# Meta-Reinforcement Learning for Fast Adaptation: Approaches and Case Study

Milan Das

Department of Electrical Engineering and Computer Science
University of California, Irvine
Email: mbdas@uci.edu

*Abstract*—This work explores fast adaptation in reinforcement learning through a streamlined architecture designed for clarity and effectiveness. This work contextualizes a simplified architecture within the broader space of meta-reinforcement learning approaches, such as gradient-based methods like Model-Agnostic Meta-Learning (MAML) and recurrent policy structures for task inference. While these techniques are powerful, they often depend on intricate optimization routines or tightly coupled architectures that can be difficult to interpret or scale.

Motivated by these challenges, this paper proposes a simpler two-part framework that separates task inference from control. A feedforward context encoder processes short interaction rollouts and generates fixed-size latent embeddings, which are used to condition a downstream policy. The encoder is trained using supervised signals derived from task-specific rollouts, while the policy learns through standard reinforcement learning. This setup allows the system to generalize to new tasks without requiring any gradient updates at test time.

Experiments are run on a suite of CartPole variants with hidden dynamics parameters. Using just 15-step sample trajectories, the method consistently achieves high returns on previously unseen environments. These findings support the idea that explicit context representation, when paired with a conditioned policy, can enable lightweight, zero-shot adaptation in meta-reinforcement learning.

## I. INTRODUCTION

Meta-Reinforcement Learning (Meta-RL) focuses on training agents that can adapt quickly to new tasks using minimal experience. Unlike standard reinforcement learning, where an agent learns a single task from scratch, meta-RL aims to develop agents that can generalize across tasks by learning how to adapt. The core idea is to expose agents to a distribution of related tasks during training so that they can rapidly fine-tune their behavior when faced with an unfamiliar scenario.

Fast adaptation becomes especially important in domains where collecting large amounts of task-specific data is costly or impractical. To address this, researchers have introduced a variety of strategies, broadly categorized into model-free and model-based approaches, and spanning different adaptation mechanisms, from gradient-based policy updates to context-conditioned behavior driven by latent task representations.

This paper reviews several representative methods that have shaped the meta-RL landscape, including gradient-based algorithms like Model-Agnostic Meta-Learning (MAML) and linear representation frameworks such as FLAP. The discussion highlights how these techniques differ in terms of architecture, optimization, and assumptions about task structure.

To ground the review in practice, the paper also examines a specific implementation: a context-conditioned policy architecture applied to a set of CartPole environments with varied dynamics. This system uses a trajectory-based encoder to infer latent task context and a policy network conditioned on that context to generate actions. The analysis considers how this setup compares to prior methods, why it achieves strong zero-shot generalization, and what lessons it offers for designing adaptable agents. Potential improvements based on recent meta-RL research are also discussed.
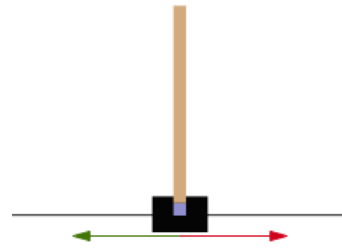


Figure 1. CartPole environment from Gym library used for meta-RL.

## II. METHODOLOGIES

Meta-reinforcement learning approaches generally fall into two broad categories based on how they approach task adaptation: model-free and model-based. Model-free methods are entirely based on direct policy learning, without an explicit model of the dynamics of the environment. Instead, they enable adaptation through mechanisms such as gradient updates to policy parameters or internal memory systems that condition behavior on past experience. In contrast, model-based methods learn a predictive model of the environment and use it to guide planning or inform policy adaptation in a data-efficient way.

Among model-free strategies, gradient-based meta-RL has received considerable attention for its ability to fine-tune policies with just a few gradient steps. One of the most well-known examples is Model-Agnostic Meta-Learning (MAML) [1], which optimizes a shared initialization across tasks such that small amounts of task-specific data lead to significant performance gains. The core idea is to train the meta-policy to be highly adaptable, functioning as an initialized learner that quickly converges to a task-optimal solution.

MAML's flexibility comes from being model-agnostic: it can be applied to any architecture trained with gradient descent. In reinforcement learning, it has shown success by accelerating policy learning across diverse task distributions. Several extensions have built on this idea, including first-order approximations like FOMAML, and alternatives such as

Reptile and ProMP. ProMP [2], in particular, integrates meta-learning with on-policy optimization methods like TRPO or PPO to improve post-adaptation performance.

