# A Survey of In-Context Reinforcement Learning

**Amir Moeini**[1] , **Jiuqi Wang**[1] , **Jacob Beck**[2] , **Ethan Blaser**[1] ,
**Shimon Whiteson**[2] , **Rohan Chandra**[1] , **Shangtong Zhang**[1]

[1]University of Virginia

[2]University of Oxford

{amoeini, shangtong}@virginia.edu

## Abstract

Reinforcement learning (RL) agents typically optimize their policies by performing expensive backward passes to update their network parameters. However, some agents can solve new tasks without updating any parameters by simply conditioning on additional context such as their action-observation histories. This paper surveys work on such behavior, known as in-context reinforcement learning.

## 1 Introduction

Reinforcement learning (RL) [Sutton and Barto, 2018] is a paradigm for solving sequential decision-making tasks via trial and error. In RL, an agent incrementally optimizes its policy as it interacts with its environment so as to maximize a reward signal in the long run. Here, the policy is a function that maps the agent's observation to the distribution from which its actions are sampled. Policies are typically represented with neural networks whose parameters are continually updated during learning. Selecting actions requires a forward pass through the network. Updating the parameters usually requires a backward pass, which can be expensive for large neural networks in terms of both memory and computation [Kingma and Ba, 2017].

Some pretrained RL agents, however, can solve new tasks without updating any network parameters. When evaluating such an agent on a new task, the input to the agent includes both the current observation and additional context that helps the agent adapt to the new task. For example, the context may include the agent's history of observations and actions in this new task up to the current time step. The remarkable ability of such agents to generalize to new tasks using context but without fine-tuning is hypothesized [Duan *et al.*, 2016; Laskin *et al.*, 2023] to arise from **the pretrained neural network implementing some (known or unknown) RL algorithm in its forward pass to process the context**. **We refer to this behavior in the forward pass as in-context reinforcement learning (ICRL)**. An immediate implication of ICRL is that the agent's performance improves as the task-related information in the context accumulates, a phenomenon called *in-context improvement*. Figure 1 illustrates ICRL's pipeline. We follow Sutton and Barto [2018] and define RL algorithms broadly as learning algorithms for solving sequential decision-making problems. This includes, for example, imitation learning [Abbeel and Ng, 2005] and temporal difference learning (TD, Sutton [1988]).

We argue that ICRL is important because it enables agents to generalize to new tasks efficiently, requiring only a forward pass without expensive parameter updates. Eliminating the need for parameter updates creates new opportunities to optimize computation and memory requirements for inference [Zhu *et al.*, 2024]. Furthermore, there is also evidence that the RL algorithms implemented in the forward pass can potentially be more sample efficient than manually engineered ones [Laskin *et al.*, 2023].

**Scope.** ICRL falls in the category of black box methods for meta RL [Beck *et al.*, 2023] and dates back to Duan *et al.* [2016]; Wang *et al.* [2016]. Early works of ICRL demonstrate only limited out-of-distribution generalization (see Section 4 for more discussion) and are well surveyed by Beck *et al.* [2023]. We view Laskin *et al.* [2023] as an important milestone of ICRL since it both coins the term and provides the first demonstration of remarkable out-of-distribution generalization of the pretrained agent. This survey, therefore, focuses on ICRL work after Laskin *et al.* [2023].

**Taxonomy.** In this paper, we survey work on ICRL along different axes. We start with different pretraining methods, such as supervised pretraining and reinforcement pretraining. We then examine different methods for constructing context at test time and the demonstrated test time performance. We then survey recent theoretical advances in understanding ICRL and, finally, neural network architecture design choices in both empirical and theoretical work on ICRL.

**Related Surveys.** The closest work to this paper is Beck *et al.* [2023], which surveys meta-RL. Beck *et al.* [2023] do not include recent ICRL advances showing strong out-of-distribution generalization since Laskin *et al.* [2023]. This survey aims to close this gap and is thus complementary to Beck *et al.* [2023]. See Beck *et al.* [2023] for a full treatment of other aspects of meta RL. The development of ICRL parallels the development of in-context learning (ICL), where supervised (instead of reinforcement) learning occurs in the forward pass. See Dong *et al.* [2023] for a full treatment of ICL. Furthermore, some ICRL pretraining methods resemble RL via supervised learning (RvS) in the offline RL community and goal-conditioned RL. However, many methods for
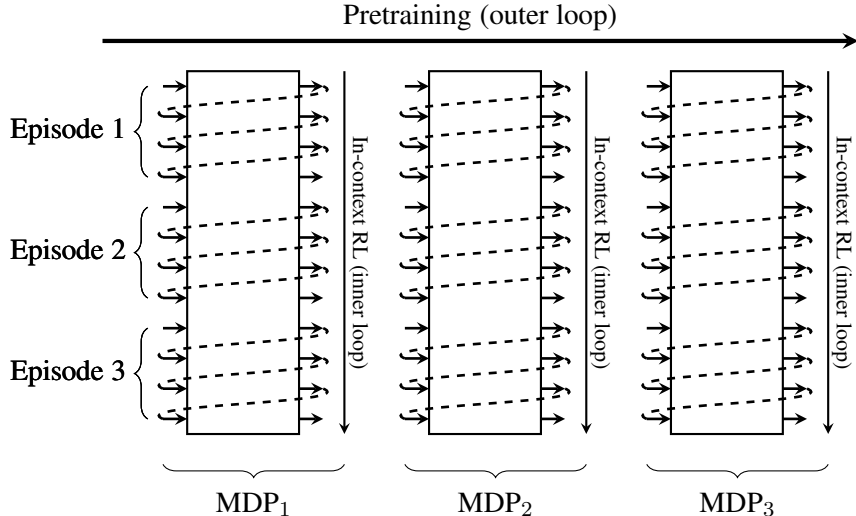
Figure 1: Overview of ICRL. After pretraining, the forward pass of the network implements some RL algorithm. The implemented RL algorithm is tested on multiple MDPs. The context in each MDP can span multiple episodes.

