



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 14    **Issue:** III    **Month of publication:** March 2026

**DOI:** <https://doi.org/10.22214/ijraset.2026.78190>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# Beyond Final Reward: A Deployment-Oriented Evaluation Framework for Online Hyperparameter Adaptation in Reinforcement Learning

Halder Albert David<sup>1</sup>, Yuhong Zhang<sup>2</sup>

<sup>1</sup>Research Scholar, <sup>2</sup>Associate Professor, College of Information Science and Engineering, School of Artificial Intelligence and Big Data, Henan University of Technology, Zhengzhou, Henan, China

**Abstract:** Online hyperparameter adaptation has become an important direction in reinforcement learning (RL) because fixed training configurations often fail to match the non-stationary dynamics of learning. Yet adaptive RL systems are still commonly judged by one dominant outcome: final episodic return. That practice is increasingly inadequate. Adaptive methods do not merely optimize a policy; they intervene in the training process itself and therefore behave as closed-loop regulators whose value depends not only on endpoint performance but also on stability, responsiveness, controller overhead, intervention quality, and deployment feasibility. This paper presents an analysis-based framework for evaluating online hyperparameter adaptation in RL without introducing new experiments. Drawing on published literature on meta-gradient reinforcement learning, AutoRL, hyperparameter optimization, population-based training, and self-tuning RL, we argue that adaptive RL should be assessed across seven complementary dimensions: adaptation effectiveness, stability and safety, computational overhead, intervention efficiency, information utility, information throughput, and deployment feasibility. We further provide a taxonomy of adaptive methods, a fair-comparison protocol for budget-matched evaluation, explicit formal definitions for several deployment-facing metrics, and a research agenda covering richer meta-state design, larger control spaces, warm-starting, transfer, and benchmark standardization. The resulting framework is intended as a practical reporting standard for future adaptive RL papers, especially those targeting single-agent or resource-constrained deployment settings.

**Keywords:** Adaptive reinforcement learning, AutoRL, online hyperparameter control, evaluation methodology, deployment feasibility, fair comparison, hyperparameter optimization

## I. INTRODUCTION

### A. Background and Motivation

Reinforcement learning has matured into a broad methodology for sequential decision making, with successful applications in control, robotics, recommendation, games, and resource management [3,5-7]. At the same time, the training of modern RL agents remains unusually sensitive to hyperparameters such as learning rate, discount factor, entropy coefficient, clipping range, update frequency, and exploration schedules [4,12-14,18]. This sensitivity is not a minor engineering inconvenience. In practical workflows it determines whether learning is stable or brittle, whether a model reaches useful performance quickly or plateaus early, and whether reported improvements survive across seeds and environments [4,13,18].

Classical hyperparameter optimization addresses this problem largely through inter-run search. Bayesian optimization, tree-structured Parzen estimators, Hyperband-style budget allocation, Optuna, and broader AutoML systems have improved sample-efficient search over static configurations [8-10,31-34]. However, these methods treat each training run as a mostly opaque object: a configuration is proposed, a run is executed, and the final score is returned. That design works well when the best configuration is approximately stationary over the entire run, but it is less well aligned with RL, where the data distribution, optimization landscape, and exploration requirements drift as the policy improves [12,13,27].

Within-run adaptivity is therefore increasingly attractive. Population-based training, meta-gradient methods, self-tuning actor-critic approaches, and decision-based controllers attempt to adjust hyperparameters as training unfolds [15-21]. Published work in meta-gradient reinforcement learning and AutoRL has made clear that hyperparameter adaptation should be studied as part of the learning process itself rather than as a purely offline selection problem [1,2]. Building on that literature, this paper argues that adaptive RL systems should be evaluated as training regulators, not only as reward maximizers, and formalizes deployment-oriented quantities such as information utility and adaptation information throughput for that purpose.

### B. Problem Statement

Despite rapid growth of AutoRL and online tuning methods, the field still lacks a shared evaluation standard for adaptive hyperparameter control. Most papers continue to emphasize final return, occasionally supplemented by sample efficiency or wall-clock speed, while giving limited and inconsistent attention to controller overhead, intervention frequency, stability envelopes, or deployment assumptions [4,11,14,18]. As a result, methods that achieve similar final scores can look indistinguishable even when one is cheap, stable, and reproducible and the other is fragile, expensive, or overly dependent on hidden infrastructure. The absence of a common framework also makes it difficult to compare single-agent methods fairly against population-based methods or to distinguish benefits caused by online adaptation from those inherited from stronger initializations [1,2,15-17].

$$\theta_{t+1} = \mathcal{A}(\theta_t, \mathcal{D}_t; \eta_t) \# (1)$$

Equation (1) formalizes the inner learning process. The policy or value parameters  $\theta_t$  are updated by the underlying RL algorithm  $\mathcal{A}$  using the current data  $\mathcal{D}_t$  and the current hyperparameter setting  $\eta_t$ . This makes explicit that the learner's trajectory depends not only on the environment data but also on the evolving control variables.

$$\eta_{t+1} = \Pi_{\Omega}(\eta_t + \pi_{\phi}(m_t)) \# (2)$$

Equation (2) defines the outer adaptation loop. The controller observes a compact meta-state  $m_t$ , produces an update through  $\pi_{\phi}$ , and projects the result back into an admissible set  $\Omega$ . The projection operator formalizes bounded actuation and makes the framework compatible with stability-aware or safety-aware adaptation.

