# SYSTEMATIC REVIEW OF MODEL-BASED REINFORCEMENT LEARNING TECHNIQUES

**Zemzem Hibet**
School of Information Technology and Engineering
Addis Ababa Institute of Technology
Addis Ababa University
`{zumihibet2}@gmail.com`

## ABSTRACT

## 1 INTRODUCTION

Model-Based Reinforcement Learning has emerged as a promising approach that integrates learned models of the environment into decision-making processes. Unlike model-free reinforcement learning methods, which rely solely on trial-and-error interactions with the environment, model-based reinforcement learning leverages predictive models to simulate future states and optimize decision-making strategies. This capability allows Model-Based Reinforcement Learning to achieve greater sample efficiency, making it particularly advantageous in scenarios where data collection is expensive or impractical Janner et al. (2021).

Traditional model-free RL techniques, such as Deep Q-Networks and Proximal Policy Optimization, have demonstrated remarkable success across various domains, including robotics, gaming, and autonomous systems Mnih et al. (2015)Schulman et al. (2017). However, these methods often require extensive interactions with the environment to learn optimal policies, leading to high sample complexity and prolonged training times. In contrast, model-based reinforcement learning aims to mitigate this limitation by constructing an internal representation of the environment, allowing the agent to plan and learn more efficiently. By utilizing models for trajectory rollouts, decision-making, and policy refinement, model-based reinforcement learning has shown improved sample efficiency and generalization capabilities Chua et al. (2018).

Despite these advantages, model-based reinforcement learning is not without its challenges. A key issue is model bias, where inaccuracies in the learned environment model lead to suboptimal decision-making. Compounding errors during multi-step rollouts can degrade policy performance, making long-horizon planning difficult. Additionally, balancing model complexity and computational feasibility is an ongoing challenge, as more sophisticated models require significant computational resources. Moreover, generalization to real-world tasks remains a concern, as policies trained in simulated environments often struggle when deployed in novel settings due to discrepancies between learned and actual dynamics Packer et al. (2019).

Given these considerations, model-based reinforcement learning research has focused on improving model accuracy, reducing compounding errors, and enhancing computational efficiency. Recent developments, such as probabilistic ensembles, latent space models, and hybrid approaches combining model-based and model-free techniques, have sought to address these issues Hafner et al. (2020)Janner et al. (2021). This systematic review aims to critically analyze the advancements in model-based reinforcement learning, evaluating their strengths, limitations, and potential directions for future research.

### 1.1 PROBLEM STATEMENT

While model-based reinforcement learning has demonstrated substantial improvements in sample efficiency and planning, several critical challenges remain unresolved. One major issue is the accuracy of learned models. Inaccurate predictions of environment dynamics can lead to poor decision-

making, particularly in complex or high-dimensional environments. Additionally, errors in the learned model can propagate over time, compounding into significant deviations from the true dynamics. This problem is exacerbated in real-world applications, where discrepancies between simulated and actual environments can hinder the transferability of model-based reinforcement learning policies.

Another challenge is the computational complexity associated with training and utilizing learned models. Many model-based reinforcement learning approaches rely on deep neural networks to model environment dynamics, which can be computationally expensive to train and evaluate. This computational overhead limits the practicality of mdel-based reinforcement learning in resource-constrained scenarios, such as real-time control systems or edge computing applications. Furthermore, the trade-off between model complexity and computational feasibility remains a critical consideration, as overly simplistic models may fail to capture the nuances of the environment, while overly complex models may be computationally prohibitive.

## 1.2 OBJECTIVE

This review critically examines recent advancements in model-based reinforcement learning, highlighting their contributions, limitations, and potential research directions. The goal is to analyze various model-based reinforcement learning techniques, their methodologies, and effectiveness, and to identify existing gaps that need to be addressed for further improvements in this field.

## 2 LITERATURE REVIEW

Recent studies have introduced a variety of innovative methodologies aimed at enhancing the performance of Model-Based Reinforcement Learning. One of the most significant contributions in this area is the Model-Based Policy Optimization framework, proposed by Janner et al. (2021). Model-Based Policy Optimization dynamically adjusts the extent to which the learned model is utilized during policy training, thereby reducing model-induced bias. This is achieved by restricting policy updates to short-horizon rollouts generated from the learned model, which helps mitigate the impact of model inaccuracies. While Model-Based Policy Optimization has demonstrated substantial improvements in sample efficiency compared to traditional model-free baselines, it faces challenges in long-horizon tasks where errors in the learned model accumulate over time, leading to suboptimal policy performance.