RvS and goal-conditioned RL, e.g., Chen *et al.* [2021], do not demonstrate key properties of ICRL, such as in-context improvement, and are therefore not included in this survey. That being said, we include RvS works that are able to demonstrate ICRL, e.g., Laskin *et al.* [2023]. Nevertheless, see Emmons *et al.* [2022]; Wen *et al.* [2023] for a full treatment of RvS, and Liu *et al.* [2022] for goal-conditioned RL.

## 2 Background

**Reinforcement Learning.** RL typically models environments or tasks as Markov Decision Processes (MDPs). An MDP consists of a state space $\mathcal{S}$, an action space $\mathcal{A}$, a reward function $r : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, a transition function $p : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \to [0,1]$, an initial distribution $p_0 : \mathcal{S} \to [0,1]$, and a discount factor $\gamma \in [0,1)$. At time step 0, an initial state $S_0$ is sampled from $p_0$. At time $t$, an agent at a state $S_t$ takes an action $A_t$ according to its policy $\pi : \mathcal{A} \times \mathcal{S} \to [0,1]$, i.e., $A_t \sim \pi(\cdot|S_t)$. The agent then receives a reward $R_{t+1} \doteq r(S_t, A_t)$ and proceeds to a successor state $S_{t+1} \sim p(\cdot|S_t, A_t)$. We use $\tau_t$ to denote the trajectory up to time $t$, defined as, $\tau_t \doteq (S_0, A_0, R_1, S_1, A_1, \ldots, S_{t-1}, A_{t-1}, R_t)$. The value function of the policy $\pi$ is defined as $v_\pi(s) \doteq \mathbb{E}\left[\sum_{i=1}^{\infty} \gamma^{i-1} R_{t+i}|S_t = s\right]$. There are two fundamental tasks in RL (given an environment). The first is policy evaluation, where the goal is to estimate $v_\pi$ for a given $\pi$. The second is control, where the goal is to find a policy $\pi$ to maximize the expected total rewards $J(\pi) \doteq \sum_s p_0(s)v_\pi(s)$. Policy evaluation is the foundation of control. Given a policy $\pi$ and its value function, a better-performing policy $\pi'$ can be computed. This is called a *policy improvement* step. Iterating between policy evaluation and policy improvement can generate the optimal policy to maximize $J(\pi)$.

RL methods usually adopt a parameterized policy $\pi_\theta$ (and a parameterized value estimate $v_\theta$). For example, $\theta$ could be the parameters of a neural network. Many RL algorithms have been proposed to update $\theta$ during the agent-environment interaction. This is referred to as online pretraining. By contrast, offline pretraining relies on a dataset $D$ consisting of previously collected transition tuples $D \doteq \{(s_i, a_i, r_i, s'_i)\}_{i=1,\ldots,K}$, where $r_i$ is the reward recevied after executing the action $a_i$ in the state $s_i$, and $s'_i$ is the corresponding successor state. This offline dataset may come from one or more (known or unknown) behavior policies. The transition tuples may or may not form a complete trajectory. Many RL algorithms are proposed to update $\theta$ using $D$ as well. Notably, both online and offline pretraining can involve learning a single policy $\pi_\theta$ from one or multiple MDPs. When a single policy is used for multiple MDPs, it needs additional input besides the current state to differentiate between and adapt to different MDPs. After pretraining, whether offline or online, the policy $\pi_\theta$ is evaluated on the same MDP or one or more new MDPs.

**In-Context Reinforcement Learning.** The key idea of ICRL is to condition the policy both on $S_t$ and some context $C_t$, such that the action $A_t$ is sampled according to $\pi_\theta(\cdot|S_t, C_t)$, rather than solely being conditioned on $S_t$. This survey considers different ways to construct $C_t$, but a simple example is to use $\tau_t$ as $C_t$. Some pretraining methods can generate $\theta$ such that the policy $\pi_\theta(\cdot|S_t, C_t)$ can obtain high rewards in test MDPs that differ from the MDPs seen during pretraining, despite $\theta$ remaining fixed. Such generalization is hypothesized to arise because the forward pass of the neural network $\theta$ implements an RL algorithm that learns from the current context $C_t$ [Duan *et al.*, 2016; Laskin *et al.*, 2023]. This is called ICRL (for control). Furthermore, the performance of $\pi_\theta$ improves with the length of the context $C_t$ (provided that $C_t$ always consists of information related to the task), which we call in-context improvement. Similarly, RL algorithms for policy evaluation can also be implemented in the forward pass if the neural network for value estimation takes as input both the state and the context [Wang *et al.*, 2024a]. This is called ICRL for policy evaluation.

# 3 Supervised Pretraining

This section surveys the first class of methods for pretraining the network parameter $\theta$ – supervised pretraining, which is commonly done through behavior cloning. The pretraining objective is typically the log-likelihood $\log \pi_\theta(a^*|s,c)$ or its variants, where $(s,c)$ is the input state-context pair and $a^*$ is the desired output action. There are multiple ways to construct the input and output.

The most common approach is to concatenate the trajectories from multiple episodes as input (called *cross-episode* input) and use the corresponding action at each step as output. Various ways are proposed to obtain the episode trajectories for constructing the input. Laskin *et al.* [2023] use trajectories generated by some existing RL algorithms across their entire lifecycle. As a result, the input includes trajectories produced by existing RL algorithms during both early pretraining stages (when agents perform poorly) and later stages (when agents achieve strong performance). Building on Laskin *et al.* [2023], Shi *et al.* [2023] propose to fill the context with a curriculum of trajectories. Specifically, they order the trajectories by task difficulty, demonstrator proficiency, or episode returns [Huang *et al.*, 2024a,b; Liu and Abbeel, 2023]. Alternative curriculum construction methods include adding decaying noise to expert demonstration trajectories [Zisman *et al.*, 2023] and using an explicit feature in the trajectory to indicate whether the current episode is better than the ones before [Dai *et al.*, 2024]. This indicator, called cross-episode return-to-go in Dai *et al.* [2024], alleviates the need to rank episodes in the context. All of these methods aim to demonstrate performance improvement in the input-out pairs, encouraging the neural network to implement **some policy improvement algorithm** in the forward pass. Kirsch *et al.* [2023] show that the performance improvement rate in the input trajectories directly influences how quickly the pretrained agent improves at test time.