### C. Scope and Positioning

This paper focuses on online hyperparameter adaptation in RL, meaning methods that modify aspects of the training configuration during a run on the basis of observed learning signals. The scope includes learned controllers, meta-gradients, decision-based adaptation, and resource-aware variants of population-based tuning. It excludes broad AutoML topics not specific to RL training dynamics, such as generic neural architecture search, dataset curation, or pure offline model selection. The paper is deliberately analysis-based: it proposes a conceptual framework, taxonomy, and protocol rather than new algorithmic benchmarks. The goal is to provide a strong methodological backbone for future experimental work.

### D. Contributions

The main contributions of this paper are as follows:

- It clarifies why final return alone is insufficient for evaluating adaptive RL systems and explains the hidden failure modes that remain invisible under endpoint-only reporting.
- It introduces a deployment-oriented evaluation framework with seven dimensions: adaptation effectiveness, stability and safety, computational overhead, intervention efficiency, information utility, information throughput, and deployment feasibility.
- It provides a taxonomy of adaptive hyperparameter methods organized by adaptation timing, resource model, update mechanism, actuation discipline, and implicit objective.
- It proposes a fair-comparison protocol that makes budget matching, initialization parity, statistical rigor, and reporting discipline explicit.
- It identifies open research problems, including richer meta-state design, larger action spaces, warm-starting, transfer-aware evaluation, and benchmark standardization.

### E. Paper Structure

The remainder of the paper follows the requested sequence for this framework paper. Section 2 explains why final reward is not enough, Section 3 presents the proposed framework, Section 4 develops the taxonomy, Section 5 proposes a fair-comparison protocol, Section 6 discusses open problems, Section 7 concludes, and a compact abstract is placed at the end of the document.

## II. WHY FINAL REWARD IS NOT ENOUGH

Final return is a necessary metric, but it is not a sufficient one for adaptive RL. The reason is structural. Static baselines and adaptive methods do not only differ in policy quality; they differ in how they shape the training trajectory that produces that quality. Two methods may converge to similar endpoint performance while having dramatically different profiles in variance, collapse risk, time-to-threshold, computational tax, and sensitivity to random seed [4,13,14,18]. In RL, where uncertainty across runs is often large and benchmark conclusions can change under small evaluation budgets, endpoint-only reporting systematically compresses information that matters for scientific and deployment decisions [4,13].

### A. Endpoint Performance Obscures Path Quality

Adaptive methods influence the path of training, not just its end state. A controller that stabilizes learning early may reach usable performance much sooner even if the final return at a fixed horizon is only modestly better. Conversely, an aggressive controller may achieve a high final score by inducing oscillations, repeated failures, or strong seed dependence along the way. For deployment-oriented research, these distinctions matter because users often care about time-to-threshold, reliability, and engineering effort rather than the final point on a single learning curve [1,4,13,18].

$$T_{\tau} = \min\{t: \bar{R}_t \geq \tau\} \quad (3)$$

Equation (3) defines time-to-threshold, which captures how quickly a method reaches a practically useful reward level  $\tau$ . Two methods may have similar final return but very different values of  $T_{\tau}$ , making this quantity especially relevant for deployment-oriented evaluation.

This observation is reinforced by reproducibility work in RL. Henderson et al. showed that seemingly minor implementation and evaluation choices can change conclusions materially [4]. Agarwal et al. further demonstrated that many benchmark claims are made with too few runs and with inadequate uncertainty accounting [13]. When adaptive controllers are added on top of already variable RL pipelines, these issues intensify rather than disappear. A comparison that reports only the best mean return therefore risks rewarding brittle methods and under-valuing robust ones.

### B. Adaptive Training Is a Control Problem

A second reason final reward is insufficient is that online hyperparameter adaptation is naturally a control problem. The base learner acts as the plant, training diagnostics act as observations, and the adaptive mechanism acts as a controller that changes the learning dynamics through bounded or unbounded actuation [1,19-23]. Once this perspective is taken seriously, standard control questions immediately arise: Is the system stable? How noisy are the observations? How delayed is the effect of an intervention? Does the controller overreact to transient disturbances? Is the actuation sparse, bounded, and interpretable, or constant and potentially destabilizing?

Final reward does not answer these questions. A method that reaches a strong endpoint while consuming frequent interventions and operating close to instability should not be considered equivalent to a method that reaches the same endpoint with modest, well-timed, bounded actions. For this reason, evaluation should reflect boundedness, timing, and selectivity of interventions rather than treating all adaptive controllers as interchangeable.

### C. Compute, Latency, and Controller Cost Are First-class Outcomes

Adaptive RL is often motivated as a way to improve learning under limited resources, yet many papers leave the resource side implicit. Population-based methods purchase adaptivity through parallel infrastructure [15-17], while gradient-based methods may incur expensive higher-order computations [19,32,35]. Even single-agent controllers add overhead through logging, feature aggregation, controller updates, and hyperparameter application. Without separating these costs, it is impossible to know whether a method improves the training process or merely shifts budget from one component to another.

