post-training, sim-to-real transfer, and VLA fine-tuning, allowing
172 each reliance parameter from Table 2 to be toggled independently. • For initialization (μ), we compare
173 starting from the offline-trained policy weights *versus* resetting to a randomly initialized network.
174 • For reference (λ), we compare maintaining conservative value penalties throughout fine-tuning
175 *versus* dropping all conservatism and using unconstrained optimization [19, 52]. • For auxiliary
176 information (β), we compare retaining the offline dataset in the replay buffer *versus* discarding it and
177 learning from online data only [7, 8]. Fig. 2 presents representative task pairs from our experiments
178 (full results across all tasks in Appendix A). In each column, the same binary choice produces
179 opposite outcomes across tasks. These reversals are not random variation: they arise because the
180 optimal reliance depends on a complex interaction between the properties of the prior, the structure
181 of the deployment task, and the degree to which they match. Since all three factors vary across
182 settings, no single configuration is reliably beneficial. Recent work has made this pattern precise
183 in offline-to-online RL by identifying distinct regimes in which the optimal strategy qualitatively
184 flips [11]. Our experiments extend this observation by showing that non-universality spans all three
185 reliance dimensions independently, not only the interaction between initialization and data retention.

186 **A Cross-Domain Pattern.** The same non-universality appears far beyond offline-to-online RL. • In
187 LLM post-training, the debate over reference regularization illustrates the point directly: removing
188 the KL penalty improves reasoning performance on verifiable tasks [13], but the penalty remains
189 essential for preventing reward hacking when rewards come from learned models [1, 14]. Recent
190 work shows that gradient regularization can outperform KL penalties entirely in some regimes while
191 failing in others [53], and that static length penalties help efficiency on easy tasks but hurt accuracy
192 on hard ones [54]. • In sim-to-real robotics, a large-scale study across three robot platforms finds
193 that widely used algorithmic defaults can be harmful on physical hardware [15], and that end-to-end
194 policy fine-tuning collapses in real-world deployment even when it succeeds in simulation [16].
195 • In vision-language-action model fine-tuning, sequential adaptation with LoRA works remarkably
196 well for large pretrained VLAs but collapses when any single ingredient is removed [55]. • In
197 model-based RL, explicit conservatism helps on high-coverage datasets but fails on low-quality data,
198 where Bayesian approaches without conservatism perform better [56]. Across all these settings, the
199 underlying pattern is the same: the optimal reliance level depends not on the method alone but on
200 the properties of the prior, the demands of the deployment task, and how well the two align. That

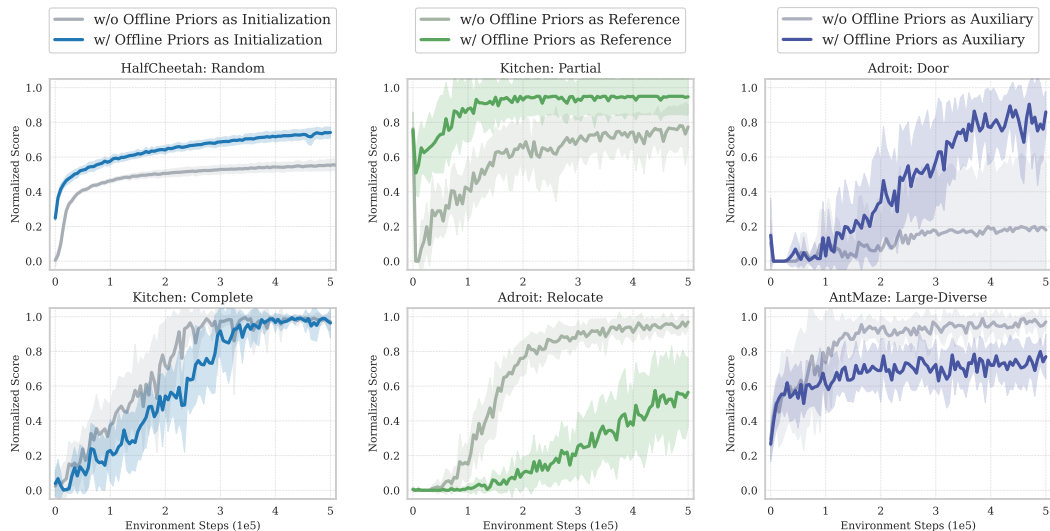


Figure 2: No single reliance configuration is universally optimal. Each column isolates one reliance parameter from Table 2: initialization (μ , left), reference (λ , middle), and auxiliary (β , right). The top row shows a task where stronger reliance on the prior helps; the bottom row shows a task where the same choice hurts. **Left:** RLPD with versus without offline-trained initialization. **Middle:** Cal-QL’s conservative penalty versus unconstrained updates, with both initialization and offline data present. **Right:** retaining versus discarding offline data, with both initialization and conservative penalty present. Full results and additional controlled comparisons are provided in Appendix A.

201 this appears independently across communities with different methods, vocabularies, and evaluation
 202 practices suggests it is structural rather than incidental.

203 **The Inherent Limits of Benchmark-Driven Evaluation.** The evidence above implies that each
 204 method implicitly encodes a particular reliance configuration. A benchmark ranking therefore reflects
 205 how well that configuration matches the evaluation conditions, not a universal ordering of methods.
 206 When the deployment setting changes, the ranking can change with it. The standard response to
 207 this fragility has been to broaden evaluation by testing on more tasks, more environments, and more
 208 data conditions. However, expanding the evaluation suite tends to multiply contradictory findings
 209 rather than eliminate them, and aggregate metrics compress these disagreements into a single ranking
 210 that explains none of them. As the empirical record grows, what accumulates is not transferable
 211 understanding but an expanding catalogue of condition-specific performance rankings.

212 Recent position papers have raised compatible concerns: that rigorous RL benchmarking is prohi-
 213 bitively expensive [57], that aggregate scores obscure fragile generalization [58], and that stan-
 214 dard protocols hide the true cost of hyperparameter selection [59] and mask deployment non-
 215 stationarity [60]. The non-universality we identify raises a more fundamental concern: even
 216 methodologically sound benchmarking cannot produce transferable conclusions when the optimal
 217 configuration is itself deployment-dependent.

218 The fundamental limitation is not the size or quality of any benchmark but the kind of question that
 219 benchmark comparisons can answer. Rankings order methods within a fixed setting; they do not
 220 reveal which properties of the prior and the deployment task govern whether a design choice helps
 221 or hurts. Progress requires a different kind of question: not which method is best, but what each
 222 deployment needs. Without infrastructure to answer that question, the field risks an indefinite cycle
 223 of benchmark expansion without convergent insight.

Takeaway: No single reliance configuration is universally optimal: the same design choice that helps in one setting can hurt in another. Benchmark comparisons cannot resolve this because they answer which method wins, not what determines when each choice helps or hurts.

224

225 4.2 From Benchmark-Driven to Diagnosis-Driven Tension Management

226 More or better benchmarks cannot overcome a limitation inherent in ranking-based evaluation itself.
 227 The field needs to change not the evaluation tools but the question those tools are designed to answer.

228 Concretely, we advocate a shift from *benchmark-driven* to *diagnosis-driven* tension management:
229 from letting fixed rankings guide method selection to letting deployment-specific evidence guide how
230 the learner relates to its priors throughout training.

231 **The Core Distinction.** In the benchmark-driven paradigm, reliance on each prior is configured before
232 deployment based on aggregate evaluation results, and remains fixed or follows a pre-determined
233 schedule once training begins. Online experience serves only to update the policy, even though
234 every transition also carries evidence about the prior itself, including whether value estimates still
235 align with the deployment environment, whether offline data still overlaps with online rollouts, and
236 whether reference behavior still serves the task. This evidence goes unused. The diagnosis-driven
237 paradigm treats it as a first-class signal. The same transitions that update the policy also reveal
238 whether value estimates remain calibrated, whether offline data still provides useful grounding, and
239 whether reference behavior continues to guide learning productively. This evidence informs reliance
240 decisions throughout training, determining when to trust the prior and when to move beyond it. In
241 short, diagnosis-driven tension management is the practice of systematically extracting evidence
242 about prior validity from online interaction and using it to guide reliance decisions.