Another influential advancement in model-based reinforcement learning is the Probabilistic Ensembles with Trajectory Sampling (PETS) method, introduced by Chua et al. (2018). PETS enhances planning capabilities by incorporating uncertainty-aware modeling through an ensemble of neural network dynamics models. This ensemble approach mitigates overfitting to specific model errors and improves the robustness of decision-making by accounting for predictive uncertainty. However, PETS assumes a Gaussian distribution for predictive uncertainty, which may not fully capture the complex, multimodal uncertainties often encountered in real-world applications. This limitation highlights the need for more sophisticated uncertainty modeling techniques to improve the reliability of model-based reinforcement learning methods in diverse and unpredictable environments.

Recent advancements have also explored the integration of deep generative models into model-based reinforcement learning frameworks. For instance, Janner et al. (2021) proposed diffusion-based planning, which leverages a generative model trained on environment transitions to predict future states more accurately. This approach has shown promising results in high-dimensional control tasks, such as robotic manipulation, due to its ability to generate precise long-horizon predictions. However, the iterative nature of diffusion models introduces significant computational overhead, making them less practical for real-time applications. Similarly, Hafner et al. (2020) introduced DreamerV2, a method that utilizes a latent space model to perform long-horizon planning. DreamerV2 compresses environment observations into a latent representation, enabling efficient imagination of future trajectories. While this approach achieves state-of-the-art performance on benchmark tasks, it remains vulnerable to inaccuracies in the latent model, which can degrade policy performance over extended planning horizons.

In addition to these advancements, researchers have explored hybrid methods that combine model-based and model-free reinforcement learning techniques. For example, Kalweit and Boedecker

introduced Deep Deterministic Policy Gradient with Model-Based Acceleration (DDPG-MB), a framework where model-based rollouts are used to augment model-free training. This hybrid approach improves sample efficiency by leveraging the strengths of both paradigms but faces challenges in maintaining stable and accurate model predictions, particularly in complex environments. Similarly, Amos et al. (2018) developed Differentiable MPC, which integrates model-based planning within deep reinforcement learning architectures. This method enhances interpretability and stability by incorporating model-based constraints into the learning process. However, the computational expense associated with differentiable planning limits its scalability, particularly in resource-constrained scenarios.

Further research has also investigated meta-learning approaches to improve the generalization capabilities of model-based reinforcement learning methods. Clavera et al. (2018) proposed Meta-RL with Probabilistic Ensembles, a framework that enables rapid adaptation to new tasks by leveraging prior experience. This approach uses probabilistic ensembles to model uncertainty and adapts quickly to novel environments, making it particularly useful for tasks requiring generalization. However, meta-learning methods often require extensive meta-training on a diverse set of tasks, which can be computationally intensive and less practical for real-world deployment. Despite these challenges, meta-learning represents a promising direction for enhancing the adaptability and versatility of model-based reinforcement learning methods.

Overall, the literature on model-based reinforcement learning reflects a growing emphasis on improving sample efficiency, robustness, and generalization through innovative methodologies. While significant progress has been made, challenges such as model bias, computational complexity, and sim-to-real transfer remain critical areas for future research. Addressing these challenges will be essential for unlocking the full potential of model-based reinforcement learning in real-world applications, ranging from robotics to autonomous systems.

## 3   DISCUSSION AND ANALYSIS

The literature on rule-based models underscores the delicate balance between interpretability, accuracy, and scalability as a central challenge in their design and application. Gal & Ghahramani (2016) demonstrated that rule extraction from neural networks can achieve high predictive accuracy while maintaining comprehensible rule sets. However, this process often results in rule explosion, a phenomenon where the complexity and volume of extracted rules grow excessively, undermining their interpretability. Rule explosion occurs when models generate an overwhelming number of rules, many of which may be redundant or overly specific, making it difficult for users to understand and trust the model's decision-making process. To address this issue, researchers have developed pruning techniques and rule clustering methods, which aim to eliminate redundant rules and group similar ones together, thereby improving the clarity and usability of the rule sets.

One notable approach to mitigating rule explosion is the use of localized rule extraction methods, such as LoRMIkA Curi et al. (2020). LoRMIkA focuses on extracting rules that are relevant to specific subsets of the data rather than generating global rules that apply to the entire dataset. By concentrating on localized decision contexts, this method reduces the number of rules while preserving their interpretability and relevance. This approach is particularly useful in scenarios where global rules may not capture the nuances of localized patterns, thereby striking a balance between interpretability and accuracy.