Recent studies on RL tuning practice and AutoRL methodology highlight that runtime, controller latency, and resource budget strongly shape whether adaptive methods are viable in practice [14-18].

This paper extends that insight by framing the systems question as one of computational-information bottlenecks: speedup is only meaningful if the controller continues to produce signals that remain useful for subsequent improvement.

*D. Adaptive Methods Should Be Judged by the Quality of Intervention*

A further limitation of endpoint-only evaluation is that it ignores the quality of intervention itself. Adaptive methods differ in how often they intervene, how large the interventions are, how correlated they are with later improvement, and whether they act when the training regime genuinely changes. This dimension is particularly important for methods whose outer-loop policy is learned, because an always-on controller can appear active without actually being informative. To formalize this issue, this paper defines adaptation information utility,  $I(U; \Delta R)$ , and adaptation information throughput,  $I(U; \Delta R)/C_{ctrl}$ , as measures of whether controller outputs carry useful reward-relevant information and how efficiently that information is produced.

These notions are valuable because they shift attention from whether a controller changed a hyperparameter to whether the change was predictive, timely, and worthwhile. In other words, a controller should be judged on intervention efficiency and signal quality, not merely on intervention existence.

TABLE I  
WHY FINAL REWARD ALONE CAN BE MISLEADING IN ADAPTIVE RL EVALUATION

Observed endpoint outcome	Hidden failure mode	Why it matters	Recommended additional evidence
Similar final return	Very different convergence path	One method may reach usable reward much earlier	Time-to-threshold, learning-curve area, early-stage stability
High final mean	High run-to-run variance	Point estimates can hide fragility across seeds	Confidence intervals, IQM/median, seed count
Strong best-seed result	Frequent collapse events	Deployment cares about worst-case failure as well as best case	Failure rate, NaN/Inf count, collapse incidence
Improved reward	Large controller tax	Gains may disappear under wall-clock or single-device budgets	Component-wise timing and total overhead
Active controller	Low intervention quality	Frequent changes can be noisy rather than useful	Intervention frequency, information utility, ablation of controller actions

**III. A DEPLOYMENT-ORIENTED EVALUATION FRAMEWORK**

The central proposal of this paper is that online hyperparameter adaptation should be evaluated through a multi-dimensional framework tailored to closed-loop RL training. The framework is deployment-oriented in the sense that it asks whether an adaptive method is not only effective in principle, but also stable, efficient, and practical under realistic computational budgets. Figure 1 summarizes the logic: an adaptive RL system generates evidence streams during training, these streams are examined through several evaluation lenses, and the resulting assessment is converted into a reporting package that makes budget and stability assumptions explicit.

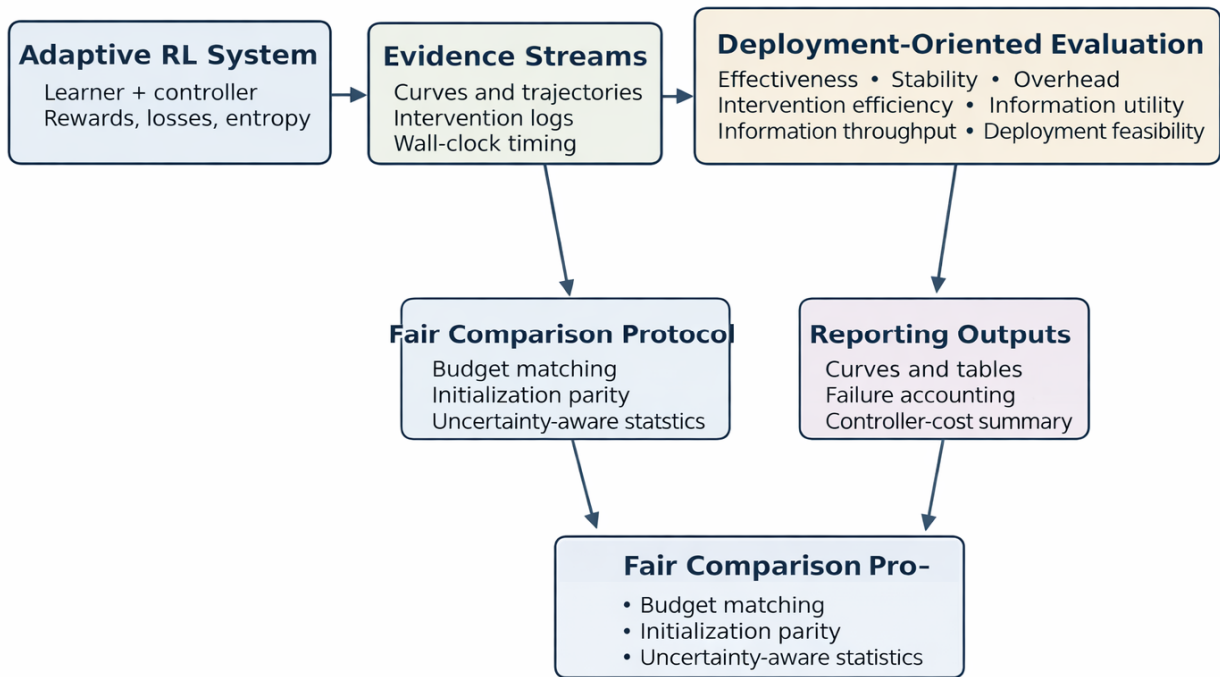


Fig. 1. Conceptual overview of the proposed evaluation pipeline for adaptive RL systems