243 This is not merely a methodological preference. However carefully a prior is constructed, its validity
244 boundary in any specific deployment remains uncertain until interaction begins. Offline evaluation
245 can estimate prior quality in general, but cannot determine which specific aspects will hold or fail
246 under new conditions. Designing better priors or more robust algorithms can reduce the frequency
247 of severe mismatches, but cannot eliminate the underlying uncertainty: the agent is always using
248 knowledge acquired elsewhere to act in an environment it has not yet fully observed. Online evidence
249 is therefore indispensable, because no other evidence about deployment-specific prior validity exists.

250 **The Diagnostic Dimensions.** The information available for making reliance decisions changes
251 fundamentally over the course of training. Before online interaction begins, the practitioner knows
252 only the prior and the task. Once online learning begins, each interaction reveals where the prior
253 holds and where it does not. These two stages require two complementary forms of diagnosis.

254 The first form is *prior-deployment match assessment*. Before online training begins, the practitioner
255 should estimate how well each prior fits the task at hand, based on properties that are observable
256 without interaction, such as data coverage, policy quality, or model accuracy [11, 61]. Recent work
257 demonstrates that even coarse estimates carry actionable information: comparing offline policy
258 quality with data quality can already determine which component the practitioner should anchor
259 on [11], and lightweight metrics can predict whether a dataset will support effective fine-tuning [61].
260 By informing the initial configuration, match assessment enables *flexibility* across deployments.

261 The second form is *tension dynamics monitoring*. Training is not static: aspects that were initially
262 valuable may become outdated as the agent’s own experience grows, turning helpful guidance into a
263 binding constraint. The field already recognizes this implicitly. Methods that anneal conservative
264 penalties [62], schedule warmup phases [8], or periodically reset the reference model [63] all assume
265 that reliance should change during training, but make these adjustments on a fixed schedule rather
266 than in response to observed learning dynamics. Emerging work shows that measurement-driven
267 adjustment is feasible: plasticity metrics can detect capacity loss during training [41, 38, 64], adaptive
268 replay buffers can rebalance data sources based on relevance signals [65], and reward-model monitors
269 can flag proxy divergence [66]. However, these tools remain fragmented across domains, each
270 targeting a specific failure mode rather than assessing prior validity as a whole. By tracking how prior
271 validity evolves during training, dynamics monitoring enables *adaptivity* within each deployment.

272 In practice, these two dimensions interact. Online evidence gathered through dynamics monitoring
273 can retrospectively validate or revise the initial match assessment, and a better initial assessment
274 reduces the burden on runtime monitoring. This loop is almost entirely absent today. Most methods
275 fix their reliance configuration at the start of training or adjust it on a predetermined schedule, without
276 using online evidence to revise the initial assessment. Closing this loop through principled diagnostic
277 infrastructure is, in our view, necessary for making online RL with offline priors reliably deployable.

Takeaway: The field should shift from benchmark-driven to diagnosis-driven tension management. Online experience carries evidence not only about the task but also about the ongoing validity of each prior. Diagnosis-driven tension management is the practice of extracting this evidence and using it to guide reliance decisions, enabling both flexible and adaptive deployment.

279 5 Objections and Counterarguments

280 We consider five potential objections to our call for diagnosis-driven online RL with offline priors.

281 ***This is an engineering challenge, not a scientific insight.*** It is fair to ask whether this paper merely
282 names a phenomenon that practitioners already navigate daily: different deployments need different
283 configurations, and finding good ones is routine engineering. Each community already handles its
284 own regime dependence, but in isolation: the offline-to-online RL community studies initialization
285 versus data retention [11], the RLHF community debates KL penalty strength [13, 53], and the
286 sim-to-real community weighs fine-tuning against freezing [15]. Our contribution is recognizing
287 that these are superficially different expressions of the same structural phenomenon and that the
288 same diagnostic principles apply across all of them. Cross-community unification of this kind
289 has consistently been treated as scientific contribution in ML, from the formalization of transfer
290 learning [67] to the systematization of catastrophic forgetting [68]. A reader may also interpret our
291 formalization as a hyperparameter tuning problem, but the decisions that μ , λ , and β abstract over are
292 not points on a continuous search grid. They abstract over architectural and algorithmic choices such
293 as using pretrained versus random initialization, imposing versus dropping a conservative penalty,
294 and retaining versus discarding offline data. Furthermore, prior validity shifts as online experience
295 accumulates, so the right reliance level at the start of training may not remain right later.

296 ***Better offline priors solve the problem at the source.*** Rather than diagnosing bounded commitment
297 during deployment, one could try to prevent it upstream. Recent work pursues this through larger
298 and more diverse datasets [5, 69], adaptive offline objectives that preserve revision capacity [51, 48],
299 world foundation models that encode broad physical priors [26, 27], and Bayesian formulations
300 that build uncertainty directly into the learned policy [49, 56]. These efforts are valuable and often
301 dramatically improve transfer and generalization. However, in every major deployment paradigm,
302 online adaptation remains a necessary stage: foundation-model priors are designed as efficient
303 initializations to be fine-tuned, not as finished policies for arbitrary deployment [70], and adaptive
304 offline objectives explicitly aim to preserve the capacity for later revision, not to remove the need
305 for it [51, 48]. The reason is structural: online RL exists in the pipeline precisely because the prior
306 does not fully solve the deployment task, and improving beyond the prior necessarily means entering
307 territory where its guidance is no longer reliable. Better priors extend the region where reliance is
308 safe; diagnosis addresses what happens at and beyond its boundary. The two are complementary, and
309 advances in either make the other more effective.

310 ***Scaling and algorithmic progress will resolve this.*** A longer-term version of the previous argument
311 holds that continued progress in model scale, data scale, and algorithm design will eventually make
312 bounded commitment negligible, rendering diagnostic infrastructure a premature investment. There
313 is real evidence for this view: at sufficient scale, some classical pathologies weaken. Large pretrained
314 VLAs show little forgetting during continual adaptation [55, 71], and larger language models exhibit
315 more efficient RL post-training [72]. However, the strongest scaling results are themselves regime-
316 dependent. The VLA recipe that eliminates forgetting requires a specific combination of large
317 model, parameter-efficient tuning, and on-policy RL; removing any ingredient causes collapse [55].
318 Fine-tuning capacity does not transfer uniformly across tasks and embodiments [73, 74], and platform-
319 dependent defaults persist even with state-of-the-art algorithms [15]. What scale changes is not
320 whether deployment-specific choices matter but which ones matter most. Meanwhile, scaling expands
321 the range of deployments the field attempts to address, introducing new embodiments, task types, and
322 deployment conditions faster than any single advance can uniformly cover. This makes principled
323 diagnosis more necessary as the field scales, not less.

324 ***The field only needs better benchmarks, not diagnostic infrastructure.*** Recent position papers have
325 proposed valuable reforms to RL evaluation: accounting for tuning costs [59], restricting lifetime
326 access [60], and testing for fragile generalization [58, 57]. This concern is well founded, and we
327 agree that evaluation methodology needs reform. However, even perfect benchmarks answer a
328 different question than diagnostics do. Benchmarks tell us which method tends to work under which
329 conditions; diagnostics tell us whether a specific deployment meets those conditions, and whether
330 the answer is changing as training proceeds. The second question requires online evidence that only
331 deployment interaction can generate, which is why benchmarks and diagnostics are complementary:
332 one narrows the space of candidate methods, the other guides their configuration during training.

