# State-of-the-Art Methods in Reinforcement Learning

AI Research Team

May 15, 2025

**Abstract**

This article explores cutting-edge methods in reinforcement learning (RL) that represent the current state of the art. We examine evolutionary approaches to RL algorithm design, meta-reinforcement learning, offline RL, and multi-agent systems. For each method, we provide a mathematical foundation, algorithmic details, and discuss practical applications. The article highlights how these advanced techniques address fundamental challenges in reinforcement learning, such as sample efficiency, generalization, and scalability. We also discuss emerging trends and future directions in the field, offering insights into how reinforcement learning continues to evolve and expand its capabilities.

## 1 Introduction

Reinforcement learning has witnessed remarkable advancements in recent years, with novel algorithms and approaches pushing the boundaries of what AI systems can achieve. As noted by [1], "The field of reinforcement learning is evolving rapidly, with new algorithms and techniques being developed to address increasingly complex challenges."

This article explores state-of-the-art methods in reinforcement learning that represent the cutting edge of the field. We focus on approaches that address key challenges such as sample efficiency, generalization, and scalability, which have traditionally limited the applicability of RL to real-world problems.

The methods discussed in this article include evolutionary approaches to RL algorithm design, meta-reinforcement learning, offline RL, and multi-agent systems. For each method, we provide a mathematical foundation, algorithmic details, and discuss practical applications.

# 2 Evolutionary Approaches to RL Algorithm Design

## 2.1 Automated Algorithm Design

Evolutionary approaches to reinforcement learning algorithm design leverage genetic algorithms and other evolutionary computation techniques to automatically discover effective RL algorithms. As described by [1], this approach "uses evolutionary algorithms to optimize the structure and parameters of reinforcement learning algorithms, potentially discovering novel approaches that outperform hand-designed methods."

### 2.1.1 PBT: Population-Based Training

Population-Based Training (PBT) combines evolutionary optimization with traditional RL training. It maintains a population of agents with different hyperparameters and periodically replaces poorly performing agents with modified copies of better-performing ones.

The PBT algorithm can be formalized as follows:

---

1: Initialize population $P = \{\theta_1, \theta_2, \ldots, \theta_n\}$ with random hyperparameters
2: **while** not converged **do**
3:    **for** each $\theta_i \in P$ in parallel **do**
4:       Train $\theta_i$ for $k$ steps using current hyperparameters
5:       Evaluate performance $f(\theta_i)$
6:    **end for**
7:    **for** each $\theta_i \in P$ **do**
8:       **if** $f(\theta_i)$ is in bottom $q\%$ of population **then**
9:          Select $\theta_j$ from top $p\%$ of population
10:          $\theta_i \leftarrow \theta_j$ (copy weights)
11:          Perturb hyperparameters of $\theta_i$
12:       **end if**
13:    **end for**
14: **end while**
15: **return** Best $\theta_i$ from population $P$

---

### 2.1.2 Evolved Policy Gradients

Evolved Policy Gradients (EPG) use evolutionary algorithms to discover novel policy gradient update rules. Instead of using hand-designed update rules like those in PPO or TRPO, EPG evolves the update rule itself.

The update rule in EPG can be represented as a neural network $U_\phi$ that maps from the current policy parameters $\theta$, states $s$, actions $a$, and rewards $r$ to parameter updates $\Delta\theta$:

$$\Delta\theta = U_\phi(\theta, s, a, r) \tag{1}$$

The parameters $\phi$ of the update rule are optimized using evolutionary strategies to maximize the expected return of the resulting policy.

## 2.2 Evolutionary Reinforcement Learning (ERL)

Evolutionary Reinforcement Learning (ERL) combines gradient-based RL methods with evolutionary algorithms to leverage the strengths of both approaches. ERL maintains a population of policies that are optimized using both gradient-based updates and evolutionary selection.