Most of these methods define a meta-objective centered around maximizing expected return after adaptation. During meta-training, the system computes gradients through the adaptation process, often requiring second-order derivatives to align the initialization with optimal task-specific learning. These techniques are typically benchmarked in continuous control environments with variations in dynamics or goals. An example would be altering target velocities or morphologies in MuJoCo tasks. What makes gradient-based meta-RL compelling is its ability to quickly adapt to new environments using only a handful of fresh trajectories, while maintaining high performance.

Black-box or memory-based meta-reinforcement learning methods offer a model-free route to adaptation by relying on internal memory systems, often implemented with recurrent neural networks, to encode task information directly from experience. These approaches skip gradient updates at test time; instead, the agent adapts by updating its internal state. A foundational example is $RL^2$ [3][4], which uses a recurrent policy, usually an LSTM, that learns to represent task-specific behavior through the evolution of its hidden state over multiple episodes. During meta-training, the recurrent network effectively learns an embedded reinforcement learning procedure. At deployment, policy parameters remain fixed; adaptation emerges entirely from how the hidden state evolves during interaction. This enables rapid within-episode adjustment, though the method often demands extensive training across diverse tasks and can suffer from instability or difficulty in credit assignment. More recent extensions explore transformers for in-context reinforcement learning and memory-augmented networks that treat entire trajectory sequences as prompts.

Context-based meta-RL methods build on similar intuitions, but introduce a more explicit mechanism for task inference. Rather than relying solely on hidden state dynamics, these methods infer a latent task representation, often called a context variable, based on recent interactions. This context vector is then used to condition the policy and, in some cases, the value function as well. A leading example is PEARL [5], which uses an inference network to encode short trajectories into a probabilistic latent variable $z$, allowing the policy $\pi(a \mid s, z)$ to adapt its behavior accordingly. PEARL blends an off-policy reinforcement learning objective with a variational inference loss, encouraging the context to capture features critical for effective task differentiation and control. Variants like VariBAD [6] apply similar principles, but often swap the variational encoder for a recurrent one or apply predictive objectives to infer context. Many of these approaches augment training with auxiliary losses, for example, learning to predict dynamics or optimizing contrastive objectives (e.g., InfoNCE), to ensure the context embedding is informative and relevant to downstream behavior. These methods are particularly effective in zero-shot adaptation settings, where agents must adjust using only context derived from brief interactions, without any test-time gradient updates.

In contrast, model-based meta-RL approaches prioritize fast adaptation by learning and leveraging environment models. The key assumption is that if an agent can quickly approximate the dynamics of a new task, it can use this model to guide efficient planning or policy adjustment. One line of work involves training dynamics models that themselves can be adapted quickly, for instance, using MAML-style updates. Another strategy involves training policies that generalize across an ensemble of learned dynamics models, with the aim of making a single policy robust to a wide range of environmental variations. A notable example is Model-Based Meta-Policy Optimization (MB-MPO) [7], which trains a policy across an ensemble of models, optimizing it so that a single gradient update allows rapid adaptation to any member of the ensemble. Other extensions combine latent variable models with dynamics learning, where a shared model is modulated by a latent context variable inferred from interactions. The major strengths of model-based approaches lie in their sample efficiency and interpretability, though their success depends heavily on how well the learned models capture real environment variability, especially in high-complexity domains.

### III. Effectiveness of Key Techniques

Across leading meta-RL approaches, a few consistent design patterns have proven particularly effective for enabling fast adaptation. One of the most important is task conditioning through learned embeddings. Regardless of whether the method is gradient-based or context-driven, most fast-adapting systems rely on representing task identity in a form that the policy can use. In gradient-based techniques like MAML, this comes through rapid adjustment of model weights, effectively shaping the policy to fit a new task in just a few steps. Context-based methods, on the other hand, make use of explicit task embeddings or memory states that are fed as inputs to the policy network. For example, PEARL introduces a latent context variable $z$ that compactly summarizes the task and enables the shared policy to generate tailored behavior by conditioning on $z$.