333 ***Deep RL is too opaque for reliable diagnosis.*** Measuring prior validity is genuinely harder than
334 measuring network capacity, and no general-purpose diagnostic toolkit exists today. However, the

335 relevant question is not why the network behaves as it does but whether the prior is still helping. The
336 former requires interpretability, which remains hard. The latter requires only observable quantities:
337 performance trends, distribution overlap, and prediction accuracy [11, 61]. Practical building blocks
338 already exist across domains: discrete regime distinctions guide method selection in offline-to-online
339 RL [11], reward-quality monitors flag proxy divergence in RLHF [66], and uncertainty estimates
340 weight synthetic data by reliability in model-based RL [75]. Each measures observable signals and
341 translates them into actionable decisions. The plasticity literature shows this trajectory is viable: from
342 informal recognition to systematic measurement to a productive diagnostic subfield with reusable
343 tools, all within a few years [41, 36]. Prior validity diagnosis is at an earlier stage of the same
344 progression. Furthermore, diagnosis need not be perfect to be productive. Even coarse measurements
345 improve on decisions that would otherwise be made without evidence, and each deployment that uses
346 diagnostic signals generates insights that sharpen future tools. This self-reinforcing cycle between
347 diagnostics and deployment is a promising path toward reliable online RL with offline priors.

348 6 Conclusion

349 This paper has argued that offline priors fundamentally reshape the structure of online RL. By in-
350 troducing knowledge whose scope of validity the agent cannot fully determine, priors transform the
351 one-sided bottlenecks of from-scratch learning into genuine two-sided tensions. We have formal-
352 ized this through the concept of bounded commitment, shown empirically that no single reliance
353 configuration is universally optimal, and argued that the field should shift from benchmark-driven to
354 diagnosis-driven tension management, in which online experience is used not only to learn the task
355 but also to assess prior validity and guide reliance decisions toward flexible and adaptive deployment.

356 **Research Opportunities.** If the field adopts diagnosis-driven tension management, several opportu-
357 nities open up that the current paradigm does not naturally support.

- 358 • *A new class of research contributions.* Under the benchmark-driven paradigm, contributions are
359 measured primarily by performance gains. Diagnosis-driven research values a different kind
360 of output: not a method that wins on a benchmark, but a signal that predicts when a design
361 choice helps or hurts, a metric that assesses prior-deployment match, or a monitor that tracks
362 prior validity during training [11, 61]. Work of this kind already exists but is typically framed as
363 supporting analysis rather than a primary contribution [76, 77, 78].
- 364 • *Cross-community knowledge transfer.* Currently, each community rediscovers similar failure
365 modes in isolation: catastrophic forgetting in offline-to-online RL, reward hacking in RLHF,
366 reality-gap collapse in sim-to-real transfer. If these are recognized as manifestations of the same
367 structural phenomenon, diagnostic tools developed in one community can inform practice in
368 others. Plasticity metrics illustrate this potential: dormant neuron ratios originated in a specific
369 experimental setting [41] but now serve as reusable diagnostics across tasks and algorithms.
- 370 • *Deployment as a source of scientific knowledge.* Under the current paradigm, deployment is the
371 endpoint of research: methods are developed, evaluated, and then applied. Diagnosis-driven
372 deployment inverts this relationship. Every deployment that uses diagnostic evidence generates
373 insights about the conditions under which each prior holds or fails, feeding back into the design
374 of better tools and more informed future deployments.
- 375 • *Methods aware of their own bounded commitments.* The diagnosis-driven paradigm also changes
376 how methods themselves are designed. Rather than seeking algorithms that perform well across
377 the board, researchers can design methods that are explicitly aware of the boundaries of their
378 prior knowledge and capable of adjusting their own reliance as those boundaries are revealed
379 during training. This represents a shift from optimizing for average-case performance to building
380 in the capacity for deployment-specific adaptation.

381 **A Broader Perspective.** Underlying these opportunities is a more fundamental question: *what kind*
382 *of knowledge should the field be accumulating?* Under the benchmark-driven paradigm, the field
383 accumulates methods: each validated under particular conditions, each adding to a growing catalogue
384 that transfers poorly across deployments. The diagnosis-driven paradigm instead accumulates
385 understanding: not which method wins where, but what determines when each approach works
386 and why. Each insight about the conditions of success informs not only current practice but future
387 method design. Individual methods will be superseded. Understanding of the fundamental structure
388 of learning persists, compounds, and shapes whatever comes next.

References

- 389
390 [1] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang,
391 Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with
392 human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- 393 [2] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang,
394 Shirong Ma, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement
395 learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 396 [3] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of
397 robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and
398 automation (ICRA)*, pages 3803–3810. IEEE, 2018.
- 399 [4] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through
400 world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 401 [5] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch,
402 Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2: A dataset for robot learning at
403 scale. In *Conference on Robot Learning*, pages 1723–1736. PMLR, 2023.
- 404 [6] Jianlan Luo, Zheyuan Hu, Charles Xu, You Liang Tan, Jacob Berg, Archit Sharma, Stefan Schaal, Chelsea
405 Finn, Abhishek Gupta, and Sergey Levine. Serl: A software suite for sample-efficient robotic reinforcement
406 learning. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 16961–16969.
407 IEEE, 2024.
- 408 [7] Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning with
409 offline data. In *International Conference on Machine Learning*, pages 1577–1594. PMLR, 2023.
- 410 [8] Zhiyuan Zhou, Andy Peng, Qiyang Li, Sergey Levine, and Aviral Kumar. Efficient online reinforcement
411 learning fine-tuning need not retain offline data. In *The Thirteenth International Conference on Learning
412 Representations*, 2025.
- 413 [9] David Silver and Richard S Sutton. Welcome to the era of experience. *Google AI*, 1:11, 2025.
- 414 [10] Alex Lewandowski, Aditya A. Ramesh, Edan Meyer, Dale Schuurmans, and Marlos C. Machado. The
415 world is bigger! a computationally-embedded perspective on the big world hypothesis. In *The Thirty-ninth
416 Annual Conference on Neural Information Processing Systems*, 2025.
- 417 [11] Lu Li, Tianwei Ni, Yihao Sun, and Pierre-Luc Bacon. The three regimes of offline-to-online reinforcement
418 learning. *arXiv preprint arXiv:2510.01460*, 2025.
- 419 [12] Shenzi Wang, Qisen Yang, Jiawei Gao, Matthieu Lin, Hao Chen, Liwei Wu, Ning Jia, Shiji Song, and
420 Gao Huang. Train once, get a family: State-adaptive balances for offline-to-online reinforcement learning.
421 *Advances in Neural Information Processing Systems*, 36:47081–47104, 2023.
- 422 [13] Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan,
423 Gaohong Liu, Juncai Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng,
424 Yuxuan Tong, Chi Zhang, Mofan Zhang, Ru Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaye Chen,
425 Jiangjie Chen, Chengyi Wang, Hongli Yu, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-
426 Ying Ma, Ya-Qin Zhang, Lin Yan, Yonghui Wu, and Mingxuan Wang. DAPO: An open-source LLM
427 reinforcement learning system at scale. In *The Thirty-ninth Annual Conference on Neural Information
428 Processing Systems*, 2025.
- 429 [14] Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In
430 *International Conference on Machine Learning*, pages 10835–10866. PMLR, 2023.
- 431 [15] Yarden As, Dhruva Tirumala, René Zurbrugg, Chenhao Li, Stelian Coros, Andreas Krause, and Markus
432 Wulfmeier. What matters for simulation to online reinforcement learning on real robots. *arXiv preprint
433 arXiv:2602.20220*, 2026.
- 434 [16] Jacob Levy, Tyler Westenbroek, Kevin Huang, Fernando Palafox, Patrick Yin, Shayegan Omidshafiei,
435 Dong-Ki Kim, Abhishek Gupta, and David Fridovich-Keil. Simulation distillation: Pretraining world
436 models in simulation for rapid real-world adaptation. *arXiv preprint arXiv:2603.15759*, 2026.
- 437 [17] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G Bellemare.
438 Reincarnating reinforcement learning: Reusing prior computation to accelerate progress. In Alice H.
439 Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information
440 Processing Systems*, 2022.