Domain generalization is another critical challenge for rule-based models. These models often struggle to transfer knowledge across different domains due to their reliance on domain-specific patterns and rules. To enhance adaptability, researchers have explored predicate logic-based rule generation, which incorporates contextual awareness into the rule formulation process. This approach allows rules to be more flexible and adaptable to varying contexts, improving their performance in cross-domain applications. Additionally, incremental learning techniques have been proposed to enable rule-based models to continuously refine their rule sets as new data becomes available. This capability not only improves the model's adaptability but also ensures that it remains relevant in dynamic environments where data distributions may shift over time.

When comparing rule-based models with non-rule-based approaches, particularly black-box models that rely on interpretability techniques like SHAP (SHapley Additive exPlanations) and LIME (Lo-

cal Interpretable Model-agnostic Explanations), the trade-offs become evident. Rule-based models offer direct interpretability through explicit, human-readable rules, which can be easily understood and validated by domain experts. In contrast, SHAP and LIME provide post-hoc explanations for black-box models, which may lack consistency or accuracy in certain contexts. For instance, post-hoc explanations can vary depending on the input data or the specific implementation of the interpretability technique, leading to potential inconsistencies. Hybrid approaches, such as rule extraction from neural networks, attempt to bridge this gap by combining the interpretability of rule-based models with the predictive power of black-box models. These hybrid methods aim to leverage the strengths of both paradigms, offering a promising direction for achieving both high accuracy and interpretability.

Scalability remains a significant challenge for rule-based systems, particularly when applied to large-scale data environments. As datasets grow in size and complexity, the computational cost of generating and evaluating rules increases, potentially limiting the practicality of rule-based models. To address this issue, researchers have developed optimization strategies, such as distributed rule mining and parallel computation. These techniques enable rule-based models to process large datasets efficiently while maintaining their interpretability. Distributed rule mining, for example, divides the rule extraction process across multiple computational nodes, reducing the time required to generate rules. Similarly, parallel computation leverages modern hardware architectures to accelerate rule evaluation, making rule-based models more scalable for big data applications.

Despite these advancements, several challenges persist in the development and deployment of rule-based models. One of the most pressing issues is the need to balance model complexity, interpretability, and real-time performance. While hybrid models that integrate rule-based reasoning with advanced machine learning techniques offer a promising solution, they often require careful tuning to ensure stability and reliability under dynamic conditions. For example, hybrid models may struggle to maintain consistent performance when the underlying data distribution changes, highlighting the need for robust adaptation mechanisms.

Another area for improvement is the incorporation of domain knowledge, contextual reasoning, and user feedback into the rule generation process. By integrating domain-specific expertise, rule-based models can generate more meaningful and actionable rules that align with real-world requirements. Contextual reasoning allows models to adapt their decision-making processes based on the specific context in which they are applied, enhancing their practical effectiveness. Additionally, incorporating user feedback into the rule generation process can help refine and improve the rules over time, ensuring that they remain relevant and accurate.

In conclusion, while significant progress has been made in addressing the challenges associated with rule-based models, further research is needed to fully realize their potential in real-world applications. By improving the stability of rule-based systems under dynamic conditions, enhancing their interpretability without sacrificing accuracy, and incorporating domain knowledge and user feedback, researchers can develop more effective and reliable rule-based models. These advancements will be critical for enabling the widespread adoption of rule-based systems in domains such as healthcare, finance, and autonomous systems, where interpretability and accuracy are paramount.

## 4 Conclusion and Future Research

Model-Based Reinforcement Learning represents a highly promising paradigm for efficient decision-making, as it integrates learned models of the environment into the decision-making process. By leveraging these models, model-based reinforcement learning methods can achieve significant improvements in sample efficiency and planning capabilities compared to purely model-free approaches. Recent advancements in model-based reinforcement learning, such as probabilistic ensembles, generative models, and hybrid approaches, have demonstrated substantial progress in addressing some of the core challenges in the field. For example, probabilistic ensembles, as seen in methods like Probabilistic Ensembles with Trajectory Sampling (PETS), have enhanced robustness by incorporating uncertainty-aware modeling, while generative models, such as diffusion-based planning and DreamerV2, have improved long-horizon prediction accuracy. Hybrid approaches, which combine model-based and model-free techniques, have further bridged the gap between sample efficiency and policy performance.

Despite these advancements, several challenges continue to hinder the widespread adoption and effectiveness of model-based reinforcement learning. One of the most significant challenges is model bias, where inaccuracies in the learned dynamics model can propagate through the planning process, leading to suboptimal policies. This issue is particularly pronounced in long-horizon tasks, where small prediction errors can compound over time, resulting in significant deviations from the true environment dynamics. Another critical challenge is computational efficiency, as many high-performing MBRL methods rely on complex neural network architectures, which can be computationally expensive to train and evaluate. This limitation is especially problematic in real-time applications or resource-constrained environments. Additionally, generalization remains a persistent issue, as model-based reinforcement learning methods often struggle to transfer policies learned in simulated environments to real-world scenarios due to discrepancies between learned and actual dynamics.