This idea of task-conditioned policy behavior is central to most modern meta-RL architectures. It allows a single network to generalize across many tasks by adjusting its outputs based on an inferred context. Architecturally, this is often done by concatenating the context vector with the current state or applying modulation techniques like FiLM within the network layers. In recurrent setups like $RL^2$ or VariBAD, the hidden state of the recurrent network evolves over time and serves as a kind of internal context, naturally integrating task inference into the agent's temporal processing.

Another key pattern that emerges, especially in context-based methods, is the separation between inference and control. Unlike MAML, where adaptation is embedded directly into the optimization loop, newer approaches often decouple the process of inferring task structure from the decision-making logic. In these systems, an encoder is responsible for

extracting a latent task representation from observed trajectories, which is then passed to a downstream policy for action selection. This modular design has a few clear advantages: it makes the architecture more interpretable and allows for greater flexibility in tuning and reuse.

A good example of this is FLAP [8], which learns a linear, shared policy representation and uses a separate adapter network to generate task-specific parameters. At test time, FLAP skips gradient updates entirely, producing new weights for the policy via a single forward pass through the adapter. This not only speeds up adaptation but also supports generalization to novel task distributions. More broadly, studies have shown that well-trained encoders can maintain robustness even under distributional shifts, especially when the latent task space is properly structured or regularized. These ideas strongly motivate the encoder-policy split used in the system studied here, where brief interaction trajectories are used to extract task information that informs policy behavior without requiring additional learning.

Another important set of factors that influence the effectiveness of meta-reinforcement learning methods involves how loss functions and training strategies are designed. In gradient-based approaches like MAML, the meta-objective is directly aligned with the adaptation goal. The model is trained to perform well after just a few gradient steps on a new task. This tight coupling helps guide the optimization toward solutions that generalize quickly. Context-based methods, in contrast, often rely on auxiliary objectives that shape the task encoder to produce informative and discriminative embeddings.

For example, in PEARL, the encoder isn't just optimized to support policy behavior. It is also trained to predict the task's value function. Empirical results have shown that conditioning on value estimates, rather than returns, can improve generalization across tasks. Other approaches, such as Meta-Q-Learning, introduce additional objectives tied to value-based representations of dynamics. Contrastive learning is another common strategy: by encouraging embeddings from the same task to cluster together while pushing apart those from different tasks, the encoder learns a more structured and useful representation space. These design choices make it easier for the policy to adapt quickly based on compact contextual signals.

Exploration also plays a crucial role in meta-RL. In $RL^2$, for instance, the agent learns exploration strategies as part of its recurrent policy. The agent's behavior in early steps of a new task influences the evolution of its hidden state, which in turn guides future actions. This results in implicit adaptation through exploration. Some other methods take a more explicit route, encouraging informative behavior through information gain objectives or using strategies like Thompson sampling to actively seek out data that helps disambiguate the task.

Training regimes also vary considerably between on-policy and off-policy meta-RL methods. On-policy algorithms, such as MAML or PPO-based $RL^2$, collect fresh data for each meta-update. This ensures that the adaptation and training data distributions are closely aligned, but it can be highly sample-

inefficient. Off-policy methods like PEARL, on the other hand, benefit from replay buffers that allow them to reuse past experience, leading to much higher sample efficiency. However, this comes with its own trade-offs. Maintaining a stable and effective inference network requires careful management of the mismatch between replayed data and current policy behavior.

Finally, the choice of benchmarking environments plays a critical role in evaluating meta-RL systems. Early research often used simple domains such as multi-armed bandits or gridworlds to study exploration across episodes. As the field progressed, continuous control benchmarks like those in the MuJoCo suite (e.g., Half-Cheetah, Ant) became the standard. These benchmarks vary task parameters such as goal velocity, mass, or friction to test adaptation in more realistic settings. The CartPole environment used in this work serves as a more tractable but still representative setup: task variation is introduced by modifying pole mass or length, creating distinct control challenges that require quick adaptation.

While simpler than locomotion or manipulation tasks, CartPole still captures the essential challenge of meta-RL: identifying underlying task dynamics and adjusting behavior with minimal data. Some benchmarks push this further by introducing sparse rewards, which force the agent to explore before it can begin optimizing. As newer suites like Meta-World continue to raise the bar in terms of complexity and sensory richness, the fundamental objective remains the same. The goal is to develop systems that can infer task structure and adapt policies efficiently in response to limited, noisy, or ambiguous information.