Some methods instead emphasize finding trajectories similar to the current task from the offline data and prepending them to the context for better generalization [Xu *et al.*, 2022; Wang *et al.*, 2024b]. These initial trajectories are referred to as a *prompt*. This encourages the network to implement **some imitation learning algorithm** in the forward pass to imitate the behavior demonstrated in the context. This approach is most effective when the prompt is collected by an expert agent [Raparthy *et al.*, 2023; Xu *et al.*, 2022]. Consequently, these methods require expert prompts at test time and cannot use their own suboptimal interactions as prompts. To mitigate this issue, Lee *et al.* [2023] construct the prompts using the suboptimal trajectories but replace the actions in the trajectories with actions predicted by a performant policy.

Some works incorporate hindsight information into context to facilitate policy optimization [Furuta *et al.*, 2022]. Hindsight information is data that can only be inferred from outcomes of future time steps. A common example is the Return-To-Go (RTG), which represents the sum of rewards from a step until the end of the episode, [Xu *et al.*, 2022; Wang *et al.*, 2024b; Dai *et al.*, 2024]. However, inspired by how RTG is estimated during testing (Section 5), some methods replace this with a target RTG, defined as either the maximum RTG in the offline dataset [Schmied *et al.*, 2024] or among episodes in the current context [Huang *et al.*, 2024a,b].

To summarize, while some works do not rely on a curriculum in the context [Lee *et al.*, 2023], they all rely on multiple trajectories in the context. Using a single trajectory in the context [Chen *et al.*, 2021] provably prevents an agent from improving on the offline-data-generating policies or "stitching" suboptimal trajectories to get an optimal one without strong assumptions [Brandfonbrener *et al.*, 2022].

Sample efficiency during supervised pretraining is also an important research area. Here we define sample efficiency broadly, including the amount of data, the amount of expert / optimal policy demonstration, and the amount of different tasks. Zisman *et al.* [2024] greatly improve the pretraining sample efficiency of Laskin *et al.* [2023] by hardcoding $n$-gram induction heads into the transformer and biasing it toward using in-context information [Akyürek *et al.*, 2024]. Dai *et al.* [2024] use importance sampling to remove the pessimistic bias that keeps the pretrained policy close to the data collection policy, thereby enabling the viable use of suboptimal data collection policies. Similarly, Dong *et al.* [2024] employ a weighted loss, where the weights, based on the observed rewards, act as importance sampling ratios to guide the suboptimal policy toward the optimal policy. With a different approach, Zisman *et al.* [2023] use a suboptimal policy and add annealing noise to its trajectories to generate learning histories similar to those of Laskin *et al.* [2023]. Kirsch *et al.* [2023] construct augmented tasks to improve sample efficiency by randomly projecting the task's observation and action spaces. While these methods enhance overall sample efficiency, robotic tasks introduce unique challenges. Due to long sequence lengths, full rollouts are inefficient. Elawady *et al.* [2024] mitigate this issue by employing partial rollouts that reduce environment interactions.

Finally, there are multiple ways to encode the context before it is presented to the neural network. For example, in Transformer-based methods, the token at each time step can be the stacked embeddings of either $(s_t, a_t, r_{t+1}, s_{t+1})$ [Kirsch *et al.*, 2023; Lee *et al.*, 2023] or $(a_{t-1}, r_t, s_t)$ [Laskin *et al.*, 2023; Raparthy *et al.*, 2023; Zisman *et al.*, 2023]. Similarly, action spaces can be embedded in distinct ways. For instance, to enable test-time adaptation to varying action spaces, Sinii *et al.* [2023] project discrete actions into random vector embeddings and train the network to output an embedding vector directly. Then, the action whose embedding is most similar to the network output is executed.

# 4 Reinforcement Pretraining

This section surveys the second class of methods for pretraining the network parameter $\theta$ – reinforcement pretraining. Instead of using the log-likelihood loss, reinforcement pretraining uses other established RL algorithms to train the policy $\pi_\theta(a|s,c)$. In contrast to supervised pretraining which uses offline datasets, reinforcement pretraining usually involves online environment interactions.

Early ICRL works with reinforcement pretraining include Duan *et al.* [2016]; Wang *et al.* [2016]; Mishra *et al.* [2018];

Ritter *et al.* [2018]; Stadie *et al.* [2019]; Zintgraf *et al.* [2020]; Melo [2022], where a sequence model (e.g., an RNN) parameterizes the policy. At test time, the policy takes the agent's online interaction history (usually across multiple episodes) as input and outputs the action, without any parameter updates. *This history-dependent policy effectively functions as an RL algorithm, as both take the complete history as input and output an action.* However, early works in this line of research demonstrate only limited out-of-distribution generalization. They only demonstrate that the learned history-dependent policy performs well in tasks similar to pretraining tasks. One hypothesis for the lack of out-of-distribution generalization in those works is that the pretrained network implements **some task identification algorithm together with certain nearest neighbor matching**. In other words, at test time, the pretrained network tries to identify pretraining tasks that are similar to the test task (based on the entire online interaction history in the test task) and acts as if the test task was one of those similar pretraining tasks. Those task identification works are well surveyed by Beck *et al.* [2023].