### A. Design Principles

Four design principles organize the framework. First, evaluation must separate policy quality from regulator quality. An adaptive method can improve the policy while still being a poor regulator if it is unstable, expensive, or overly interventionist. Second, evaluation must distinguish task-time budgets from wall-clock budgets because adaptive methods often trade one against the other [1,2,15-17]. Third, controller usefulness should be measured by evidence, not assumed from activity: interventions should be connected to subsequent training gains. Fourth, the reporting standard must support fair cross-family comparison among static, population-based, gradient-based, and decision-based methods.

### B. Evaluation Dimensions

- 1) Adaptation effectiveness. This dimension asks whether online adaptation improves the learning process in substantive terms. Relevant outcomes include final return, time-to-threshold, performance consistency across seeds, robustness across environments, and retention under transfer [1,11,12,24]. The important point is that effectiveness should be assessed over the trajectory, not only at the endpoint.
- 2) Stability and safety. Adaptive methods intervene in already noisy RL pipelines, so stability is a primary outcome rather than a secondary diagnostic. Evaluation should report numerical failures, reward collapses, divergence triggers, safeguard activations, smoothness of parameter evolution, and whether the controller respects explicit actuation bounds [1,4,13,18].

$$\text{InstabRate}_\delta = \frac{1}{T} \sum_{t=1}^T \mathbf{1}! [\text{NaN}_t \vee \text{Inf}_t \vee |\bar{R}_t - \bar{R}_{t-1}| > \delta] \#(4)$$

Equation (4) turns stability into an explicit measurable quantity. It counts the proportion of training points associated with numerical failure or abrupt reward discontinuity. The threshold delta can be set after reward normalization so that the framework distinguishes smooth adaptive control from brittle or collapse-prone behavior.

- 3) Computational overhead. A controller that improves reward but doubles wall-clock time may be unattractive for practical deployment. Computational overhead should therefore be decomposed into base training, controller forward and backward cost, hyperparameter application, logging, and evaluation. This paper makes that decomposition explicit so that reported gains can be interpreted against the real cost of adaptation.

$$\rho_{oh} = \frac{C_{adaptive} - C_{base}}{C_{base}} \#(5)$$

Equation (5) defines relative controller overhead. It quantifies the extra computational tax introduced by online adaptation beyond the base learner alone. Reporting rho\_oh prevents papers from claiming adaptive gains without disclosing their operational cost.

- 4) Intervention efficiency. This dimension captures whether the controller intervenes economically. Reporting should include intervention frequency, average update magnitude, event-trigger ratio, and the fraction of interventions associated with later improvement. Sparse, well-targeted control is often more useful than constant reconfiguration.

$$IE_h = \frac{1}{N_{int}} \sum_{k=1}^{N_{int}} \max(0, \overline{R_{t_k+h}} - \overline{R_{t_k}}) \#(6)$$

Equation (6) measures the average positive return improvement associated with an intervention over a short horizon h. It therefore distinguishes sparse, useful interventions from dense but low-value controller activity. In a deployment-oriented setting, high intervention efficiency is preferable to constant but weak control.

- 5) Information utility. Information utility measures whether controller actions are informative about subsequent performance change. Even approximate proxies - such as mutual information, lagged correlation, or signed improvement association - are valuable because they separate meaningful control from noisy activity [2].

$$U_h = I(U_t; \Delta R_{t:t+h}), \quad \Delta R_{t:t+h} = \overline{R_{t+h}} - \overline{R_t} \#(7)$$

Equation (7) defines information utility as the mutual information between controller outputs and subsequent reward improvement over horizon h. A high value indicates that the adaptation signal is meaningfully related to future learning benefit rather than being arbitrary or noisy.

- 6) Information throughput. Throughput combines utility with elapsed time, asking how much reward-relevant control evidence is generated per unit wall-clock time. This metric is especially important when comparing methods with very different infrastructure assumptions or outer-loop costs [2].

$$\mathcal{J}_h = \frac{U_h}{C_{ctrl}} \#(8)$$

Equation (8) defines adaptation information throughput as the amount of reward-relevant control information obtained per unit control cost. This allows a method to be evaluated not only by whether it adapts, but by how efficiently it converts computation into informative adaptation.

- 7) Deployment feasibility. Finally, an adaptive method should be evaluated for its feasibility under the intended operating regime. Can it run on a single GPU? Does it require a population? Does it assume frequent checkpoints, long profiling phases, or specialized kernels? Does it remain understandable enough for engineering teams to trust it? These questions determine practical adoption even when benchmark scores look strong.