## IV. CASE STUDY: CONTEXT-CONDITIONED POLICY

The system [9] examined in this work is a meta-reinforcement learning (meta-RL) setup built around two main components: (1) a trajectory-based context encoder, and (2) a policy network conditioned on the inferred context. The system is trained and tested on a family of CartPole environments, each differing in underlying dynamics, specifically categorized into light, medium, and heavy variants. These variants correspond to changes in pole or cart mass, affecting the difficulty of maintaining balance.
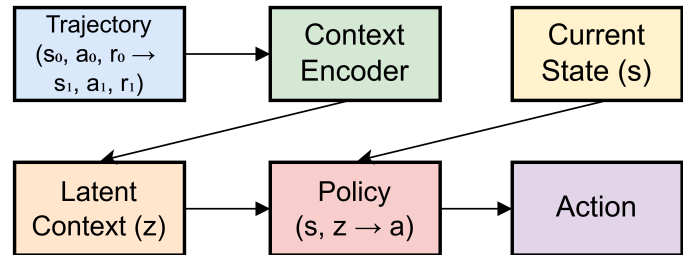


Figure 2. System architecture: A trajectory is encoded into a latent context vector, which is concatenated with the current state and passed to the policy to generate an action. This enables fast adaptation via forward inference.

At inference time, the context encoder takes a short trajectory, typically a partial episode, and encodes it into a fixed-length vector representing the latent task. This context vector captures the environment-specific dynamics observed during

interaction. The policy network receives the current state (such as cart position, velocity, pole angle, and angular velocity) along with the context vector and produces an action, applying force to the left or right.

Training uses the REINFORCE algorithm, an on-policy Monte Carlo policy gradient method. Both the encoder and the policy are trained jointly. During each training episode, the agent samples from one of the known task variants and interacts briefly with the environment. The encoder learns to summarize these interactions into informative task embeddings, while the policy learns to act optimally given this latent representation. Over time, the system develops the ability to differentiate task identities based on trajectory information and adjust its behavior accordingly.

In the evaluation phase, the system is tested on CartPole variants it has not seen during training, with altered pole lengths or masses. Importantly, there are no gradient updates at test time. Instead, the agent performs a forward pass through the trained encoder using a short rollout in the new environment, generating a context vector that immediately informs the policy. This design allows the system to achieve fast, zero-shot adaptation, relying purely on inference to generalize across task variations with minimal online interaction.
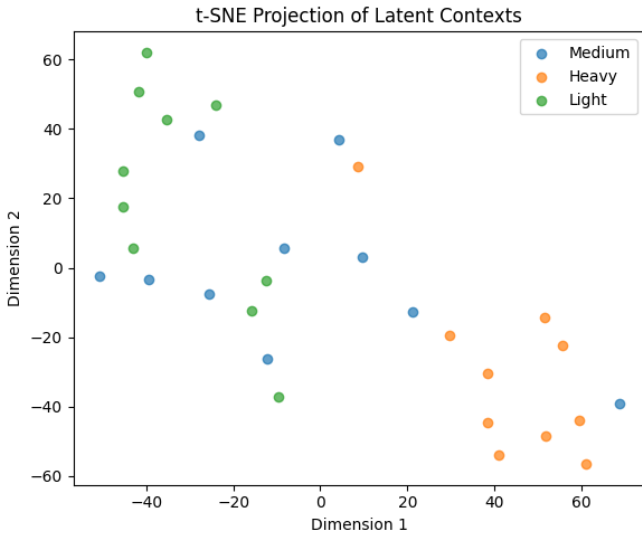


Figure 3. 2D projection of learned context embeddings using t-SNE.

This architecture exemplifies a common structure in context-based meta-RL, where task inference and control are cleanly separated. The encoder performs implicit system identification from short trajectories, while the policy uses this information to produce context-aware actions. Even in this relatively simple domain, the system demonstrates strong generalization. The latent structure learned by the encoder is visualized in Figure 3, where embeddings form well-separated clusters based on pole mass. Embeddings from the unseen environment are positioned between those of the trained environments, suggesting the encoder generalizes via interpolation.

This setup reflects core principles seen across context-based meta-reinforcement learning. Architecturally, it follows a familiar pattern: a context encoder produces a latent task embedding, which then conditions a downstream policy. This structure is similar to methods like PEARL, where separating task inference from control improves both modularity and interpretability. In this system, the encoder handles the inference of task identity from past experience, while the policy focuses solely on action selection based on current observations and the inferred context.