- 441 [18] Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online reinforcement
442 learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 443 [19] Mitsuhiko Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral Kumar,
444 and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. *Advances in
445 Neural Information Processing Systems*, 36:62244–62269, 2023.
- 446 [20] Andrew Wagenmaker, Kevin Huang, Liyiming Ke, Kevin Jamieson, and Abhishek Gupta. Overcoming the
447 sim-to-real gap: Leveraging simulation to learn to explore for real-world RL. In *The Thirty-eighth Annual
448 Conference on Neural Information Processing Systems*, 2024.
- 449 [21] Nicklas Hansen, Hao Su, and Xiaolong Wang. TD-MPC2: Scalable, robust world models for continuous
450 control. In *The Twelfth International Conference on Learning Representations*, 2024.
- 451 [22] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan
452 Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic
453 control. In *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023.
- 454 [23] Yanjiang Guo, Jianke Zhang, Xiaoyu Chen, Xiang Ji, Yen-Jen Wang, Yucheng Hu, and Jianyu Chen.
455 Improving vision-language-action model with online reinforcement learning. In *2025 IEEE International
456 Conference on Robotics and Automation (ICRA)*, pages 15665–15672. IEEE, 2025.
- 457 [24] Scott Fujimoto and Shixiang (Shane) Gu. A minimalist approach to offline reinforcement learning. In
458 M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in
459 Neural Information Processing Systems*, volume 34, pages 20132–20145. Curran Associates, Inc., 2021.
- 460 [25] Bohan Hou, Gen Li, Jindou Jia, Tuo An, Xinying Guo, Sicong Leng, Haoran Geng, Yanjie Ze, Tatsuya
461 Harada, Philip Torr, Oier Mees, Marc Pollefeys, Zhuang Liu, Jiajun Wu, Pieter Abbeel, Jitendra Malik,
462 Yilun Du, and Jianfei Yang. World model for robot learning: A comprehensive survey. *arXiv preprint
463 arXiv:2605.00080*, 2026.
- 464 [26] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay,
465 Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv
466 preprint arXiv:2501.03575*, 2025.
- 467 [27] GigaBrain Team, Boyuan Wang, Bohan Li, Chaojun Ni, Guan Huang, Guosheng Zhao, Hao Li, Jie Li,
468 Jindi Lv, Jingyu Liu, et al. Gigabrain-0.5 m*: a vla that learns from world model-based reinforcement
469 learning. *arXiv preprint arXiv:2602.12099*, 2026.
- 470 [28] Seonghyeon Ye, Yunhao Ge, Kaiyuan Zheng, Shenyuan Gao, Sihyun Yu, George Kurian, Suneel Indupuru,
471 You Liang Tan, Chuning Zhu, Jiannan Xiang, Ayaan Malik, Kyungmin Lee, William Liang, Nadun
472 Ranawaka, Jiasheng Gu, Yinzhen Xu, Guanzhi Wang, Fengyuan Hu, Avnish Narayan, Johan Bjorck, Jing
473 Wang, Gwanghyun Kim, Dantong Niu, Ruijie Zheng, Yuqi Xie, Jimmy Wu, Qi Wang, Ryan Julian, Danfei
474 Xu, Yilun Du, Yevgen Chebotar, Scott Reed, Jan Kautz, Yuke Zhu, Linxi "Jim" Fan, and Joel Jang. World
475 action models are zero-shot policies, 2026.
- 476 [29] Yiren Lu, Justin Fu, George Tucker, Xinlei Pan, Eli Bronstein, Rebecca Roelofs, Benjamin Sapp, Brandyn
477 White, Aleksandra Faust, Shimon Whiteson, et al. Imitation is not enough: Robustifying imitation with
478 reinforcement learning for challenging driving scenarios. In *2023 IEEE/RSJ International Conference on
479 Intelligent Robots and Systems (IROS)*, pages 7553–7560. IEEE, 2023.
- 480 [30] Kenneth Shaw, Shikhar Bahl, and Deepak Pathak. Videodex: Learning dexterity from internet videos. In
481 *6th Annual Conference on Robot Learning*, 2022.
- 482 [31] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv
483 preprint arXiv:1911.11361*, 2019.
- 484 [32] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
485 reinforcement learning. *Advances in neural information processing systems*, 33:1179–1191, 2020.
- 486 [33] Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. First return, then explore.
487 *Nature*, 590(7847):580–586, 2021.
- 488 [34] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement learning:
489 A survey. *Information Fusion*, 85:1–22, 2022.
- 490 [35] Evgenii Nikishin, Max Schwarzer, Pierluca D’Oro, Pierre-Luc Bacon, and Aaron Courville. The primacy
491 bias in deep reinforcement learning. In *International conference on machine learning*, pages 16828–16847.
492 PMLR, 2022.

- 493 [36] Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney.
494 Understanding plasticity in neural networks. In *International Conference on Machine Learning*, pages
495 23190–23211. PMLR, 2023.
- 496 [37] Shibhansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mahmood,
497 and Richard S Sutton. Loss of plasticity in deep continual learning. *Nature*, 632(8026):768–774, 2024.
- 498 [38] Guozheng Ma, Lu Li, Sen Zhang, Zixuan Liu, Zhen Wang, Yixin Chen, Li Shen, Xueqian Wang, and
499 Dacheng Tao. Revisiting plasticity in visual reinforcement learning: Data, modules and training stages. In
500 *The Twelfth International Conference on Learning Representations*, 2024.
- 501 [39] Jianye Hao, Tianpei Yang, Hongyao Tang, Chenjia Bai, Jinyi Liu, Zhaopeng Meng, Peng Liu, and Zhen
502 Wang. Exploration in deep reinforcement learning: From single-agent to multiagent domain. *IEEE*
503 *transactions on neural networks and learning systems*, 35(7):8762–8782, 2023.
- 504 [40] Evgenii Nikishin, Junhyuk Oh, Georg Ostrovski, Clare Lyle, Razvan Pascanu, Will Dabney, and Andre
505 Barreto. Deep reinforcement learning with plasticity injection. In *Thirty-seventh Conference on Neural*
506 *Information Processing Systems*, 2023.
- 507 [41] Ghada Sokar, Rishabh Agarwal, Pablo Samuel Castro, and Utku Evci. The dormant neuron phenomenon
508 in deep reinforcement learning. In *International Conference on Machine Learning*, pages 32145–32168.
509 PMLR, 2023.
- 510 [42] Timo Klein, Christoph Luther, Manus McAuliffe, Lukas Miklautz, Claudia Plant, and Sebastian Tschit-
511 atschek. Plasticity loss in deep reinforcement learning: A survey. *arXiv preprint arXiv:2411.04832*,
512 2024.
- 513 [43] Yicheng Luo, Jackie Kay, Edward Grefenstette, and Marc Peter Deisenroth. Finetuning from offline
514 reinforcement learning: Challenges, trade-offs and practical solutions. *arXiv preprint arXiv:2303.17396*,
515 2023.
- 516 [44] Maciej Wolczyk, Bartłomiej Cupiał, Mateusz Ostaszewski, Michał Borkiewicz, Michał Zajac, Razvan
517 Pascanu, Łukasz Kuciński, and Piotr Miłoś. Fine-tuning reinforcement learning models is secretly a
518 forgetting mitigation problem. In *Forty-first International Conference on Machine Learning*, 2024.
- 519 [45] Clare Lyle, Zeyu Zheng, Khimya Khetarpal, Hado Van Hasselt, Razvan Pascanu, James Martens, and Will
520 Dabney. Disentangling the causes of plasticity loss in neural networks. *arXiv preprint arXiv:2402.18762*,
521 2024.
- 522 [46] Rosie Zhao, Alexandru Meterez, Sham M. Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach.
523 Echo chamber: RL post-training amplifies behaviors learned in pretraining. In *Second Conference on*
524 *Language Modeling*, 2025.
- 525 [47] Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang.
526 Does reinforcement learning really incentivize reasoning capacity in LLMs beyond the base model? In
527 *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- 528 [48] Dibya Ghosh, Anurag Ajay, Pulkit Agrawal, and Sergey Levine. Offline rl policies should be trained to be
529 adaptive. In *International Conference on Machine Learning*, pages 7513–7530. PMLR, 2022.
- 530 [49] Hao Hu, Yiqin Yang, Jianing Ye, Chengjie Wu, Ziqing Mai, Yujing Hu, Tangjie Lv, Changjie Fan,
531 Qianchuan Zhao, and Chongjie Zhang. Bayesian design principles for offline-to-online reinforcement
532 learning. In *Forty-first International Conference on Machine Learning*, 2024.
- 533 [50] Wen Zhang, Lingfei Deng, Lei Zhang, and Dongrui Wu. A survey on negative transfer. *IEEE/CAA Journal*
534 *of Automatica Sinica*, 10(2):305–329, 2022.
- 535 [51] Tianwei Ni, Vineet Jain, Akash Karthikeyan, and Pierre-Luc Bacon. From static policies to adaptive priors
536 in offline reinforcement learning. *Preprint*, 2026.
- 537 [52] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum
538 entropy deep reinforcement learning with a stochastic actor. In *International conference on machine*
539 *learning*, pages 1861–1870. Pmlr, 2018.
- 540 [53] Johannes Ackermann, Michael Noukhovitch, Takashi Ishida, and Masashi Sugiyama. Gradient regular-
541 ization prevents reward hacking in reinforcement learning from human feedback and verifiable rewards.
542 *arXiv preprint arXiv:2602.18037*, 2026.