TABLE II  
PROPOSED EVALUATION DIMENSIONS AND REPRESENTATIVE METRICS

Dimension	Core question	Representative metrics	Typical reporting unit
Adaptation effectiveness	Does adaptation improve learning outcomes?	Final return; time-to-threshold; curve area; transfer retention	Mean and uncertainty across seeds
Stability and safety	Does the controller avoid destabilizing training?	NaN/Inf counts; collapse rate; safeguard activations; bounded-update violations	Counts or rates per run
Computational overhead	What controller tax is added?	Controller forward/backward time; total wall-clock overhead	Seconds and percent overhead
Intervention efficiency	Are interventions economical and timely?	Intervention frequency; update magnitude; event-trigger ratio	Per checkpoint or per episode
Information utility	Do controller actions carry useful signal?	Mutual information; lagged correlation; reward-improvement association	Per adaptation window
Information throughput	How fast is useful control information produced?	Utility divided by wall-clock time	Bits/time or normalized score/time
Deployment feasibility	Is the method realistic for the target setting?	Hardware requirement; population size; profiling cost; fallback dependence	Narrative plus resource table

### C. Reference Reporting Template

A strong adaptive RL paper should ideally report a compact package of evidence: (i) learning curves with confidence intervals over multiple seeds, (ii) time-to-threshold and/or performance-over-time summaries, (iii) controller overhead broken down by component, (iv) explicit stability accounting, (v) intervention statistics, and (vi) a short deployment note describing hardware assumptions and initialization procedures. This template is not intended to overburden authors. Rather, it is meant to prevent a common failure mode in which a method is described as adaptive and practical while crucial information about controller cost or reliability is missing.

### D. Relationship to Adjacent Published Literature

The framework is not invented in a vacuum. Meta-gradient reinforcement learning, self-tuning actor-critic approaches, population-based training, and broader AutoRL surveys all indicate that adaptive methods should be compared on more than final reward [1,2,19,20,26]. The present paper synthesizes these strands into a broader reporting standard and adds explicit definitions for information utility and adaptation information throughput as deployment-oriented quantities.

## IV. TAXONOMY OF ADAPTIVE HYPERPARAMETER METHODS

Evaluation is easier when the space of methods is organized clearly. The literature can be taxonomized along five complementary axes: when adaptation occurs, what resource model it assumes, how updates are generated, what safety discipline constrains actuation, and what objective the method implicitly optimizes. These axes are intentionally orthogonal. A single method can be within-run, single-agent, gradient-based, bounded, and throughput-aware, while another can be within-run, population-based, unconstrained, and reward-only.

### A. Adaptation Timing: Inter-run Versus Within-run

Inter-run methods search over static configurations across repeated trials and then hold the chosen configuration fixed during each run [8-10,31-34]. Within-run methods modify training parameters on the fly, responding to regime shifts and transient instability [1,11,15-21].

The first family is often easier to analyze, but the second family is more faithful to the non-stationary character of RL.

*B. Resource Model: Single-agent Versus Population-based*

Resource assumptions are central rather than incidental. Single-agent methods preserve the structure of one training run and are therefore attractive in resource-constrained settings [1,2,12,24]. Population-based methods achieve exploration and adaptivity by maintaining multiple agents in parallel, periodically copying and mutating policies and hyperparameters [15-17]. These methods can be effective, but their wall-clock and infrastructure requirements change the comparison regime fundamentally.

*C. Update Mechanism: Schedules, Gradients, Decisions, and Hybrids*

A third axis concerns how updates are generated. Fixed schedules encode a time-indexed plan chosen in advance. Meta-gradient and bilevel approaches differentiate through learning dynamics to optimize hyperparameters [19,20,35,36]. Decision-based methods learn or design a policy over updates using bandits, RL controllers, or heuristics [1,18,21]. Hybrid systems combine strong offline initialization with lightweight online refinement, a pattern increasingly common in practice [1,8,10,12].

*D. Actuation Discipline: Unconstrained Versus Bounded Control*

Adaptive methods also differ in how tightly they constrain interventions. Some permit effectively unconstrained continuous updates; others enforce clipping, projection, smooth switching, sparse checkpoints, or dynamic fallback [1,2]. This axis matters because safe deployment depends not only on what a method can do, but on what it is prevented from doing when training becomes noisy or unstable.

*E. Implicit Objective: Reward, Stability, Compute, or Information*

Finally, methods should be classified by their real objective. Many papers optimize reward or sample efficiency alone. Others explicitly target reproducibility, robustness, or wall-clock cost [1,2,13,14]. Recent systems work suggests an additional category: information-aware methods that preserve or maximize the usefulness of control signals under compute constraints [2]. Stating the objective clearly prevents mismatched comparisons, such as criticizing a low-overhead single-agent controller for not matching the raw final score of a compute-intensive population method while ignoring the budget difference.