Unlike PEARL, which uses a probabilistic inference network trained with variational objectives, this implementation opts for a simpler approach, training the encoder via policy gradients using the REINFORCE algorithm. While this reduces architectural complexity, it still achieves the key goal of generating informative context vectors that reflect environment-specific behavior. The encoder's role here is conceptually similar to the adapter module in FLAP [8], which learns to produce task-specific policy parameters. In the CartPole setup, the encoder doesn't generate weights directly but instead outputs a latent context based on short interaction histories. Both systems rely on forward inference rather than test-time optimization, allowing rapid behavioral adaptation.

This type of black-box adaptation, where behavior shifts purely through inference, also parallels the approach taken in RL$^2$, where a recurrent policy adapts via changes in its hidden state. The key difference here is architectural: instead of embedding task memory internally within the policy, this system externalizes it through a dedicated encoder. This design adds flexibility, as it opens the door to alternative encoder structures such as recurrent or transformer-based models that could better handle temporal reasoning or longer trajectories.

One practical aspect worth highlighting is the choice of REINFORCE for training. While straightforward to implement, REINFORCE is known for its high variance and slower convergence. More recent work often prefers algorithms like PPO or off-policy methods for better sample efficiency and stability. However, in a structured, low-dimensional domain like CartPole, REINFORCE proves adequate. It enables the encoder to learn representations that align with task-specific reward signals, implicitly guiding the context embedding to reflect meaningful variations in dynamics. This echoes findings from other meta-RL approaches that use auxiliary losses or contrastive learning to ensure task relevance is encoded effectively.

Performance results, shown in Figure 4, indicate that the system achieves near-optimal returns on all training variants and sustains similarly high returns on the unseen task (average > 195 with low variance). This confirms that the policy can exploit latent task embeddings to generalize without retraining. The encoder successfully extracts dynamics-relevant structure, and the policy responds robustly to novel inputs.

First, the context encoder implicitly performs system identification. Although the different CartPole environments vary only by a physical parameter such as mass or length, these differences show up in how the system behaves. For example, heavier poles fall faster and require more force to stabilize.
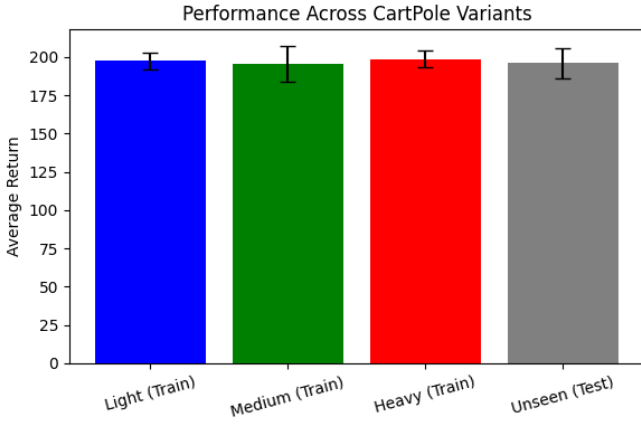
## Performance Across CartPole Variants

Figure 4. Average episode returns across variants over 20 trials.

By observing a short trajectory, the encoder picks up on these cues and encodes them into a latent vector that reflects the underlying task dynamics. This task representation is generated without any explicit labels, yet it provides the policy with enough information to adjust its behavior. For instance, in a "heavy pole" scenario, the policy might learn to apply larger or more frequent corrective forces. This interaction, where the encoder interprets behavioral patterns and the policy adapts accordingly, allows the system to perform well after observing just a single short episode.

Second, the architectural separation between inference and control proves essential. The encoder is focused solely on generating a concise and informative task representation based on recent experience, while the policy maps this representation, combined with the current state, to an appropriate action. This modular structure mirrors a probabilistic reasoning process: infer what kind of environment the agent is in, then act optimally based on that inference. Over training, the encoder becomes better at distilling meaningful patterns from trajectories, and the policy becomes more adept at leveraging those patterns to specialize its behavior. This design supports generalization by allowing one policy network to emulate a wide range of behaviors, with the appropriate one selected dynamically through context input. This kind of context-based conditioning has become a recurring design pattern in meta-RL because of its effectiveness in disambiguating tasks and improving performance compared to multitask policies that treat all situations uniformly.