The ERL framework can be described as follows:

---

1: Initialize population of policies $P = \{\pi_1, \pi_2, \ldots, \pi_n\}$
2: Initialize replay buffer $D$
3: **while** not converged **do**
4:     **for** each $\pi_i \in P$ **do**
5:         Collect experiences by executing $\pi_i$ in the environment
6:         Add experiences to replay buffer $D$
7:         Evaluate fitness $f(\pi_i)$ based on accumulated rewards
8:     **end for**
9:     // Gradient-based updates
10:     **for** selected policies $\pi_i \in P$ **do**
11:         Sample mini-batch from $D$
12:         Update $\pi_i$ using gradient-based RL algorithm (e.g., DQN, DDPG)
13:     **end for**
14:     // Evolutionary updates
15:     Select elite policies based on fitness
16:     Generate new policies through crossover and mutation of elite policies
17:     Replace worst-performing policies with new policies
18: **end while**
19: **return** Best policy from population $P$

---

# 3 Meta-Reinforcement Learning

## 3.1 Learning to Learn

Meta-reinforcement learning aims to develop algorithms that can quickly adapt to new tasks by leveraging experience from related tasks. The key idea is to learn a learning algorithm rather than a specific policy for a single task.

### 3.1.1 Model-Agnostic Meta-Learning (MAML)

Model-Agnostic Meta-Learning (MAML) is a meta-learning algorithm that can be applied to reinforcement learning. MAML learns an initialization of policy parameters that can be quickly adapted to new tasks with a few gradient steps.

The MAML objective for RL can be formulated as:

$$\max_{\theta} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T})} \left[ \mathbb{E}_{\pi_{\theta'}} \left[ \sum_{t=0}^{H} \gamma^t r_t \right] \right] \tag{2}$$

where $\theta'$ is the adapted parameters for task $\mathcal{T}$, computed as:

$$\theta' = \theta + \alpha \nabla_{\theta} \mathbb{E}_{\pi_{\theta}} \left[ \sum_{t=0}^{H} \gamma^t r_t \right] \tag{3}$$

The MAML algorithm for RL can be described as follows:

---
1: Initialize parameters $\theta$
2: **while** not converged **do**
3:     Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:     **for** each task $\mathcal{T}_i$ **do**
5:         Sample trajectories $\tau_i$ using policy $\pi_{\theta}$ in task $\mathcal{T}_i$
6:         Compute adapted parameters: $\theta_i' = \theta + \alpha \nabla_{\theta} J_{\mathcal{T}_i}(\pi_{\theta})$
7:         Sample new trajectories $\tau_i'$ using policy $\pi_{\theta_i'}$ in task $\mathcal{T}_i$
8:     **end for**
9:     Update $\theta \leftarrow \theta + \beta \nabla_{\theta} \sum_i J_{\mathcal{T}_i}(\pi_{\theta_i'})$
10: **end while**
11: **return** $\theta$
---

### 3.1.2 Recurrent Meta-Reinforcement Learning

Recurrent Meta-Reinforcement Learning uses recurrent neural networks to implicitly perform meta-learning. The recurrent network learns to adapt its hidden state based on the history of states, actions, and rewards, effectively implementing an adaptive policy.

The policy in recurrent meta-RL can be formulated as:

$$\pi_{\theta}(a_t|s_t, h_t) = p(a_t|s_t, h_t; \theta) \tag{4}$$

where $h_t$ is the hidden state of the recurrent network, updated as:

$$h_t = f_{\theta}(h_{t-1}, s_{t-1}, a_{t-1}, r_t, s_t) \tag{5}$$

The hidden state $h_t$ implicitly encodes information about the current task, allowing the policy to adapt its behavior accordingly.

## 3.2 Context-Based Meta-RL

Context-based meta-RL approaches use a context vector to encode task-specific information. The policy conditions on this context vector to adapt its behavior to the current task.

### 3.2.1 Variational Context-Conditioned Policies

Variational Context-Conditioned Policies (VCCP) use a variational inference approach to learn a context encoder that maps task-specific trajectories to a context vector. The policy then conditions on this context vector to generate actions.