- 543 [54] Keqin Peng, Yuanxin Ouyang, Xuebo Liu, Zhiliang Tian, Ruijian Han, Yancheng Yuan, and Liang Ding.
544 Think dense, not long: Dynamic decoupled conditional advantage for efficient reasoning. *arXiv preprint*
545 *arXiv:2602.02099*, 2026.
- 546 [55] Jiaheng Hu, Jay Shim, Chen Tang, Yoonchang Sung, Bo Liu, Peter Stone, and Roberto Martin-Martin.
547 Simple recipe works: Vision-language-action models are natural continual learners with reinforcement
548 learning. *arXiv preprint arXiv:2603.11653*, 2026.
- 549 [56] Tianwei Ni, Esther Derman, Vineet Jain, Vincent Taboga, Siamak Ravanbakhsh, and Pierre-Luc Ba-
550 con. Long-horizon model-based offline reinforcement learning without conservatism. *arXiv preprint*
551 *arXiv:2512.04341*, 2025.
- 552 [57] Emma Jordan, Adam White, Bruno Castro da Silva, Martha White, and Philip S. Thomas. Position:
553 Benchmarking is limited in reinforcement learning research. In *Forty-first International Conference on*
554 *Machine Learning*, 2024.
- 555 [58] Zihan Chen, Yiming Zhang, Hengguang Zhou, Zenghui Ding, Yining Sun, and Cho-Jui Hsieh. Rethinking
556 rl evaluation: Can benchmarks truly reveal failures of rl methods? *arXiv preprint arXiv:2510.10541*, 2025.
- 557 [59] Ziqi Tang and Xuezhou Zhang. Position: Ignoring hyperparameter tuning costs misleads the development
558 of efficient rl algorithms. *preprint*, 2025.
- 559 [60] Golnaz Mesbahi, Parham Mohammad Panahi, Olya Mastikhina, Steven Tang, Martha White, and Adam
560 White. Position: Lifetime tuning is incompatible with continual reinforcement learning. In *Forty-second*
561 *International Conference on Machine Learning Position Paper Track*, 2025.
- 562 [61] Arip Asadulaev, Fakhri Karray, and Martin Takac. Expert or not? assessing data quality in offline
563 reinforcement learning. *arXiv preprint arXiv:2510.12638*, 2025.
- 564 [62] Geonwoo Cho, Jaegyun Im, Doyoon Kim, and Lexin Li. Annealing bridges offline and online RL, 2025.
- 565 [63] Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. ProRL:
566 Prolonged reinforcement learning expands reasoning boundaries in large language models. In *The Thirty-*
567 *ninth Annual Conference on Neural Information Processing Systems*, 2025.
- 568 [64] Guowei Xu, Ruijie Zheng, Yongyuan Liang, Xiyao Wang, Zhecheng Yuan, Tianying Ji, Yu Luo, Xiaoyu
569 Liu, Jiaxin Yuan, Pu Hua, Shuzhen Li, Yanjie Ze, Hal Daumé III, Furong Huang, and Huazhe Xu. Drm:
570 Mastering visual reinforcement learning through dormant ratio minimization. In *The Twelfth International*
571 *Conference on Learning Representations*, 2024.
- 572 [65] Chihyeon Song, Jaewoo Lee, and Jinkyoo Park. Adaptive replay buffer for offline-to-online reinforcement
573 learning. In *The 29th International Conference on Artificial Intelligence and Statistics*, 2026.
- 574 [66] Yuchun Miao, Liang Ding, Sen Zhang, Rong Bao, Lefei Zhang, and Dacheng Tao. Information-theoretic
575 reward modeling for stable rlhf: Detecting and mitigating reward hacking. *arXiv preprint arXiv:2510.13694*,
576 2025.
- 577 [67] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and*
578 *data engineering*, 22(10):1345–1359, 2009.
- 579 [68] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu,
580 Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic
581 forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- 582 [69] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti,
583 Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A
584 large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024.
- 585 [70] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael
586 Rafailov, Ethan P Foster, Pannag R Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ
587 Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. OpenVLA: An open-source
588 vision-language-action model. In *8th Annual Conference on Robot Learning*, 2024.
- 589 [71] Huihan Liu, Changyeon Kim, Bo Liu, Minghuan Liu, and Yuke Zhu. Pretrained vision-language-action
590 models are surprisingly resistant to forgetting in continual learning. *arXiv preprint arXiv:2603.03818*,
591 2026.
- 592 [72] Zelin Tan, Hejia Geng, Xiaohang Yu, Mulei Zhang, Guancheng Wan, Yifan Zhou, Qiang He, Xiangyuan
593 Xue, Heng Zhou, Yutao Fan, et al. Scaling behaviors of llm reinforcement learning post-training: An
594 empirical study in mathematical reasoning. *arXiv preprint arXiv:2509.25300*, 2025.

- 595 [73] Donghoon Kim, Minji Bae, Unghui Nam, Gyeonghun Kim, Suyun Lee, Kyuhong Shim, and Byonghyo
596 Shim. Adaptive capacity allocation for vision language action fine-tuning. *arXiv preprint arXiv:2603.07404*,
597 2026.
- 598 [74] Xinghang Li, Peiyan Li, Long Qian, Minghuan Liu, Dong Wang, Jirong Liu, Bingyi Kang, Xiao Ma,
599 Xinlong Wang, Di Guo, et al. What matters in building vision–language–action models for generalist
600 robots. *Nature Machine Intelligence*, pages 1–15, 2026.
- 601 [75] Mehran Aghabozorgi, Alireza Moazeni, Yanshu Zhang, and Ke Li. WIMLE: Uncertainty-aware world
602 models with IMLE for sample-efficient continuous control. In *The Fourteenth International Conference on*
603 *Learning Representations*, 2026.
- 604 [76] Guozheng Ma, Lu Li, Haoyu Wang, Zixuan Liu, Pierre-Luc Bacon, and Dacheng Tao. What makes value
605 learning efficient in residual reinforcement learning? *arXiv preprint arXiv:2602.10539*, 2026.
- 606 [77] Johan Obando-Ceron, Walter Mayor, Samuel Lavoie, Scott Fujimoto, Aaron Courville, and Pablo Samuel
607 Castro. Simplicial embeddings improve sample efficiency in actor–critic agents. In *The Fourteenth*
608 *International Conference on Learning Representations*, 2026.
- 609 [78] Isaac Han, Sangyeon Park, Seungwon Oh, Donghu Kim, Hojoon Lee, and KyungJoong Kim. FIRE:
610 Frobenius-isometry reinitialization for balancing the stability–plasticity tradeoff. In *The Fourteenth*
611 *International Conference on Learning Representations*, 2026.
- 612 [79] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
613 data-driven reinforcement learning, 2020.

614 **A Illustrative Experiments: Full Results**

615 Section 4.1 presents representative task pairs to illustrate the non-universality of tension management. This
 616 appendix provides the complete experimental setup and full results across all tasks and reliance dimensions.