In this paper, we instead focus on more recent advances in ICRL, where pretrained networks demonstrate stronger out-of-distribution generalization by implementing more advanced RL algorithms in the forward pass. These include Grigsby *et al.* [2024, 2023]; Lu *et al.* [2023]; Bauer *et al.* [2023]; Park *et al.* [2024]; Wang *et al.* [2024a]; Elawady *et al.* [2024]; Xu *et al.* [2024]; Cook *et al.* [2024]. Although their pretraining is still performed by (modification of) standard RL algorithms (Table 1), using long-context neural networks such as Transformers to parameterize the policy introduces substantial learning stability challenges [Grigsby *et al.*, 2023]. To improve stability, several modifications to the pretraining process have been proposed. Grigsby *et al.* [2023] employ multiple discount rates simultaneously to stabilize long-horizon credit assignment. Bauer *et al.* [2023] introduce a task selection strategy that prioritizes tasks slightly beyond the agent's current expertise, significantly enhancing sample efficiency. For scenarios involving learning across tasks with highly varied return scales, Grigsby *et al.* [2024] utilize actor-critic objectives decoupled from return magnitudes, thereby improving convergence. That being said, why the recent works are able to demonstrate stronger generalization remains an open problem (Section 9).

## 5 Test Time Context

Having surveyed the two main pretraining paradigms, we now turn to test time design choices, beginning with the context construction. While context construction is often similar during both pretraining and testing, some information provided in the context construction during pretraining is not available at test time. This section starts with examining how to address these differences.

One example of such information is expert demonstrations. Methods using such demonstrations in pretraining [Rakelly *et al.*, 2019; Lee *et al.*, 2023; Wang *et al.*, 2024b] often experience significant performance drops when

| Method | Pretraining Algorithm |
|---|---|
| Grigsby *et al.* [2023] | Modified DDPG [Lillicrap *et al.*, 2019] |
| Grigsby *et al.* [2024] | Modified DDPG |
| Lu *et al.* [2023] | Muesli [Hessel *et al.*, 2021] |
| Bauer *et al.* [2023] | Muesli |
| Elawady *et al.* [2024] | Modified PPO [Schulman *et al.*, 2017] |
| Wang *et al.* [2024a] | TD |
| Park *et al.* [2024] | Regret minimization |
| Cook *et al.* [2024] | PPO |
| Xu *et al.* [2024] | Modified DQN [Mnih *et al.*, 2015] |

Table 1: Algorithms used for reinforcement pretraining.

prompted with suboptimal interactions in test time. This decline occurs because the model has limited exposure to the near-optimal trajectory space during testing, leading to a mismatch between the context distribution in pretraining and testing.

If at test time the agent outputs actions using only offline trajectories as context, the performance will heavily depend on the quality of the offline demonstrations [Lee *et al.*, 2023]. This is also the case if the agent is expected to output good actions after obtaining only one or a few online trajectories [Raparthy *et al.*, 2023; Wang *et al.*, 2024b]. But when the agent is allowed to interact more extensively with the environment at test time before it is expected to output good actions, the need for expert demonstrations in the context is reduced [Laskin *et al.*, 2023; Huang *et al.*, 2024a,b; Sinii *et al.*, 2023; Kirsch *et al.*, 2023; Lee *et al.*, 2023]. Lee *et al.* [2023] demonstrate the trade-off between the allowed online interaction budget and the need for expert demonstrations.

Another example of context information that is not available at test time is RTG from Section 3. Various methods are used to estimate RTG at test time. Huang *et al.* [2024b,a]; Xu *et al.* [2022]; Schmied *et al.* [2024]; Wang *et al.* [2024b] approximate RTG once per task based on the offline trajectories. At the beginning of the episode, an initial RTG is given. This RTG is iteratively updated based on the observed reward. Alternatively, Dai *et al.* [2024] train a secondary network to predict RTG dynamically based on the interaction history.

During testing, methods that do not rely on demonstrations must learn the task solely from their own previous interactions. However, we can selectively choose what to include in the context. For instance, Cook *et al.* [2024] divide the total interaction horizon into generations, with each generation comprising several agents. These agents systematically use the best-performing agent from the previous generation to generate interactions that are then incorporated into the current context. This approach allows the current agent to build upon prior experience. The framework, referred to as cultural accumulation, achieves superior test-time performance scaling compared to the base single-generation method.

# 6 Test Time Performance

In this section, we survey the test-time performance of pre-trained agents from two aspects, generalization and sample efficiency. Since comparing generalization across different benchmarks is challenging, we consider generalization benchmark by benchmark. Notably, unlike goal-conditioned methods that explicitly condition the pretrained agent on the test task type, the ICRL agent must infer it implicitly by itself.

The first remarkable out-of-distribution generalization in ICRL is demonstrated by Laskin *et al.* [2023] in multi-armed bandit problems. Their pretrained agents learn new bandit problems with adversarial rewards (i.e., engineered so that pretraining-optimal policies perform poorly) and achieve regret nearly equivalent to standard bandit algorithms that involve parameter updates.

ICRL's out-of-distribution generalization improves as models, pretraining duration, and experience diversity scale [Bauer *et al.*, 2023; Kirsch *et al.*, 2023]. For instance, Bauer *et al.* [2023] propose XLand 2.0, a procedurally generated 3D environment featuring diverse goals, rules, and configurations. They demonstrate generalization on this challenging benchmark using ICRL, enabled by large-scale supervised pretraining with a curriculum and other improvements.

Other commonly used benchmarks in the ICRL literature include Dark Room [Laskin *et al.*, 2023], DMLab Watermaze [Laskin *et al.*, 2023], Procgen [Cobbe *et al.*, 2020], Meta-World [Yu *et al.*, 2021], and Mujoco Control [Todorov *et al.*, 2012]. Each benchmark requires a specific form of generalization which is demonstrated by different works. Test tasks can differ from pretraining tasks across various task-dependent factors (such as goal location or object types). Accordingly, the difficulty of generalization can be better understood by considering the types of factors and the extent of variations [Kirk *et al.*, 2021].