To address these challenges, future research should focus on improving the robustness and generalization of MBRL methods. One promising direction is the development of hybrid approaches that combine the strengths of model-based and model-free techniques. For instance, methods that adaptively switch between model-based planning and model-free policy execution based on the confidence of the learned model could mitigate the impact of model inaccuracies. Luo et al. (2021) proposed such an adaptive framework, where the agent dynamically adjusts its reliance on the model based on its predictive accuracy. This approach not only reduces the risk of model-induced errors but also leverages the efficiency of model-free methods in well-understood regions of the state space.

Another critical area for future research is the development of more expressive and uncertainty-aware modeling techniques. Accurate uncertainty estimation is essential for robust planning, as it allows the agent to account for potential errors in its predictions. Recent advancements in Bayesian deep learning and normalizing flows offer promising opportunities for better capturing epistemic uncertainty in learned dynamics models. For example, Gal & Ghahramani (2016) demonstrated that dropout in neural networks can be used as a Bayesian approximation to estimate model uncertainty. Similarly, normalizing flows, which are capable of modeling complex probability distributions, could be employed to better represent multimodal uncertainties in real-world environments. By incorporating these techniques, MBRL methods can improve their reliability and robustness in decision-making.

Meta-learning also presents a compelling avenue for future research in model-based reinforcement learning. Meta-learning, or learning to learn, enables models to adapt quickly to new tasks with minimal retraining, making it particularly valuable for applications requiring generalization across diverse environments. Finn et al. (2017) introduced Model-Agnostic Meta-Learning (MAML), a framework that has been adapted in model-based reinforcement learning to enable rapid adaptation to novel tasks. By leveraging meta-learning, model-based reinforcement learning methods can become more versatile and capable of handling a wider range of real-world scenarios. For example, meta-learning could be used to train dynamics models that generalize across different robotic platforms or environmental conditions, reducing the need for extensive retraining in new contexts.

In addition to these technical advancements, future research should also explore ways to improve the scalability and practical applicability of model-based reinforcement learning methods. This includes developing lightweight algorithms that maintain high performance while reducing computational overhead, as well as addressing the challenges of sim-to-real transfer. Techniques such as domain randomization and system identification have shown promise in bridging the gap between simulated and real-world environments, but further innovation is needed to achieve reliable performance in real-world applications.

Finally, incorporating domain knowledge and user feedback into the MBRL framework could enhance its practical effectiveness. By integrating domain-specific expertise, model-based reinforcement learning methods can generate more meaningful and actionable policies that align with real-world requirements. User feedback can also play a crucial role in refining and improving the learned models over time, ensuring that they remain relevant and accurate in dynamic environments. While model-based reinforcement learning has made significant strides in recent years, challenges such as model bias, computational efficiency, and generalization remain critical barriers to its widespread adoption. Future research should focus on developing hybrid approaches, uncertainty-aware modeling techniques, and meta-learning strategies to improve the robustness and adaptability of model-based reinforcement learning methods. By addressing these challenges, model-based reinforcement

learning has the potential to revolutionize decision-making in a wide range of applications, from robotics and autonomous systems to healthcare and finance.

## REFERENCES

Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models, 2018. URL `https://arxiv.org/abs/1805.12114`.

Sebastian Curi, Felix Berkenkamp, and Andreas Krause. Efficient model-based reinforcement learning through optimistic policy search and planning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 14156–14170. Curran Associates, Inc., 2020. URL `https://proceedings.neurips.cc/paper_files/paper/2020/file/a36b598abb934e4528412e5a2127b931-Paper.pdf`.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks, 2017. URL `https://arxiv.org/abs/1703.03400`.

Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR.

Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination, 2020. URL `https://arxiv.org/abs/1912.01603`.

Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization, 2021. URL `https://arxiv.org/abs/1906.08253`.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei Rusu, Joel Veness, Marc Bellemare, Alex Graves, Martin Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518:529–33, 02 2015. doi: 10.1038/nature14236.

Charles Packer, Katelyn Gao, Jernej Kos, Philipp Krähenbühl, Vladlen Koltun, and Dawn Song. Assessing generalization in deep reinforcement learning, 2019. URL `https://arxiv.org/abs/1810.12282`.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL `https://arxiv.org/abs/1707.06347`.