The VCCP objective can be formulated as:

$$\max_{\theta,\phi} \mathbb{E}_{\mathcal{T}\sim p(\mathcal{T})} \left[ \mathbb{E}_{c\sim q_\phi(c|\tau_{\mathcal{T}})} \left[ \mathbb{E}_{\pi_\theta(a|s,c)} \left[ \sum_{t=0}^{H} \gamma^t r_t \right] \right] - \beta D_{KL}(q_\phi(c|\tau_{\mathcal{T}})||p(c)) \right]$$

(6)

where $q_\phi(c|\tau_{\mathcal{T}})$ is the context encoder that maps trajectories $\tau_{\mathcal{T}}$ from task $\mathcal{T}$ to context vector $c$, and $p(c)$ is a prior over context vectors.

# 4 Offline Reinforcement Learning

## 4.1 Learning from Static Datasets

Offline reinforcement learning, also known as batch RL, aims to learn optimal policies from static datasets of experiences without additional environment interaction. This is particularly valuable in domains where online interaction is costly, dangerous, or impractical.

### 4.1.1 Conservative Q-Learning (CQL)

Conservative Q-Learning (CQL) addresses the challenge of overestimation in offline RL by learning a conservative Q-function that lower-bounds the true Q-value.

The CQL objective can be formulated as:

$$\min_\theta \alpha\mathbb{E}_{s\sim D,a\sim\mu(a|s)}[Q_\theta(s,a)] - \alpha\mathbb{E}_{s\sim D,a\sim\pi_\beta(a|s)}[Q_\theta(s,a)] + \frac{1}{2}\mathbb{E}_{(s,a,r,s')\sim D}[(Q_\theta(s,a) - \hat{\mathcal{B}}^\pi Q_{\bar\theta}(s,a))^2]$$

(7)

where $D$ is the offline dataset, $\mu(a|s)$ is the current policy, $\pi_\beta(a|s)$ is the behavior policy that generated the dataset, and $\hat{\mathcal{B}}^\pi$ is the empirical Bellman operator.

### 4.1.2 Batch-Constrained deep Q-learning (BCQ)

Batch-Constrained deep Q-learning (BCQ) addresses the extrapolation error in offline RL by constraining the learned policy to stay close to the behavior policy that generated the dataset.

BCQ uses a generative model $G_\omega(s)$ to model the behavior policy and selects actions by maximizing the Q-function over a set of actions sampled from this model:

$$\pi(s) = \arg \max_{a \in \{a_i \sim G_\omega(s)\}_{i=1}^n} Q_\theta(s, a) \tag{8}$$

The BCQ algorithm can be described as follows:

---

1: Initialize Q-networks $Q_{\theta_1}, Q_{\theta_2}$, target networks $Q_{\bar{\theta}_1}, Q_{\bar{\theta}_2}$, and generative model $G_\omega$
2: Train generative model $G_\omega$ on offline dataset $D$ to model behavior policy
3: **while** not converged **do**
4:     Sample mini-batch $(s, a, r, s') \sim D$
5:     Sample $n$ actions $\{a_i \sim G_\omega(s')\}_{i=1}^n$
6:     Compute target: $y = r + \gamma \min_{j=1,2} Q_{\bar{\theta}_j}(s', \arg \max_{a_i} Q_{\bar{\theta}_1}(s', a_i))$
7:     Update Q-networks to minimize $(Q_{\theta_j}(s, a) - y)^2$ for $j = 1, 2$
8:     Update target networks: $\bar{\theta}_j \leftarrow \tau \theta_j + (1 - \tau)\bar{\theta}_j$ for $j = 1, 2$
9: **end while**
10: **return** Policy $\pi(s) = \arg \max_{a \in \{a_i \sim G_\omega(s)\}_{i=1}^n} Q_{\theta_1}(s, a)$

---

## 4.2 Offline-to-Online Reinforcement Learning

Offline-to-online RL aims to leverage offline data to bootstrap online learning, reducing the amount of online interaction needed to learn effective policies.

### 4.2.1 Advantage-Weighted Regression (AWR)

Advantage-Weighted Regression (AWR) is a simple algorithm that can be applied in both offline and online settings. AWR learns a policy by weighted regression, where the weights are derived from the advantage function.