To succeed in the Dark Room benchmark and its variants, the agent should learn to efficiently find held-out (i.e., not used during pretraining) invisible goal locations. Each unique goal location (or combination of them in the key-to-door variant) represents a new task. This generalization is demonstrated by Laskin *et al.* [2023]; Lee *et al.* [2023]; Huang *et al.* [2024a,b]; Sinii *et al.* [2023]; Zisman *et al.* [2023, 2024]; Kirsch *et al.* [2023]; Dai *et al.* [2024]; Grigsby *et al.* [2023]; Elawady *et al.* [2024] to different degrees. Notably, Grigsby *et al.* [2023] adapts to a new key-to-door environment in only 300 interactions.

Similarly, to succeed in DMLab Watermaze, where the input to the agent is raw pixels, the agent needs to find a trapdoor in new locations of a maze. This generalization is demonstrated by Laskin *et al.* [2023]; Zisman *et al.* [2023]; Shi *et al.* [2023]; Ritter *et al.* [2018] to different degrees.

To succeed in the Mujoco Control benchmark, the agent must control simulated robots to achieve given tasks (e.g., make a HalfCheetah run or an Ant navigate) with variations in dynamics (e.g., altered friction or mass) or target parameters (e.g., desired speed or direction). This generalization is demonstrated by Xu *et al.* [2022]; Wang *et al.* [2024b]; Grigsby *et al.* [2023]; Mishra *et al.* [2018]; Melo [2022]; Kirsch *et al.* [2023] to different degrees. In particular,

Kirsch *et al.* [2023] show that after pretraining on Ant tasks, the agent can solve the Cartpole task in DeepMind Control Suite.

Procgen is a benchmark consisting of 16 procedurally generated 2D games (e.g., platformers, puzzles) with pixel observations. To demonstrate generalization, the agent should learn held-out games [Raparthy *et al.*, 2023]. The games differ across many factors (such as objects, objectives, and types of movement), which have proved difficult to generalize to, especially when expert demonstrations are not available during the test. Grigsby *et al.* [2024]; Schmied *et al.* [2024] show initial progress in an easier setting, where they test on the same pretraining games with limited modifications (e.g., changes to the starting location or textures)

Meta-World consists of robotic manipulation challenges. In Meta-Learning 1 (ML1), the variations are continuous (e.g., different object or goal positions) within a single manipulation category. By contrast, ML45 uses 45 manipulation categories (e.g., opening drawers or turning faucets) for pretraining and 5 new categories for testing. Several studies have shown models generalizing on ML1 [Xu *et al.*, 2022; Wang *et al.*, 2024b; Grigsby *et al.*, 2023; Mishra *et al.*, 2018; Melo, 2022], and Grigsby *et al.* [2024] show generalization on the 45 pretraining manipulation categories of ML45 in a setting similar to the limited one described earlier for Procgen.

We now turn to the sample efficiency of the pretrained agents in the test time. Laskin *et al.* [2023] demonstrate that a pretrained network with fixed parameters needs fewer samples at test time to achieve similar performance to that of baseline RL algorithms that require gradient updates. Kirsch *et al.* [2023] successfully control test-time sample efficiency by manipulating how much an episode improves upon the previous one when constructing the cross-episode pretraining contexts. Lee *et al.* [2023] show that the forward pass of their pretrained network is an efficient implementation of posterior sampling, a sample-efficient RL algorithm, under specific conditions during the test. Likewise, Xu *et al.* [2024] propose an end-to-end framework for learning an agent that performs Bayesian inference in context, thereby improving test-time sample efficiency on out-of-distribution tasks. Sample efficiency remains a challenge in sparse-reward tasks when the agent is not sufficiently biased toward thorough exploration. Stadie *et al.* [2019]; Norman and Clune [2023] propose addressing this issue by modifying the objective to maximize only the cumulative reward of later exploitive episodes, thereby allowing the initial explorative episodes to focus on better exploration for subsequent exploitive episodes.

# 7 Theory

We now consider recent advances in the theoretical understanding of ICRL.

**Supervised Pretraining.** Supervised pretraining can be understood through the lens of behavior cloning. In canonical behavior cloning, the goal is to learn a policy. The policy usually depends on only the current state or the history within the current episode. In supervised pretraining,

the goal is to learn an algorithm similar to the source algorithm used to generate the offline dataset. In other words, to learn a policy that depends on the entire history of previous episodes. Lin *et al.* [2023] derive a general bound on the behavioral similarity and performance gap between the learned algorithm (in the forward pass of the neural network) and the source algorithm. Behavioral similarity is the similarity between the action distributions generated by the learned and source algorithms given the same input. The performance gap is their difference in episode return. They further demonstrate how Transformers can approximate Lin-UCB [Chu *et al.*, 2011], Thompson sampling [Russo *et al.*, 2018], and UCB-VI [Azar *et al.*, 2017] in the forward pass and provide the respective regret bounds. A follow-up work by Shi *et al.* [2024] presents an analogous behavioral similarity guarantee of supervised pretraining for decentralized and centralized learning in two-player zero-sum Markov games. Shi *et al.* [2024] further prove by construction that there exist Transformers that can realize V-learning [Jin *et al.*, 2024] for decentralized learning and VI-ULCB [Bai and Jin, 2020] for centralized learning in the forward passes, accompanied with upper bounds of the approximation error of Nash equilibria for both settings.

**Reinforcement Pretraining.** Park *et al.* [2024] propose to pretrain language models directly by minimizing regret without requiring action labels. They show theoretically that by minimizing regret with sufficiently many pretraining trajectories, the pretrained language models can demonstrate no-regret learning at test time. Lastly, they prove that the global minimizer of the (surrogate) regret loss with a single-layer linear attention transformer implements the known no-regret algorithm Follow-The-Regularized-Leader (FTRL) [Shalev-Shwartz and Singer, 2007] in the forward pass. Wang *et al.* [2024a] consider ICRL for policy evaluation. They prove by construction that Transformers can precisely implement temporal difference methods in the forward pass for policy evaluation, including TD($\lambda$) [Sutton, 1988] and average reward TD [Tsitsiklis and Roy, 1999]. They also show that those parameters naturally emerge when they train a value estimation transformer with TD on multiple policy evaluation tasks. Theoretical understanding of this emergence of TD is provided from an invariant set perspective.