Third, the system achieves zero-shot generalization by structuring its latent task representations in a meaningful way. The encoder is trained to generate embeddings that reflect task differences, and over time, it organizes these embeddings so that similar tasks are close together in latent space. As a result, when the agent encounters a task that falls between two known variants, such as one with medium-heavy pole mass, the resulting context vector lies between the two corresponding embeddings. Because the policy is a continuous function of this context input, it can interpolate smoothly and produce appropriate behavior even for unseen tasks. This ability to

generalize beyond the training distribution is a strength of context-based approaches like this one and is also emphasized in methods such as FLAP, which achieve fast adaptation through linear parameter modulation. In contrast, methods like MAML require test-time optimization and may struggle on out-of-distribution tasks. The current system instead adapts through a single inference pass.

Finally, the system's performance on both seen and unseen tasks provides indirect but strong evidence that the context embeddings are informative and well-structured. If the encoder were producing weak or noisy representations, the policy would be left trying to map ambiguous state-context combinations to actions, resulting in inconsistent performance. Instead, the consistent high returns (with means near 200 and low standard deviation across trials) suggest that the context vector carries enough task-specific information to guide reliable behavior. During training, the policy gradients encourage the encoder to refine its representations so that they align with task identity. This ongoing feedback between encoder and policy leads to co-adaptation: the encoder learns to extract the right signals, and the policy learns how to respond to them. The fact that the system performs well even on CartPole variants it was never trained on indicates that it has captured a generalizable mapping from observed dynamics to effective behavior, which is a defining goal of meta-reinforcement learning.

To assess the contribution of the context encoder to adaptation, an ablation study was conducted on the held-out environment, comparing the full system to two degraded variants: (1) a policy conditioned on a random latent vector sampled from $\mathcal{N}(0, I)$, and (2) a policy with no context signal (i.e., a zero vector input). These tests help clarify whether the policy truly relies on the encoder's output, and how it behaves without meaningful context.

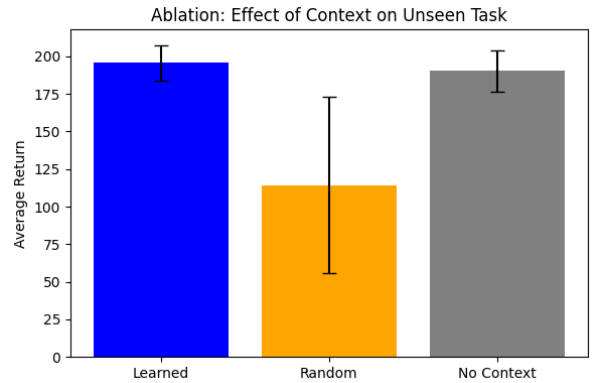## Ablation: Effect of Context on Unseen Task

Figure 5. Ablation results on the unseen test environment.

As shown in Figure 5, the full system with learned context achieves the highest average return (mean $\approx 195$). In contrast, conditioning on a random vector leads to a sharp drop in performance (mean $\approx 115$) and high variance, indicating that the policy is sensitive to task embedding quality. Notably, the no-context condition (mean $\approx 190$) outperforms the random case, suggesting the policy defaults to a robust but generic

behavior in the absence of context. These results confirm that the encoder is not ignored during training; rather, it plays a meaningful role in helping the policy specialize to unseen dynamics with minimal interaction.

## V. POTENTIAL IMPROVEMENTS AND VARIANTS

Several potential enhancements drawn from recent meta-RL research could improve the system's performance, scalability, and robustness. These suggestions address both conceptual considerations and practical bottlenecks, offering directions for continued experimentation.

One area worth exploring is the use of probabilistic context representations. In its current form, the encoder likely produces a deterministic latent vector, which works well in a clean, structured environment like CartPole. However, in more complex or ambiguous settings, a deterministic encoding may be too rigid. Introducing a probabilistic latent space, as done in PEARL, would allow the model to represent uncertainty over the task identity when only limited or noisy trajectory data is available. This uncertainty could guide the policy to behave more cautiously or to seek out information that helps disambiguate the task. For example, if two tasks produce similar initial dynamics, a deterministic encoder might commit prematurely to the wrong embedding. In contrast, a probabilistic one could maintain a distribution over possible contexts and act accordingly. This could be implemented using stochastic encoders trained with a variational objective like the ELBO, or via regularization techniques such as dropout. While this level of flexibility may not be necessary for CartPole, it becomes increasingly important in partially observable or multimodal environments.