The AWR policy update can be formulated as:

$$\pi_{\text{new}} = \arg \max_\pi \mathbb{E}_{s \sim D, a \sim \pi_\beta(a|s)} \left[ \log \pi(a|s) \exp\left(\frac{1}{\beta} A^\pi(s, a)\right) \right] \tag{9}$$

where $A^\pi(s, a)$ is the advantage function and $\beta$ is a temperature parameter.

# 5 Multi-Agent Reinforcement Learning

## 5.1 Cooperative Multi-Agent RL

Cooperative multi-agent reinforcement learning (MARL) involves multiple agents working together to achieve a common goal. The challenge lies in coordinating the agents' actions to maximize the team's overall performance.

### 5.1.1 Value Decomposition Networks (VDN)

Value Decomposition Networks (VDN) address the multi-agent credit assignment problem by decomposing the team's value function into a sum of individual agent value functions:

$$Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{a}) = \sum_{i=1}^{n} Q_i(\tau_i, a_i) \tag{10}$$

where $\boldsymbol{\tau} = (\tau_1, \tau_2, \ldots, \tau_n)$ are the action-observation histories of all agents, and $\mathbf{a} = (a_1, a_2, \ldots, a_n)$ are their actions.

### 5.1.2 QMIX

QMIX extends VDN by using a mixing network that combines individual agent Q-values in a more expressive way while still maintaining the monotonicity property that enables efficient maximization:

$$Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{a}) = f_{\text{mix}}(Q_1(\tau_1, a_1), Q_2(\tau_2, a_2), \ldots, Q_n(\tau_n, a_n); \boldsymbol{\tau}) \tag{11}$$

where $f_{\text{mix}}$ is a mixing network that satisfies $\frac{\partial Q_{\text{tot}}}{\partial Q_i} \geq 0$ for all $i$.

## 5.2 Competitive and Mixed Multi-Agent RL

### 5.2.1 Multi-Agent Deep Deterministic Policy Gradient (MADDPG)

MADDPG extends the DDPG algorithm to multi-agent settings by using a centralized training with decentralized execution paradigm. Each agent has its own actor and critic, but the critic has access to all agents' actions and observations during training.

The MADDPG objective for agent $i$ can be formulated as:

$$J(\theta_i) = \mathbb{E}_{s, a_1, \ldots, a_n}[R_i(s, a_1, \ldots, a_n)] \tag{12}$$

The policy gradient for agent $i$ is:

$$\nabla_{\theta_i} J(\theta_i) = \mathbb{E}_{s, a_1, \ldots, a_n}[\nabla_{\theta_i} \log \pi_i(a_i|o_i) Q_i^{\pi}(s, a_1, \ldots, a_n)] \tag{13}$$

where $Q_i^{\pi}$ is the centralized action-value function for agent $i$.

### 5.2.2 Emergent Multi-Agent Coordination

Recent research has focused on how coordination emerges in multi-agent systems through self-play and population-based training. These approaches have led to the discovery of complex cooperative behaviors without explicit coordination mechanisms.

The key idea is to train a population of agents that learn to coordinate through repeated interactions:

$$\pi_i^* = \arg\max_{\pi_i} \mathbb{E}_{\pi_{-i} \sim P_{-i}, s \sim \rho^{\pi_i, \pi_{-i}}}[R_i(s, a_i, a_{-i})] \tag{14}$$

where $P_{-i}$ is a distribution over the policies of other agents, and $\rho^{\pi_i, \pi_{-i}}$ is the state distribution induced by the joint policy.

# 6 Reinforcement Learning with Human Feedback

## 6.1 Preference-Based RL

Preference-based RL learns from human preferences rather than explicit reward signals. This approach is particularly valuable when reward functions are difficult to specify or when aligning AI systems with human values is important.

### 6.1.1 Preference-Based Reward Learning

Preference-based reward learning infers a reward function from human preferences between pairs of trajectories. The Bradley-Terry model is commonly used to model these preferences:

$$P(\tau_1 \succ \tau_2) = \frac{\exp(r(\tau_1))}{\exp(r(\tau_1)) + \exp(r(\tau_2))} \tag{15}$$

where $r(\tau)$ is the inferred reward for trajectory $\tau$.