# 8 Architectures

A central design choice in ICRL is the architecture of the neural network used to process context. The neural network must be able to handle long context lengths, often containing multiple episodes of interaction, and effectively use information from many past interactions.

Although earlier meta RL works [Duan *et al.*, 2016; Wang *et al.*, 2016; Ritter *et al.*, 2018] use RNN and its variants to parameterize history-dependent policies, most surveyed ICRL works employ a causal transformer backbone [Laskin *et al.*, 2023; Lee *et al.*, 2023; Raparthy *et al.*, 2023; Sinii *et al.*, 2023; Zisman *et al.*, 2023; Shi *et al.*, 2023; Kirsch *et al.*, 2023; Xu *et al.*, 2022; Grigsby *et al.*, 2024, 2023; Bauer *et al.*, 2023; Melo, 2022; Elawady *et al.*, 2024; Zisman *et al.*, 2024], given transformer's demonstrated ef-

ficacy in handling long sequences [Vaswani *et al.*, 2017]. However, the inference time of Transformers is quadratic w.r.t. the input length. To speed it up, state space models [Gu and Dao, 2023], whose inference time is linear, are used [Cook *et al.*, 2024]. Huang *et al.* [2024b] employ a state space model for their high-level decision maker, which processes long histories, and a transformer for their low-level decision maker, which processes shorter sequences. Lu *et al.* [2023] modify an existing state space model, S5 [Smith *et al.*, 2023], such that it becomes compatible with cross-episode context. Schmied *et al.* [2024] use an xLSTM [Beck *et al.*, 2024] for similar purposes.

Hierarchical structures are also designed with different objectives. For instance, Wang *et al.* [2024b] improve Xu *et al.* [2022] by incorporating additional modules that extract both task-level and step-specific prompts relevant to the current task and step, which are then used to augment the context provided to a Transformer. Dai *et al.* [2024] use a secondary network to predict the RTG required for the context during inference. To improve computational efficiency by processing fewer tokens in a transformer without sacrificing overall historical information, Huang *et al.* [2024a,b] split decision-making into two levels. Specifically, the high-level module processes tokens sampled at fixed intervals, while the low-level module predicts the intervening tokens corresponding to each high-level token.

Compared to supervised pretraining, reinforcement pretraining introduces engineering challenges regarding stability during pretraining [Grigsby *et al.*, 2023]. To address these challenges, Grigsby *et al.* [2023] share a single sequence model across both actor and critic networks and demonstrate that preventing the critic's objective from minimizing the actor's objective can ensure pretraining stability. They also modify the transformer architecture to preserve plasticity over long pretraining durations and avoid performance collapse, also adopted by Xu *et al.* [2024]. Similarly, Elawady *et al.* [2024] append learnable key and value vectors as "sinks" to the transformer's attention mechanism to provide the flexibility of not attending to any input token. This modification makes learning faster and more stable in scenarios involving long but low-information observation sequences, such as those encountered in robotics [Elawady *et al.*, 2024].

Regarding theoretical analysis, having a full-sized multi-layer transformer with arbitrary nonlinear activations is prohibitively challenging due to the complexity of the network structure. Lin *et al.* [2023] and Shi *et al.* [2024] use masked attentions with ReLU activations, and Park *et al.* [2024]; Wang *et al.* [2024a] use linear attentions.

# 9 Open Problems and Opportunities

ICRL is an emerging area with many open problems. First, we draw attention to multi-agent RL. Generalization to unseen agents (teammates or opponents) during the deployment time is a fundamental challenge in multi-agent RL and meta RL has been applied to address this challenge [Charakorn *et al.*, 2021; Gerstgrasser and Parkes, 2022]. However, the demonstrated generalization is only limited in small-scale problems and only limited out-of-

distribution generalization is demonstrated. Recent advances in multi-agent RL with large sequence models [Meng *et al.*, 2023] provide a new opportunity to address this challenge with ICRL.

Second, we draw attention to robotics. ICRL is now only demonstrated in simulated environments. The sim-to-real gap is a well-known generalization challenge in robotics. It is a promising direction to investigate whether ICRL will emerge in recent internet-scale robot pretraining [Brohan *et al.*, 2023] and whether ICRL can help close the sim-to-real gap.

Lastly, we draw attention to reinforcement pretraining. Krishnamurthy *et al.* [2024] criticize ICRL saying "they are explicitly trained to produce this behavior using data from reinforcement learning agents or expert demonstrations on related tasks." This criticism might be true for supervised pretraining but clearly does not hold for reinforcement pretraining, where the network is only trained to output good actions or to output good value estimates without constraints on how the network achieves this. The network itself discovers that implementing certain RL algorithms in the forward pass is a good solution. In this sense, ICRL truly emerges during reinforcement pretraining. Existing theoretical analyses on reinforcement pretraining [Park *et al.*, 2024; Wang *et al.*, 2024a] use simplified models and simplified pretraining algorithms. Fully white-boxing the emergence of ICRL during reinforcement pretraining in more realistic settings remains an open problem, both theoretically and empirically.

## 10 Conclusion

This paper presented the first comprehensive survey of ICRL, an emerging and flourishing area. We surveyed ICRL from different aspects, including both pretraining and testing, both empirical and theoretical analyses. We hope this survey will stimulate the growth of the ICRL community.

## Acknowledgements

## References

Pieter Abbeel and Andrew Y. Ng. Exploration and apprenticeship learning in reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, 2005.

Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-Context Language Learning: Architectures and Algorithms. *ArXiv preprint*, 2024.

Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, 2017.

Yu Bai and Chi Jin. Provable self-play algorithms for competitive reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, 2020.