Another promising direction is active exploration for better context acquisition. Currently, the encoder receives trajectories generated by a default or previously learned policy, but these may not always contain the most informative signals. Some recent work, such as Contextual Meta-RL with Contrastive Learning (CCM), has shown that training an exploration policy specifically to reduce task uncertainty can significantly improve adaptation. In CartPole, this might mean deliberately applying varied control inputs at the start of an episode to test how the environment reacts. Rather than immediately trying to solve the task, the agent would begin with a brief "probing" phase to gather useful context. This kind of explore-then-exploit strategy can lead to better performance, especially in tasks where the distinguishing features are subtle or only appear under certain conditions.

The learning algorithm itself also presents an opportunity for improvement. Although REINFORCE served its purpose in this implementation, it is widely known to suffer from high variance and sample inefficiency. More stable alternatives like Proximal Policy Optimization (PPO) or Soft Actor-Critic (SAC) are commonly used in modern reinforcement learning and would likely accelerate training and improve policy robustness. Additionally, moving to an off-policy framework, similar to what PEARL employs, would make it possible to reuse past experiences through a replay buffer, improving sample efficiency. This is especially useful when training context encoders that need to generalize across a wide variety of trajectories. While coordinating encoder and policy updates in an off-policy setting is more complex, it has been shown to work in practice and could substantially improve training efficiency.

The encoder itself could also benefit from stronger training signals. In the current system, it is optimized indirectly through the policy's reward, which is a relatively weak supervisory signal. Incorporating auxiliary objectives could help shape more structured and informative latent representations. If environment parameters like pole mass are available during training, the encoder could be trained to predict them directly. Alternatively, it could learn to estimate returns, predict future states, or apply contrastive learning to encourage similar tasks to cluster in latent space. These techniques help enforce structure in the embedding space and make the resulting context vectors more meaningful, especially under noisy or ambiguous conditions.

Taken together, these improvements highlight a number of practical pathways for scaling the system to more complex, less predictable domains. While the current implementation performs well in CartPole, more realistic environments will demand richer representations, better data efficiency, and more deliberate exploration strategies.

Complementing the encoder-level improvements are architectural modifications that could help scale the system to more complex or partially observable environments. While CartPole is fully observable and operates in a low-dimensional space, real-world tasks often involve hidden variables or noisy, incomplete inputs. One natural extension would be to redesign the encoder to handle temporal sequences using recurrent structures like LSTMs or GRUs. A recurrent encoder could accumulate information over time, enabling it to infer latent task characteristics that only become apparent after several steps of interaction. For example, subtle differences in inertia or delayed feedback may not be visible in a single short trajectory but could be detected through longer-term patterns. This kind of temporal integration maintains the existing modular split between inference and control while improving robustness to sequential variability.

In scenarios where the input is high-dimensional, such as tasks involving visual observations, the encoder could also be extended with convolutional layers or transformer-based perception modules. These perception-driven encoders have been increasingly adopted in modern meta-RL benchmarks like Meta-World, where agents must adapt to diverse robotic manipulation tasks in visually rich environments. Incorporating such modules would bring the architecture closer to real-world applicability.

Another promising direction is to explore hybrid approaches that combine context-based inference with lightweight gradient-based fine-tuning. In the current design, all adaptation occurs through the latent context. However, some recent work has shown that treating the context as a prior or initialization for further gradient updates can strike a useful balance between

speed and flexibility. For example, the context vector might modulate the initial policy parameters, with one or two quick updates refining the behavior for edge-case tasks that fall outside the training distribution. While this introduces some test-time compute overhead, it can be valuable in harder settings where inference alone is insufficient.

Evaluating out-of-distribution (OOD) generalization is also an important next step. While the system generalizes well to interpolated tasks, like pole masses between the "light" and "heavy" variants, it is less clear how well it extrapolates to more extreme or atypical dynamics. Running evaluations on tasks with ultra-heavy poles or extremely short lengths could help reveal whether the encoder has learned a meaningful latent structure or is simply interpolating between memorized classes. In the FLAP framework, this kind of extrapolation is supported by enforcing linearity in policy space. Although this system does not explicitly model such structure, a well-trained latent space, especially one shaped using domain randomization, could support similar generalization. By training on a continuous range of pole masses rather than a few fixed categories, the encoder would be encouraged to learn a more general task representation.