TABLE III  
TAXONOMY OF ADAPTIVE HYPERPARAMETER METHODS IN RL

Method family	Adaptation timing	Resource model	Update mechanism	Typical strengths / blind spots
Offline HPO / AutoML	Inter-run	Single run per trial, many repeated trials	Bayesian optimization, TPE, Hyperband, Optuna	Strong initialization; weak responsiveness to within-run shifts
Fixed schedules	Within-run, predetermined	Single-agent	Hand-designed time-based rules	Cheap and stable; cannot react to atypical dynamics
Population-based training	Within-run	Population parallelism	Copy-mutate-exploit/explore	Strong adaptivity; high infrastructure cost
Meta-gradient / bilevel	Within-run	Single-agent or limited parallelism	Differentiation through training dynamics	Principled updates; sensitive to gradient noise and differentiability assumptions
Decision-based controllers	Within-run	Usually single-agent	Bandits, RL controllers, recurrent meta-policies	Context-sensitive updates; quality depends on observability and safeguards
Systems-aware adaptive methods	Within-run	Single-agent or specialized runtime	Profiling, scheduling, semantics-preserving acceleration	Deployment-aware efficiency; harder to benchmark with reward alone

## V. FAIR COMPARISON PROTOCOL FOR ADAPTIVE RL PAPERS

A framework is only useful if it leads to better comparisons. This section converts the earlier analysis into a practical protocol for benchmarking adaptive RL methods. The protocol is especially important for comparisons across families because different methods consume different types of budget, rely on different levels of prior knowledge, and expose different kinds of risk.

### A. Match Interaction Budgets and Account for Infrastructure

Any fair comparison must make the total training budget explicit. Environment interactions, number of actors, hardware parallelism, and wall-clock budget should all be reported. A single-agent controller operating on one device should not be compared to a population method as though both consumed the same budget merely because they were run for the same number of environment steps [1,2,15-17].

### B. Equalize Initialization and Prior Information

Methods should also be compared under comparable initialization quality. If one adaptive method starts from a carefully tuned TPE or Optuna configuration while another begins from an arbitrary default, the resulting difference confounds initialization benefit with online adaptation benefit. Similar discipline should become standard, because otherwise comparisons entangle warm-start quality with the online controller itself.

### C. Report Timestep and Wall-clock Views Together

Adaptive methods often improve learning in one accounting system while looking weaker in another. For example, a controller may yield better reward per environment step but incur large supervisory overhead, or conversely reduce wall-clock time without changing sample efficiency [2]. A strong paper should therefore present at least one timestep-based and one wall-clock-based comparison whenever the controller incurs non-trivial cost.

### D. Separate Learning Benefit from Controller Tax

Authors should decompose the total training pipeline into base learner computation, controller update cost, hyperparameter-application cost, and logging/evaluation overhead. This enables readers to see whether gains come from improved learning dynamics, implementation optimizations, or hidden costs transferred elsewhere in the system [2].

### E. Use Uncertainty-aware Statistics

Because RL results are noisy, adaptive RL papers should report multiple seeds, uncertainty intervals, and robust aggregate statistics such as medians or interquartile-based summaries when appropriate [4,13,18]. Significance testing should be used cautiously and with correction for multiple comparisons when many tasks or metrics are analyzed [2,13]. A controller that wins on a single mean number but lacks uncertainty accounting should not be treated as decisively better.

### F. Include a Deployment Note

Finally, every adaptive RL paper should include a short deployment note describing the intended operating regime: single-GPU training, edge device, offline cluster search, real-time control, or another environment. This note should state why the chosen evaluation metrics match the deployment scenario. Such transparency reduces the risk of mismatched criticism and makes the contribution easier to interpret.

TABLE IV  
MINIMUM FAIR-COMPARISON CHECKLIST FOR ADAPTIVE RL PAPERS

Checklist item	What should be reported	Why it matters
Budget parity	Environment steps, actors, GPUs/CPUs, wall-clock budget	Prevents hidden infrastructure advantages
Initialization parity	Common warm-start or explicit explanation of differences	Separates offline tuning benefit from online adaptation
Seed reporting	Number of seeds and uncertainty intervals	Reduces over-interpretation of noisy

		averages
Controller tax	Component-wise runtime overhead	Shows the real cost of adaptation
Stability accounting	NaN/Inf, collapse events, safeguard activations	Adaptive systems must be safe as well as effective
Intervention accounting	How often and how strongly the controller intervenes	Distinguishes meaningful control from noisy activity
Deployment note	Target operating regime and assumptions	Aligns evaluation with intended use case

## VI. OPEN PROBLEMS AND RESEARCH AGENDA

A shared evaluation framework should do more than standardize reporting; it should clarify where the field needs better ideas. Several open problems follow naturally from the analysis above and from gaps repeatedly noted across the AutoRL, meta-gradient, and RL tuning literatures [1,2,14,16].

### A. Richer Meta-state and Observability

Many current controllers operate on compact summaries such as recent reward, loss, gradient norm, or entropy. This keeps overhead low, but may leave the controller partially blind in high-dimensional or strongly non-stationary settings [1]. Future work should explore richer representations of training phase, uncertainty, advantage distributions, and non-stationarity while keeping the outer loop lightweight enough for deployment.