Jakob Bauer, Kate Baumli, Feryal M. P. Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmieg, Michael Chang, Natalie Clay, Adrian Collister, Vibhavari Dasagi, Lucy Gonzalez, Karol Gregor, Edward Hughes, Sheleem Kashem, Maria Loks-Thompson, Hannah Openshaw, Jack Parker-Holder, Shreya Pathak, Nicolas Perez Nieves, Nemanja Rakicevic, Tim Rocktäschel, Yannick Schroecker, Satinder Singh, Jakub Sygnowski, Karl Tuyls, Sarah York, Alexander Zacherl, and Lei M. Zhang. Human-timescale adaptation in an open-ended task space. In *Proceedings of the International Conference on Machine Learning*, 2023.

Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon Whiteson. A Survey of Meta-Reinforcement Learning. *ArXiv preprint*, 2023.

Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xLSTM: Extended Long Short-Term Memory. *ArXiv preprint*, 2024.

David Brandfonbrener, Alberto Bietti, Jacob Buckman, Romain Laroche, and Joan Bruna. When does return-conditioned supervised learning work for offline reinforcement learning? In *Advances in Neural Information Processing Systems*, 2022.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.

Rujikorn Charakorn, Poramate Manoonpong, and Nat Dilokthanakul. Learning to cooperate with unseen agent via meta-reinforcement learning. *arXiv preprint arXiv:2111.03431*, 2021.

Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In *Advances in Neural Information Processing Systems*, 2021.

Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2011.

Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, 2020.

Jonathan Cook, Chris Lu, Edward Hughes, Joel Z. Leibo, and Jakob Nicolaus Foerster. Artificial Generational Intelligence: Cultural Accumulation in Reinforcement Learning.

In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

Zhenwen Dai, Federico Tomasi, and Sina Ghiassian. In-context Exploration-Exploitation for Reinforcement Learning. *ArXiv preprint*, 2024.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A Survey on In-context Learning. *ArXiv preprint*, 2023.

Juncheng Dong, Moyang Guo, Ethan X. Fang, Zhuoran Yang, and Vahid Tarokh. In-Context Reinforcement Learning Without Optimal Action Labels. In *ICML 2024 Workshop on In-Context Learning*, 2024.

Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. RL$^2$: Fast Reinforcement Learning via Slow Reinforcement Learning. *ArXiv preprint*, 2016.

Ahmad Elawady, Gunjan Chhablani, Ram Ramrakhya, Karmesh Yadav, Dhruv Batra, Zsolt Kira, and Andrew Szot. ReLIC: A Recipe for 64k Steps of In-Context Reinforcement Learning for Embodied AI. *ArXiv preprint*, 2024.

Scott Emmons, Benjamin Eysenbach, Ilya Kostrikov, and Sergey Levine. Rvs: What is essential for offline RL via supervised learning? In *Proceedings of the International Conference on Learning Representations*, 2022.

Hiroki Furuta, Yutaka Matsuo, and Shixiang Shane Gu. Generalized decision transformer for offline hindsight information matching. In *Proceedings of the International Conference on Learning Representations*, 2022.

Matthias Gerstgrasser and David C Parkes. Meta-rl for multi-agent rl: Learning to adapt to evolving agents. In *Advances in Neural Information Processing Systems*, 2022.

Jake Grigsby, Linxi Fan, and Yuke Zhu. AMAGO: Scalable In-Context Reinforcement Learning for Adaptive Agents. *ArXiv preprint*, 2023.

Jake Grigsby, Justin Sasek, Samyak Parajuli, Daniel Adebi, Amy Zhang, and Yuke Zhu. AMAGO-2: Breaking the Multi-Task Barrier in Meta-Reinforcement Learning with Transformers. *ArXiv preprint*, 2024.

Albert Gu and Tri Dao. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. *ArXiv preprint*, 2023.

Matteo Hessel, Ivo Danihelka, Fabio Viola, Arthur Guez, Simon Schmitt, Laurent Sifre, Theophane Weber, David Silver, and Hado Van Hasselt. Muesli: Combining improvements in policy optimization. In *International conference on machine learning*, 2021.

Sili Huang, Jifeng Hu, Hechang Chen, Lichao Sun, and Bo Yang. In-Context Decision Transformer: Reinforcement Learning via Hierarchical Chain-of-Thought. *ArXiv preprint*, 2024.

Sili Huang, Jifeng Hu, Zhejian Yang, Liwei Yang, Tao Luo, Hechang Chen, Lichao Sun, and Bo Yang. Decision Mamba: Reinforcement Learning via Hybrid Selective Sequence Modeling. *ArXiv preprint*, 2024.

Chi Jin, Qinghua Liu, Yuanhao Wang, and Tiancheng Yu. V-learning—a simple, efficient, decentralized algorithm for multiagent reinforcement learning. *Mathematics of Operations Research*, 2024.

Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *ArXiv preprint*, 2017.

Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A Survey of Zero-shot Generalisation in Deep Reinforcement Learning. *ArXiv preprint*, 2021.

Louis Kirsch, James Harrison, C. Daniel Freeman, Jascha Sohl-Dickstein, and Jürgen Schmidhuber. Towards General-Purpose In-Context Learning Agents. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.

Akshay Krishnamurthy, Keegan Harris, Dylan J. Foster, Cyril Zhang, and Aleksandrs Slivkins. Can large language models explore in-context? *ArXiv preprint*, 2024.

Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald, DJ Strouse, Steven Stenberg Hansen, Angelos Filos, Ethan A. Brooks, Maxime Gazeau, Himanshu Sahni, Satinder Singh, and Volodymyr Mnih. In-context reinforcement learning with algorithm distillation. In *Proceedings of the International Conference on Learning Representations*, 2023.

Jonathan Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma Brunskill. Supervised pretraining can learn in-context reinforcement learning. In *Advances in Neural Information Processing Systems*, 2023.

Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *ArXiv preprint*, 2019.

Licong Lin, Yu Bai, and Song Mei. Transformers as Decision Makers: Provable In-Context Reinforcement Learning via Supervised Pretraining. *ArXiv preprint*, 2023.

Hao Liu and Pieter Abbeel. Emergent agentic transformer from chain of hindsight experience. In *Proceedings of the International Conference on Machine Learning*, 2023.