Finally, introducing multi-episode adaptation could further strengthen performance in complex domains. While current results show strong zero-shot behavior from a single short rollout, some tasks may benefit from incremental refinement over time. The meta-RL literature frequently explores "N-shot" adaptation scenarios, where agents are allowed a few interactions to adjust before final evaluation. Extending the system to support iterative context updates, either by re-encoding across episodes or integrating new data incrementally, would bring it closer to this broader evaluation standard. PEARL already follows this approach by updating its latent context posterior with new episodes, and $RL^2$ naturally adapts through recurrent updates to its internal state. Although CartPole is relatively forgiving, more challenging environments could make a compelling case for such sequential adaptation mechanisms.

## VI. CONCLUSIONS

This work presents a context-conditioned policy architecture for CartPole variants, demonstrating a clear and practical application of context-based meta-reinforcement learning. By combining a trajectory-based encoder with a policy network conditioned on latent task embeddings, the system reflects core ideas from recent meta-RL research, particularly the use of learned task representations to drive fast adaptation. Strong performance on held-out environments, achieved without any test-time gradient updates, highlights the system's ability to generalize through forward inference alone.

The results align closely with broader trends in the field. Approaches such as PEARL and FLAP illustrate the advantages of separating task inference from control, including quicker adaptation, reduced computational overhead at deployment, and greater flexibility in encoding task structure. Although this system is relatively minimal in its design, using REINFORCE and low-dimensional observations, it still captures many of the

benefits observed in more complex frameworks. The fact that it performs well under such constraints underscores the effectiveness of context-based adaptation, even when implemented in a lightweight form.

More fundamentally, the system captures the essence of fast adaptation: recognizing the task at hand and responding appropriately, without the need to re-learn from scratch. The encoder distills critical features from short interaction histories, while the policy behaves as though it were specialized for each new task. This division of labor not only simplifies computation, but also enables consistent, task-sensitive behavior in unfamiliar environments.

From a research standpoint, this case study provides evidence that many of the architectural insights found in state-of-the-art meta-RL approaches can be realized in simplified experimental settings. It also lays the groundwork for meaningful extensions. Directions such as probabilistic embeddings, auxiliary training signals, active exploration strategies, and hybrid adaptation mechanisms represent clear opportunities for scaling the system to more realistic or partially observable domains.

In closing, this study reinforces the potential of context-conditioned meta-learning as a fast and scalable approach to reinforcement learning in variable environments. By synthesizing key ideas from the literature and grounding them in a working prototype, the system serves both as a functional demonstration and a conceptual proof point. With further refinement, architectures of this kind could help bridge the gap between controlled benchmarks and the dynamic, uncertain conditions encountered in real-world applications.

## REFERENCES

[1] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2017.

[2] J. Rothfuss *et al.*, "Promp: Proximal meta-policy search," in *International Conference on Learning Representations (ICLR)*, 2019.

[3] Y. Duan *et al.*, "Rl$^2$: Fast reinforcement learning via slow reinforcement learning," *arXiv preprint arXiv:1611.02779*, 2016.

[4] J. X. Wang *et al.*, "Learning to reinforcement learn," *arXiv preprint arXiv:1611.05763*, 2016.

[5] K. Rakelly, A. Zhou, D. Quillen, C. Finn, and S. Levine, "Efficient off-policy meta-reinforcement learning via probabilistic context variables," in *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019.

[6] L. Zintgraf, K. Shi, K. Hofmann, and S. Whiteson, "Varibad: A very good method for bayes-adaptive deep rl via meta-learning," in *International Conference on Learning Representations (ICLR)*, 2020.

[7] I. Clavera *et al.*, "Model-based reinforcement learning via meta-policy optimization," in *Conference on Robot Learning (CoRL)*, 2018.

[8] M. Peng, B. Zhu, and J. Jiao, "Linear representation meta-reinforcement learning for instant adaptation," *arXiv preprint arXiv:2109.11096*, 2021.

[9] M. Das, *Meta-rl-adaptation: A context-conditioned policy for cartpole variants*, https://github.com/mbd888/meta-rl-adaptation, Accessed: 2025-06-03, 2025.