### B. Larger Control Spaces with Safe Actuation

Most methods adapt only a small number of hyperparameters. Extending control to clipping ranges, target-network rates, replay parameters, optimizer internals, or architecture-level knobs could make adaptation more powerful, but it also increases the risk of instability [1,19,20]. Research is needed on scalable safety mechanisms, especially for mixed discrete-continuous control spaces.

### C. Warm-starting and Amortized Initialization

Even methods designed for single-agent settings often incur an upfront initialization cost. Warm-start or amortized strategies would improve repeated deployment scenarios by reducing the cost of re-tuning across related tasks. This is an important methodological gap because many practical teams retrain related agents repeatedly rather than once.

### D. Transfer-aware and Continual Evaluation

A controller that is effective on one benchmark but fails to transfer across tasks may have limited practical value. Future work should evaluate whether adaptation policies are reusable across related environments, reward structures, or hardware regimes [1,21-24]. Transfer-aware reporting could become a major differentiator between hand-crafted schedules and learned controllers.

### E. Better Measures of Controller Usefulness

Information utility and throughput offer an appealing beginning, but more work is needed on causal attribution. When a controller changes a hyperparameter, how much later improvement can legitimately be attributed to that action rather than to background learning momentum? Methods from lagged analysis, counterfactual evaluation, and causal time-series modeling may help answer this more rigorously [2].

### F. Benchmark Standardization

The field would benefit from benchmark suites and reporting cards designed specifically for adaptive RL. Such benchmarks should include settings with varying non-stationarity, sparse reward, wall-clock constraints, and resource asymmetry. Without standardization, adaptive RL papers will continue to mix incomparable assumptions about infrastructure, initialization, and controller cost.

TABLE V  
OPEN PROBLEMS AND THE KIND OF EVIDENCE FUTURE WORK SHOULD PROVIDE

Open problem	Why it matters	Useful future evidence
Richer meta-state design	Better observability may improve controller timing and safety	Ablations on signal sets and overhead trade-offs
Expanded control spaces	More knobs may improve performance but increase instability	Safety analysis and bounded-update studies
Warm-starting / amortization	Repeated deployment depends on total time-to-solution	Cross-run reuse experiments and initialization accounting
Transfer-aware control	Reusable controllers are more valuable than task-specific heuristics	Cross-task retention and domain-shift studies
Controller usefulness metrics	Need to separate meaningful from noisy intervention	Lagged association, counterfactual analyses, utility metrics
Standardized reporting	Comparability remains weak without common benchmarks	Benchmark cards and agreed minimum reporting sets

## VII. CONCLUSION

This paper argued that online hyperparameter adaptation in reinforcement learning should be evaluated as a closed-loop training regulation problem rather than as an endpoint-only optimization problem. Final return remains important, but it captures only one part of the value - and risk - of adaptive systems. By synthesizing published work in meta-gradient RL, AutoRL, HPO, population-based, and meta-learning literatures, we proposed a deployment-oriented framework built around seven evaluation dimensions, a taxonomy of method families, and a fair-comparison protocol for future studies. The broader message is methodological: adaptive RL papers should make controller cost, stability, intervention quality, and deployment assumptions explicit. Doing so will improve scientific rigor, make cross-family comparisons more honest, and increase the practical value of research on adaptive RL systems.

## REFERENCES

- [1] Z. Xu, H. P. van Hasselt, and D. Silver, "Meta-Gradient Reinforcement Learning," in *Advances in Neural Information Processing Systems 31 (NeurIPS)*, 2018.
- [2] J. Parker-Holder, R. Rajan, X. Song, A. Biedenkapp, Y. Miao, T. Eimer, B. Zhang, V. Nguyen, R. Calandra, A. Faust, F. Hutter, and M. Lindauer, "Automated Reinforcement Learning (AutoRL): A Survey and Open Problems," *Journal of Artificial Intelligence Research*, vol. 74, pp. 517-568, 2022.
- [3] X. Wang, S. Wang, X. Liang, D. Zhao, J. Huang, X. Xu, B. Dai, and Q. Miao, "Deep reinforcement learning: A survey," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 4, pp. 5064-5078, 2022.
- [4] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep reinforcement learning that matters," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, 2018.
- [5] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. Cambridge, MA: MIT Press, 2018.
- [6] V. Mnih et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529-533, 2015.
- [7] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," arXiv:1707.06347, 2017.
- [8] J. Bergstra and Y. Bengio, "Random search for hyper-parameter optimization," *Journal of Machine Learning Research*, vol. 13, no. 10, pp. 281-305, 2012.
- [9] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kegl, "Algorithms for hyper-parameter optimization," in *Advances in Neural Information Processing Systems*, vol. 24, 2011.
- [10] M. Feurer and F. Hutter, "Hyperparameter optimization," in *Automated Machine Learning: Methods, Systems, Challenges*. Cham: Springer, pp. 3-33, 2019.
- [11] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, "Hyperband: A novel bandit-based approach to hyperparameter optimization," *Journal of Machine Learning Research*, vol. 18, no. 185, pp. 1-52, 2018.
- [12] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 2623-2631, 2019.
- [13] B. Bischl et al., "Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 13, no. 2, e1484, 2023.
- [14] T. Eimer, M. Lindauer, and R. Raileanu, "Hyperparameters in reinforcement learning and how to tune them," in *International Conference on Machine Learning*, pp. 9104-9149, 2023.