Minghuan Liu, Menghui Zhu, and Weinan Zhang. Goal-conditioned reinforcement learning: Problems and solutions. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 2022.

Chris Lu, Yannick Schroecker, Albert Gu, Emilio Parisotto, Jakob N. Foerster, Satinder Singh, and Feryal M. P. Behbahani. Structured state space models for in-context reinforcement learning. In *Advances in Neural Information Processing Systems*, 2023.

Luckeciano C. Melo. Transformers are meta-reinforcement learners. In *Proceedings of the International Conference on Machine Learning*, 2022.

Linghui Meng, Muning Wen, Chenyang Le, Xiyun Li, Dengpeng Xing, Weinan Zhang, Ying Wen, Haifeng Zhang, Jun

Wang, Yaodong Yang, et al. Offline pre-trained multi-agent decision transformer. *Machine Intelligence Research*, 2023.

Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. In *Proceedings of the International Conference on Learning Representations*, 2018.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 2015.

Ben Norman and Jeff Clune. First-Explore, then Exploit: Meta-Learning to Solve Hard Exploration-Exploitation Trade-Offs. *ArXiv preprint*, 2023.

Chanwoo Park, Xiangyu Liu, Asuman Ozdaglar, and Kaiqing Zhang. Do LLM Agents Have Regret? A Case Study in Online Learning and Games. *ArXiv preprint*, 2024.

Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In *Proceedings of the International Conference on Machine Learning*, 2019.

Sharath Chandra Raparthy, Eric Hambro, Robert Kirk, Mikael Henaff, and Roberta Raileanu. Generalization to New Sequential Decision Making Tasks with In-Context Learning. *ArXiv preprint*, 2023.

Samuel Ritter, Jane X. Wang, Zeb Kurth-Nelson, Siddhant M. Jayakumar, Charles Blundell, Razvan Pascanu, and Matthew Botvinick. Been there, done that: Meta-learning with episodic recall. In *Proceedings of the International Conference on Machine Learning*, 2018.

Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial on thompson sampling. *Foundations and Trends® in Machine Learning*, 2018.

Thomas Schmied, Thomas Adler, Vihang Patil, Maximilian Beck, Korbinian Pöppel, Johannes Brandstetter, Günter Klambauer, Razvan Pascanu, and Sepp Hochreiter. A Large Recurrent Action Model: xLSTM enables Fast Inference for Robotics Tasks. *ArXiv preprint*, 2024.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *ArXiv preprint*, 2017.

Shai Shalev-Shwartz and Yoram Singer. A primal-dual perspective of online learning algorithms. *Machine Learning*, 2007.

Lucy Xiaoyang Shi, Yunfan Jiang, Jake Grigsby, Linxi Fan, and Yuke Zhu. Cross-episodic curriculum for transformer agents. In *Advances in Neural Information Processing Systems*, 2023.

Chengshuai Shi, Kun Yang, Jing Yang, and Cong Shen. Transformers as Game Players: Provable In-context Game-playing Capabilities of Pre-trained Models. *ArXiv preprint*, 2024.

Viacheslav Sinii, Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, and Sergey Kolesnikov. In-Context Reinforcement Learning for Variable Action Spaces. *ArXiv preprint*, 2023.

Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. Simplified state space layers for sequence modeling. In *Proceedings of the International Conference on Learning Representations*, 2023.

Bradly C. Stadie, Ge Yang, Rein Houthooft, Xi Chen, Yan Duan, Yuhuai Wu, Pieter Abbeel, and Ilya Sutskever. Some Considerations on Learning to Explore via Meta-Reinforcement Learning. *ArXiv preprint*, 2019.

Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction (2nd Edition)*. MIT press, 2018.

Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 1988.

Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *Proceedings of the International Conference on Intelligent Robots and Systems*, 2012.

John N. Tsitsiklis and Benjamin Van Roy. Average cost temporal-difference learning. *Automatica*, 1999.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.

Jane X. Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z. Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *ArXiv preprint*, 2016.

Jiuqi Wang, Ethan Blaser, Hadi Daneshmand, and Shangtong Zhang. Transformers Learn Temporal Difference Methods for In-Context Reinforcement Learning. *ArXiv preprint*, 2024.

Zhe Wang, Haozhu Wang, and Yanjun Qi. Hierarchical Prompt Decision Transformer: Improving Few-Shot Policy Generalization with Global and Adaptive Guidance. *ArXiv preprint*, 2024.

Muning Wen, Runji Lin, Hanjing Wang, Yaodong Yang, Ying Wen, Luo Mai, Jun Wang, Haifeng Zhang, and Weinan Zhang. Large Sequence Models for Sequential Decision-Making: A Survey. *ArXiv preprint*, 2023.

Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua B. Tenenbaum, and Chuang Gan. Prompting decision transformer for few-shot policy generalization. In *Proceedings of the International Conference on Machine Learning*, 2022.

Tengye Xu, Zihao Li, and Qinyuan Ren. Meta-Reinforcement Learning Robust to Distributional Shift Via Performing Lifelong In-Context Learning. In *Proceedings of the International Conference on Machine Learning*, 2024.

Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Avnish Narayan, Hayden Shively, Adithya Bellathur, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-World:

A Benchmark and Evaluation for Multi-Task and Meta Reinforcement Learning. *ArXiv preprint*, 2021.

Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. *Transactions of the Association for Computational Linguistics*, 2024.

Luisa M. Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep RL via meta-learning. In *Proceedings of the International Conference on Learning Representations*, 2020.

Ilya Zisman, Vladislav Kurenkov, Alexander Nikulin, Viacheslav Sinii, and Sergey Kolesnikov. Emergence of In-Context Reinforcement Learning from Noise Distillation. *ArXiv preprint*, 2023.

Ilya Zisman, Alexander Nikulin, Andrei Polubarov, Nikita Lyubaykin, and Vladislav Kurenkov. N-Gram Induction Heads for In-Context RL: Improving Stability and Reducing Data Needs. *ArXiv preprint*, 2024.