617 **Setup.** All experiments are conducted in the offline-to-online RL setting on D4RL benchmarks [79], covering
 618 Adroit manipulation (Pen, Relocate, Door), Kitchen (Complete, Partial, Mixed), AntMaze navigation (Large-
 619 Diverse, Large-Play, Ultra-Diverse), and MuJoCo locomotion (HalfCheetah, Hopper, Walker2D, each with
 620 Random, Medium-Replay, Medium, and Medium-Expert datasets). Each experiment isolates one reliance
 621 parameter from Table 2 by toggling it while holding the others fixed. We organize the experiments into three
 622 groups corresponding to the three functional roles, with two groups further split to reveal interactions between
 623 parameters. Table 3 summarizes the experimental conditions. In all cases, the “with” and “without” conditions
 624 differ in exactly one reliance dimension, enabling controlled comparison.

Table 3: Summary of experimental conditions. Each row isolates one reliance parameter while holding the others fixed. Checkmarks indicate active components.

Experiment	Comparison	Init (μ)	Ref (λ)	Aux (β)
Initialization	w/ init vs w/o init	varies	×	✓
Reference (no aux)	w/ ref vs w/o ref	✓	varies	×
Reference (with aux)	w/ ref vs w/o ref	✓	varies	✓
Auxiliary (no ref)	w/ aux vs w/o aux	✓	×	varies
Auxiliary (with ref)	w/ aux vs w/o aux	✓	✓	varies

625 **Initialization (μ).** This experiment tests whether initializing the online learner from the offline-trained policy
 626 and value function improves over random initialization. Both conditions retain the offline dataset in the replay
 627 buffer following the RLPD protocol [7] and use standard SAC updates without conservative penalties. The
 628 comparison thus isolates the effect of initialization reliance: starting from prior parameters ($\mu = 1$) versus
 629 starting from random parameters ($\mu = 0$), with auxiliary support held constant.

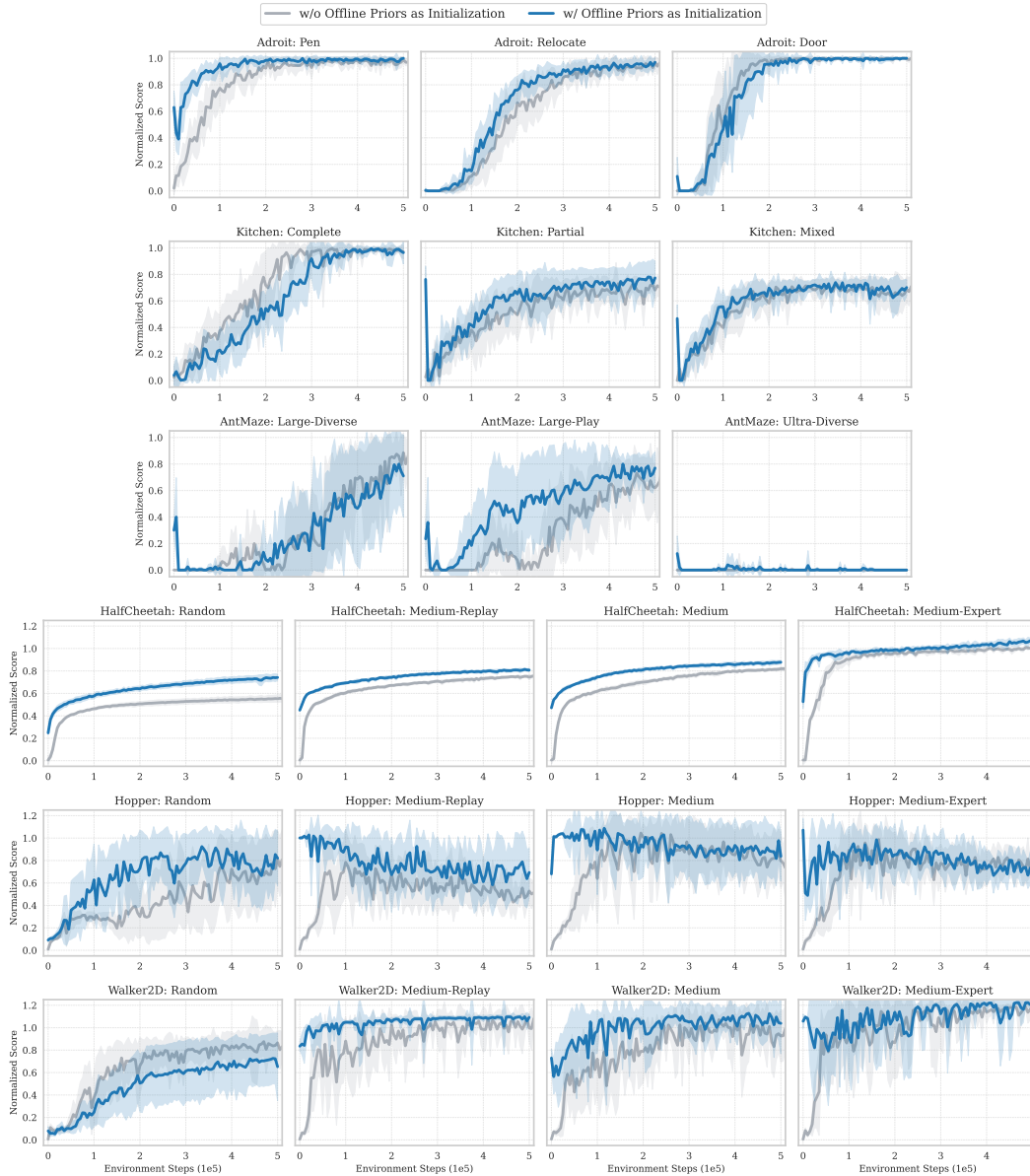


Figure 3: Effect of initialization reliance (μ), with auxiliary data retained in both conditions. Using offline-trained initialization helps on some tasks but hurts on others, illustrating that the optimal μ is task-dependent.

630 **Reference (λ), without auxiliary data.** This experiment tests whether maintaining a conservative value
 631 penalty during online learning (Cal-QL [19]) improves over unconstrained SAC [52]. Both conditions start from
 632 the offline-trained initialization and do not retain offline data in the replay buffer. The comparison isolates the
 633 effect of reference reliance: conservative penalty active ($\lambda > 0$) versus no penalty ($\lambda = 0$), without auxiliary
 634 support.