- [15] R. Agarwal, M. Schwarzer, P. S. Castro, A. Courville, and M. G. Bellemare, "Deep reinforcement learning at the edge of the statistical precipice," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [16] J. Parker-Holder et al., "Automated reinforcement learning (AutoRL): A survey and open problems," *Journal of Artificial Intelligence Research*, vol. 74, pp. 517-568, 2022.
- [17] L. Franke, D. K. I. Weidele, N. Dehmamy, L. Ning, and D. Haehn, "AutoRL X: Automated reinforcement learning on the web," *ACM Transactions on Interactive Intelligent Systems*, vol. 14, no. 4, pp. 1-30, 2024.
- [18] J. Adkins, M. Bowling, and A. White, "A method for evaluating hyperparameter sensitivity in reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 37, 2024.
- [19] T. Zahavy, Z. Xu, V. Veeriah, M. Hessel, J. Oh, H. P. van Hasselt, D. Silver, and S. Singh, "A self-tuning actor-critic algorithm," *Advances in Neural Information Processing Systems*, vol. 33, pp. 20913-20924, 2020.
- [20] S. Flennerhag, T. Zahavy, B. O'Donoghue, H. P. van Hasselt, A. Gyorgy, and S. Singh, "Optimistic meta-gradients," *Advances in Neural Information Processing Systems*, vol. 36, 2023.
- [21] D. Maclaurin, D. Duvenaud, and R. P. Adams, "Gradient-based hyperparameter optimization through reversible learning," in *Proceedings of the 32nd International Conference on Machine Learning*, 2015.
- [22] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil, "Bilevel programming for hyperparameter optimization and meta-learning," in *Proceedings of the 35th International Conference on Machine Learning*, pp. 1568-1577, 2018.
- [23] T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey, "Meta-learning in neural networks: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 9, pp. 5149-5169, 2021.
- [24] A. Vettoruzzo, M.-R. Bouguelia, J. Vanschoren, T. Rognvaldsson, and K. Santosh, "Advances and challenges in meta-learning: A technical review," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 7, pp. 4763-4779, 2024.
- [25] J. Vanschoren, "Meta-learning: A survey," arXiv:1810.03548, 2018.
- [26] M. Jaderberg et al., "Population based training of neural networks," arXiv:1711.09846, 2017.
- [27] H. Bai and R. Cheng, "Generalized population-based training for hyperparameter optimization in reinforcement learning," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 8, no. 5, pp. 3450-3462, 2024.
- [28] A. Dushatskiy, A. Chebykin, T. Alderliesten, and P. Bosman, "Multi-objective population based training," in *International Conference on Machine Learning*, pp. 8969-8989, 2023.
- [29] J. Albrechts, H. M. Martin, and M. Tavakol, "Model-based meta-reinforcement learning for hyperparameter optimization," in *International Conference on Intelligent Data Engineering and Automated Learning*, pp. 27-39, 2024.
- [30] Z. Yang and L. Ma, "Adaptive step size rules for stochastic optimization in large-scale learning," *Statistics and Computing*, vol. 33, no. 2, article 45, 2023.
- [31] H. Pourshamsaei and A. Nobakhti, "Predictive reinforcement learning in non-stationary environments using weighted mixture policy," *Applied Soft Computing*, vol. 153, article 111305, 2024.
- [32] S. Li, S. Su, and X. Lin, "Optimizing the hyper-parameters of deep reinforcement learning for building control," *Building Simulation*, vol. 18, no. 4, pp. 765-789, 2025.
- [33] M. Yuan, B. Li, X. Jin, and W. Zeng, "Ultho: Ultra-lightweight yet efficient hyperparameter optimization in deep reinforcement learning," arXiv:2503.06101, 2025.
- [34] F. Zhao, F. Ji, T. Xu, and N. Zhu et al., "Hierarchical parallel search with automatic parameter configuration for particle swarm optimization," *Applied Soft Computing*, vol. 151, article 111126, 2024.
- [35] J. Lemobayo et al., "Hyperparameter tuning in machine learning: A comprehensive review," *Journal of Engineering Research and Reports*, vol. 26, no. 6, pp. 388-395, 2024.
- [36] L. Acerbi and W. J. Ma, "Practical Bayesian optimization for model fitting with Bayesian adaptive direct search," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [37] E. Ntentos, S. J. Warnett, and U. Zdun, "Supporting architectural decision making on training strategies in reinforcement learning architectures," in *2024 IEEE 21st International Conference on Software Architecture*, pp. 90-100, 2024.
- [38] R. Zhao et al., "Maximum entropy population-based training for zero-shot human-AI coordination," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, pp. 6145-6153, 2023.
- [39] S. Guan et al., "Adaptive multi-agent HVAC control for thermal comfort using multi-agent PPO with population-based training," *Energy and Buildings*, article 116882, 2025.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)