The reward function can be learned by maximum likelihood estimation:

$$\max_{r_\phi} \sum_{(\tau_1, \tau_2) \in D} \log P(\tau_1 \succ \tau_2 | r_\phi) \tag{16}$$

where $D$ is a dataset of preference pairs, and $r_\phi$ is a parameterized reward function.

## 6.2 Reinforcement Learning from Human Feedback (RLHF)

RLHF combines preference-based reward learning with reinforcement learning to train agents that align with human preferences. This approach has been particularly successful in training large language models.

The RLHF pipeline typically consists of three stages:

1. Supervised fine-tuning on human demonstrations

2. Reward model training from human preferences

3. RL optimization to maximize the learned reward

The RL optimization objective can be formulated as:

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}}[r_{\phi}(\tau)] - \beta D_{KL}(\pi_{\theta}||\pi_{\text{ref}}) \tag{17}$$

where $r_{\phi}$ is the learned reward function, $\pi_{\text{ref}}$ is a reference policy (typically the supervised fine-tuned model), and $\beta$ is a coefficient that controls the strength of the KL divergence penalty.

# 7 Emerging Trends and Future Directions

## 7.1 Scaling Laws for RL

Recent research has begun to investigate scaling laws for reinforcement learning, similar to those established for supervised learning. These studies aim to understand how performance improves with increases in model size, data, and computation.

Preliminary findings suggest that RL performance may follow power-law scaling with respect to model size and training compute, but with different exponents than supervised learning. Understanding these scaling laws could guide future research and development efforts.

## 7.2 Foundation Models for RL

The success of foundation models in natural language processing and computer vision has inspired similar approaches in reinforcement learning. Foundation models for RL would be pre-trained on diverse tasks and environments, then fine-tuned for specific applications.

Potential approaches include:

- Pre-training world models on diverse environments

- Learning general-purpose representations from large datasets of agent experiences

- Developing multi-task policies that can be adapted to new tasks with minimal fine-tuning

## 7.3 Causal Reinforcement Learning

Causal reinforcement learning incorporates causal reasoning into the RL framework to improve generalization and transfer. By learning causal models of the

environment, agents can better predict the consequences of their actions and generalize to new situations.

Causal RL approaches include:

- Learning structural causal models from agent experiences

- Using causal discovery algorithms to identify causal relationships between variables

- Incorporating causal knowledge into policy learning and planning algorithms

# 8    Conclusion

This article has explored state-of-the-art methods in reinforcement learning, highlighting recent advancements that push the boundaries of what AI systems can achieve. From evolutionary approaches to algorithm design to meta-reinforcement learning, offline RL, and multi-agent systems, these methods address key challenges in reinforcement learning and open new possibilities for applications.

As noted by [1], "The evolution of reinforcement learning algorithms is accelerating, driven by a combination of theoretical insights, empirical discoveries, and computational advances." This rapid progress suggests a bright future for reinforcement learning, with increasingly capable and versatile algorithms that can tackle complex real-world problems.

The emerging trends discussed in this article, including scaling laws, foundation models, and causal reinforcement learning, point to exciting directions for future research. As these areas develop, we can expect reinforcement learning to continue its trajectory of innovation and impact across diverse domains.

# References

[1] Google Research (2023). *Evolving Reinforcement Learning Algorithms.* Retrieved from https://research.google/blog/evolving-reinforcement-learning-algorithms/

[2] Raschka, S. (2025). *The State of Reinforcement Learning for LLM Reasoning.* Retrieved from https://sebastianraschka.com/blog/2025/the-state-of-reinforcement-learning-for-llm-reasoning.html

[3] Towards Data Science (2023). *RL from One Example:   Why 1-Shot RLVR Might Be the Breakthrough We've Been Waiting For.* Retrieved from https://towardsdatascience.com/rl-from-one-example-why-1-shot-rlvr-might-be-the-breakthrough-weve-been-waiting-for/