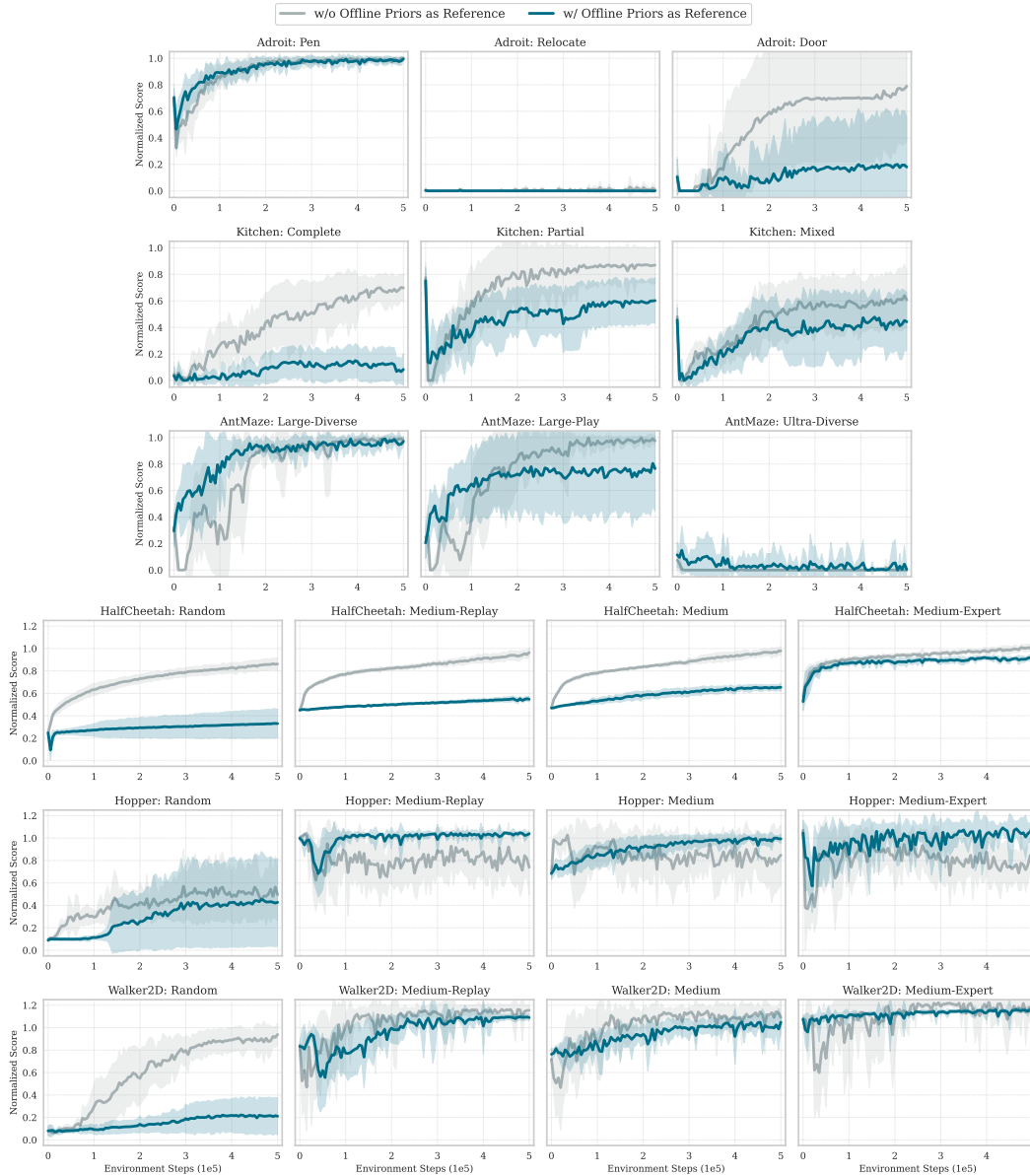


Figure 4: Effect of reference reliance (λ) without auxiliary data. Both conditions use offline-trained initialization. Maintaining conservative value penalties helps on some tasks but hurts on others.

635 **Reference (λ), with auxiliary data.** This experiment repeats the reference comparison with auxiliary data
 636 present: the offline dataset is retained in the replay buffer in both conditions. The comparison isolates whether
 637 using Cal-QL’s conservative penalty (as opposed to standard SAC updates) remains beneficial when offline data
 638 is also available as auxiliary support. Comparing Figures 4 and 5 reveals how the effect of a reference constraint
 639 can itself depend on whether auxiliary data is available, illustrating the coupling between reliance parameters
 640 discussed in Section 3.

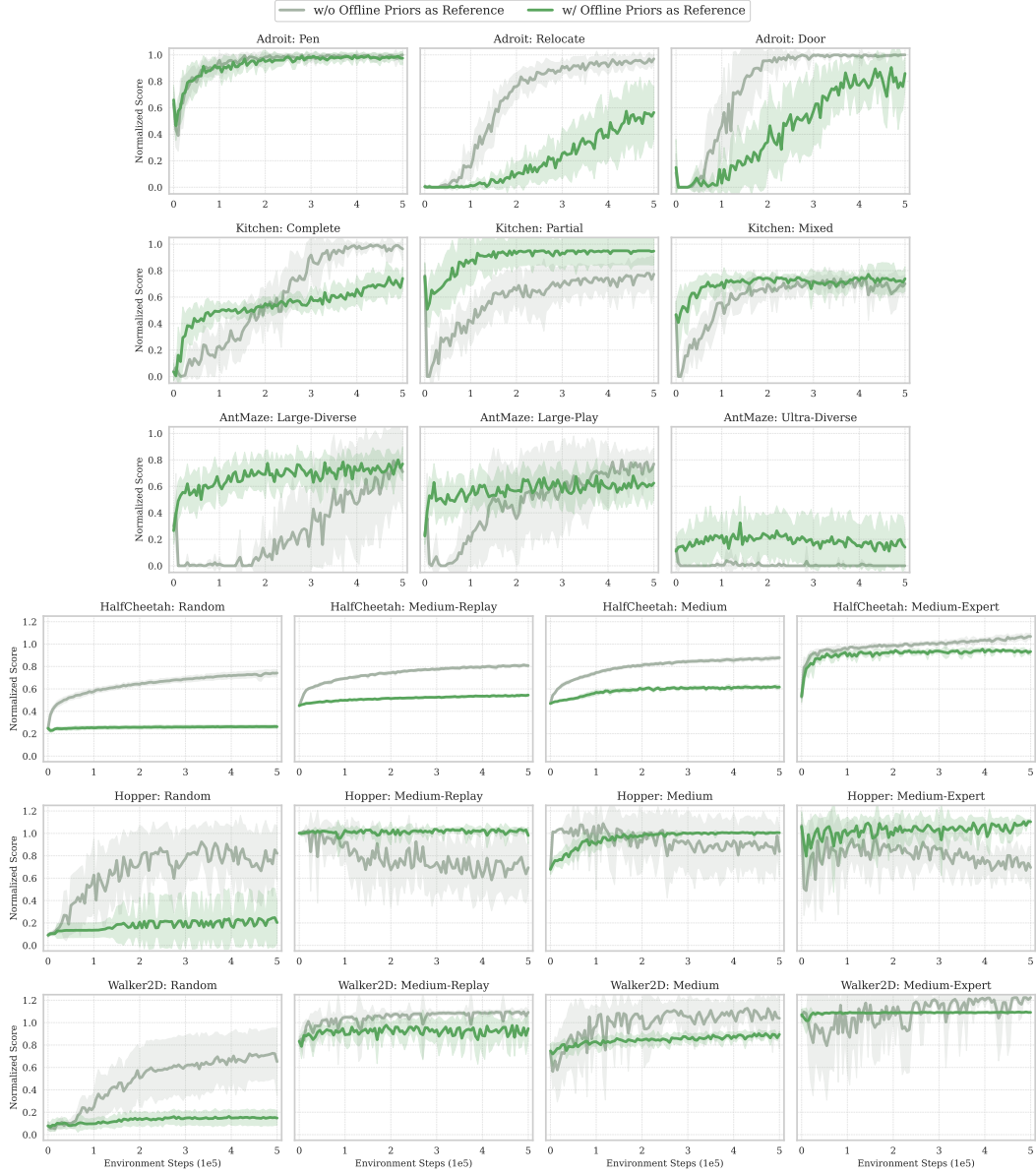


Figure 5: Effect of reference reliance (λ) with auxiliary data ($\beta > 0$). Both conditions use offline-trained initialization and retain offline data. The same conservative penalty that helps or hurts in Figure 4 can behave differently when offline data is also retained, demonstrating interaction between reliance parameters.

641 **Auxiliary (β), without reference constraint.** This experiment tests whether retaining the offline dataset
 642 in the replay buffer improves over discarding it, when the online learner uses standard SAC without conservative
 643 penalties. Both conditions start from the offline-trained initialization. The comparison isolates the effect of
 644 auxiliary reliance: offline data retained ($\beta > 0$) versus online data only ($\beta = 0$), without reference constraint.

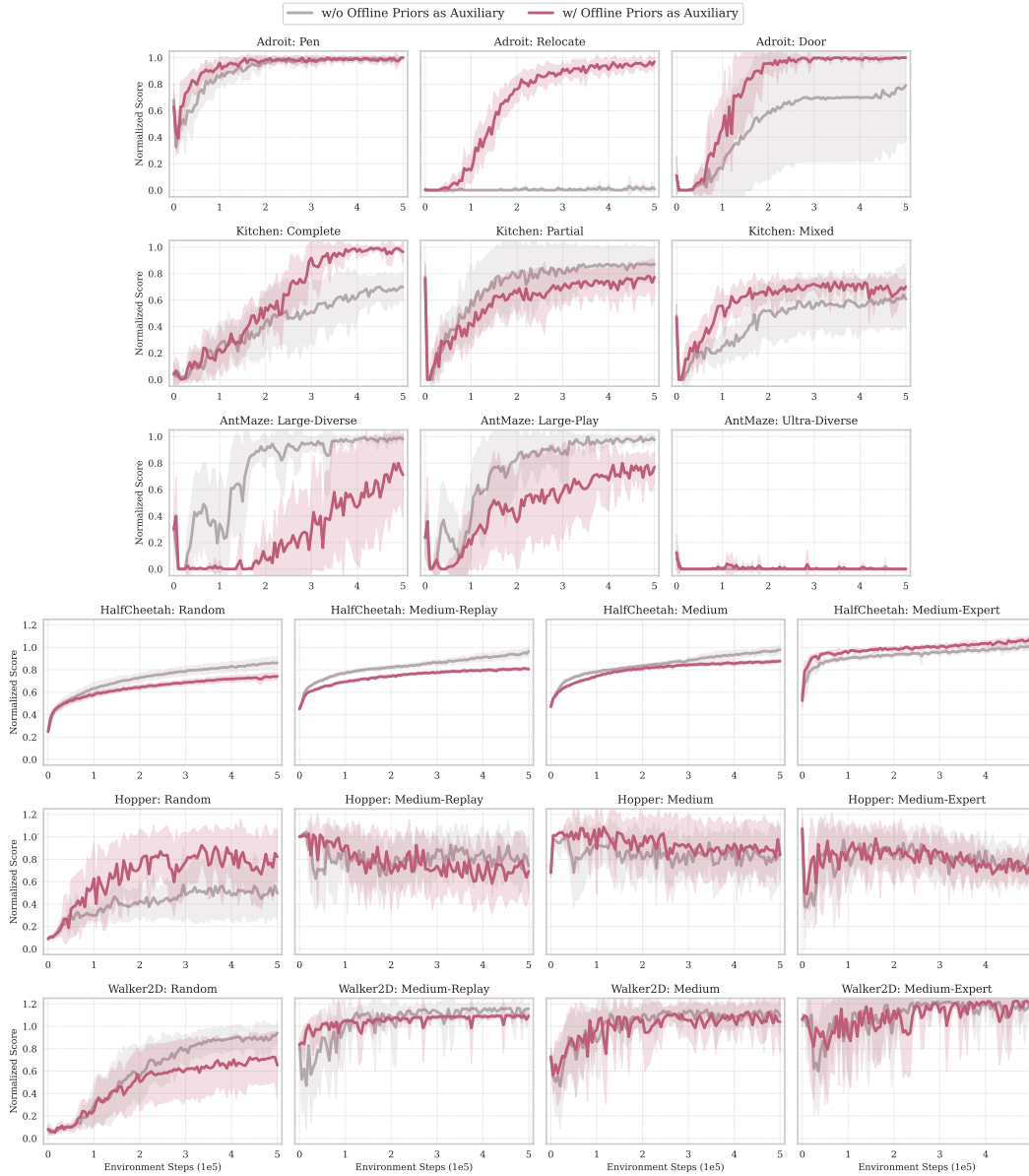


Figure 6: Effect of auxiliary reliance (β) without reference constraint. Both conditions use offline-trained initialization and standard SAC. Retaining offline data helps on some tasks but hurts on others.

645 **Auxiliary (β), with reference constraint.** This experiment repeats the auxiliary comparison with a
 646 reference constraint present: both conditions use Cal-QL’s conservative penalty during online learning. The
 647 comparison isolates whether retaining offline data remains beneficial when conservative value estimation is
 648 also active. Comparing Figures 6 and 7 reveals how the effect of auxiliary data depends on whether a reference
 649 constraint is in place, further confirming that the reliance parameters interact.

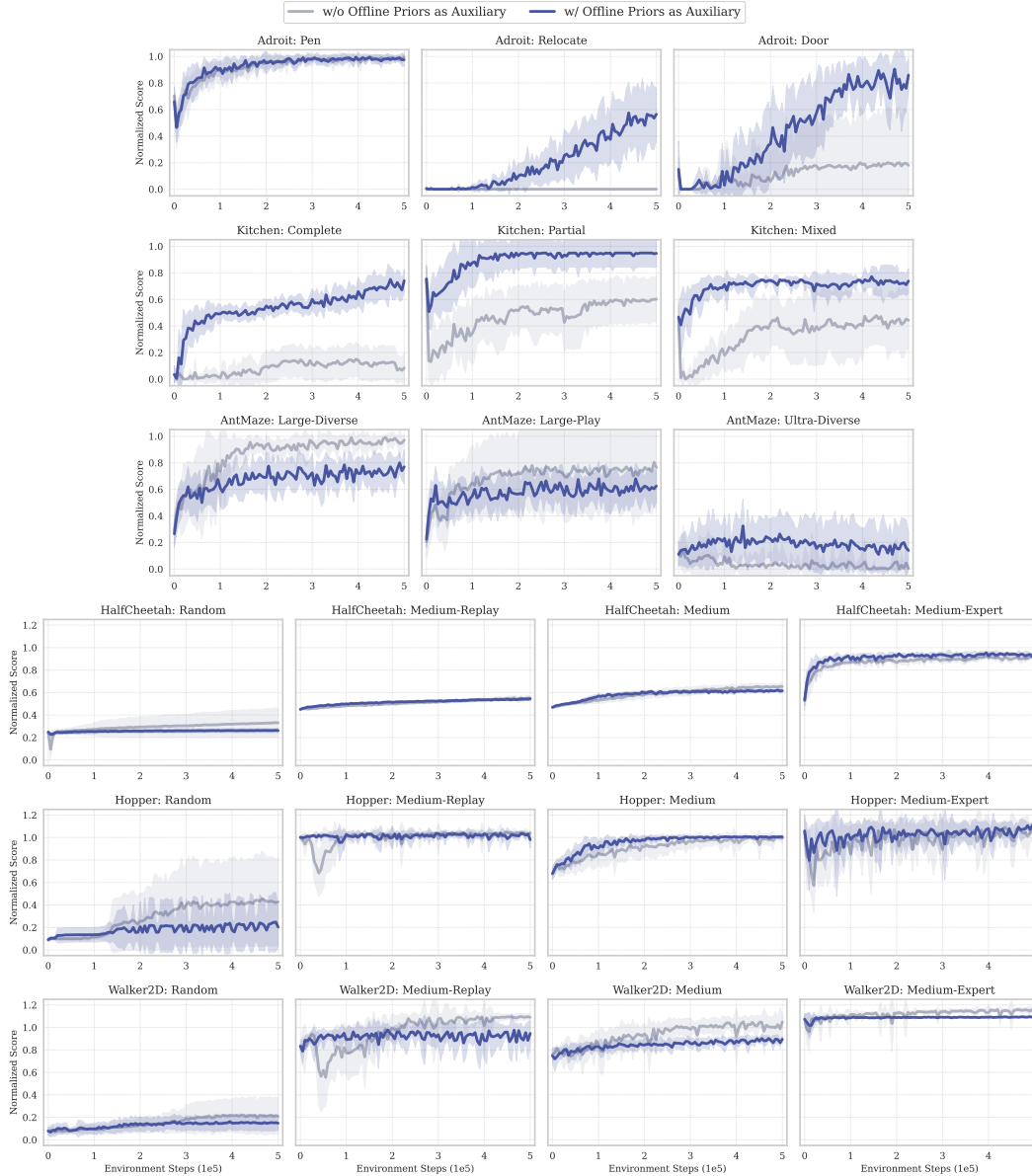


Figure 7: Effect of auxiliary reliance (β) with reference constraint ($\lambda > 0$). Both conditions use offline-trained initialization and Cal-QL’s conservative penalty. The effect of retaining offline data can differ from Figure 6, demonstrating interaction between reliance parameters.

650 **Summary.** Across all five experiments and all task domains, the same pattern emerges: toggling any single
 651 reliance parameter produces help-or-hurt reversals across tasks. No setting of μ , λ , or β is uniformly beneficial.
 652 The paired experiments further reveal that the effect of one parameter depends on the configuration of the others:
 653 the impact of a reference constraint differs depending on whether auxiliary data is present (Figures 4 vs 5), and
 654 the impact of auxiliary data differs depending on whether a reference constraint is active (Figures 6 vs 7). This
 655 confirms that the reliance parameters are not independently tunable and supports the argument in Section 4.1
 656 that tension management has no universal optimum.