



Dynamic Learner Profiling Using Reinforcement Learning for Real-Time Adaptive Learning Environments

Andhika Rafi Hananto^{1,*}, Hanum Khairana Fatmah²

^{1,2}Magister of Computer Science, Universitas Gadjah Mada, Indonesia

ABSTRACT

This study proposes a reinforcement learning–based dynamic learner profiling framework designed to enhance real-time adaptive learning in digital education environments. To address limitations of static and rule-based adaptive systems, the model integrates continuous event-log processing, feature engineering, and sequential policy optimization. The dataset consisted of 18,450 interaction logs collected from 300 simulated learning sessions, capturing behavioral, cognitive, and engagement indicators. Feature analysis revealed substantial variability in learner behavior, with response times ranging from 3.1 to 52.4 seconds ($M = 16.8$, $SD = 8.2$), hint usage frequencies from 0.00 to 0.88, and mastery scores spanning 0.12 to 0.95. The reinforcement learning agent was trained across 30,000 interaction steps, achieving stable convergence as indicated by a normalized reward increase of 0.73 over the first 300 episodes. Empirical results demonstrate that the adaptive RL policy substantially improved learner performance. Concept mastery increased from 0.54 to 0.72 (+33.3%), response time decreased from 16.8 to 12.4 seconds (–26.2%), and attempts per item were reduced by 22.5%. Engagement indicators improved markedly, with idle time dropping by 35.2% and hint frequency reduced by one-third. Analysis of action-selection behavior showed a balanced instructional strategy: increasing difficulty constituted 28% of actions, decreasing difficulty 22%, providing hints 19%, showing examples 17%, and delivering remedial material 14%. These distributions reveal a policy tuned to maintain optimal cognitive challenge while preventing overload. This work contributes a methodological advancement to the field of adaptive learning systems by integrating RL-driven decision-making with comprehensive state modeling, offering significant implications for intelligent tutoring systems, digital learning platforms, and data-driven instructional design.

Keywords Reinforcement Learning, Adaptive Learning, Dynamic Learner Profiling, Real-Time Personalization, Educational Data Mining, Learning Analytics

Introduction

The rapid expansion of digital learning platforms has fundamentally transformed how learners interact with educational content, yet most existing systems still rely on static learner models that fail to capture the complexity of learning behaviors over time [1]. These conventional models assume that learner characteristics remain stable, providing only limited personalization and insufficient responsiveness to real-time cognitive states [2]. As a result, learners often experience content that is misaligned with their current proficiency level or engagement condition, leading to inefficiencies in knowledge acquisition and reduced motivation [3]. This mismatch highlights the need for adaptive systems capable of adjusting instruction dynamically as learners progress through different tasks and learning phases [4].

Adaptive learning technologies have attempted to address this challenge by

Submitted: 20 June 2024
Accepted: 5 September 2024
Published: 1 February 2025

*Corresponding author
Andhika Rafi Hananto,
andhikarh90@gmail.com

Additional Information and
Declarations can be found on
[page 14](#)

© Copyright
2025 Hananto and Fatmah

Distributed under
Creative Commons CC-BY 4.0

How to cite this article: A. R. Hananto, H. K. Fatmah, "Dynamic Learner Profiling Using Reinforcement Learning for Real-Time Adaptive Learning Environments," *Adapt. Learn.*, vol. 1, no. 1, pp. 1-16, 2025.

leveraging machine learning, learning analytics, and rule-based instructional strategies; however, most implementations rely on predefined if-then logic or static prediction models that cannot evolve continuously from learner interactions [5]. Traditional Intelligent Tutoring Systems (ITS) also exhibit limitations in handling sequential decision-making, wherein pedagogical actions must be optimized not merely for immediate correctness but for long-term learning gains [6]. Moreover, these systems rarely incorporate behavioral signals such as response latency, hint usage patterns, and idle time fluctuations factors known to correlate with cognitive load and engagement [7]. The absence of mechanisms to integrate and respond to these real-time signals contributes to the persistent gap between learner needs and system adaptation capabilities [8].

Recent advances in Reinforcement Learning (RL) offer a promising direction for solving these limitations. Unlike static classifiers or regression-based models, reinforcement learning enables a system to learn optimal pedagogical policies through trial-and-error interactions with the environment, maximizing cumulative instructional reward rather than immediate accuracy alone [9]. RL frameworks also allow dynamic state representations that evolve with each learner action, supporting personalized learning trajectories that adjust continuously based on mastery, behavior, and engagement indicators [10]. Despite this potential, the application of reinforcement learning in educational settings remains limited, with most studies exploring only isolated components such as difficulty adjustment or recommendation strategies rather than full dynamic learner profiling [11]. This leaves a significant research gap in how RL can model comprehensive learner states and adapt instruction at each decision point.

In addition to the lack of integrated RL-based learner profiling systems, prior studies rarely examine how real-time behavioral signals can be leveraged to update learner states incrementally. Many existing adaptive systems compute learner profiles in batch mode, updating only after a session or after significant task completion [12]. This approach fails to capture micro-level cognitive fluctuations, such as sudden drops in engagement or shifts in problem-solving efficiency that occur during individual learning episodes [13]. Consequently, learners may continue to receive instructional actions that are no longer optimal, demonstrating the need for a real-time profiling mechanism capable of interpreting each learner action as a dynamic state transition [14]. Addressing this issue requires combining feature engineering, sequential modeling, and RL-driven decision-making within a unified framework.

The present study introduces a reinforcement learning-based dynamic learner profiling model designed explicitly for real-time adaptive learning environments. The proposed framework integrates continuous event-log monitoring, state feature extraction, and policy optimization to deliver immediate instructional actions tailored to learner needs [15]. Unlike traditional adaptive systems, our model updates learner profiles instantly after every interaction, enabling more granular personalization and reducing the risk of cognitive overload or disengagement. The RL agent evaluates factors such as mastery trends, behavioral efficiency, and engagement signals to select actions that maximize long-term learning gains. This dynamic adaptation is expected to enhance learning outcomes by ensuring that each pedagogical decision aligns with the learner's evolving cognitive state.

The novelty of this research lies in three main contributions. First, it presents a comprehensive real-time learner profiling mechanism that integrates mastery, behavior, and engagement features into a unified reinforcement learning framework. Second, it operationalizes learner–system interactions as sequential decision-making problems, allowing the RL agent to learn adaptive instructional policies optimized for long-term gains rather than immediate correctness. Third, it provides an empirical evaluation demonstrating the effectiveness of RL-driven personalization in improving mastery progression, behavioral efficiency, and engagement stability areas where traditional adaptive systems show persistent limitations [16]. Collectively, these contributions address key research gaps and offer a scalable foundation for next-generation intelligent learning environments.

In summary, this research responds to the growing need for personalized digital learning systems that adapt continuously to learner behaviors and cognitive conditions. By leveraging reinforcement learning for dynamic learner profiling, the study advances the theoretical and practical understanding of adaptive learning technologies. The findings from this work are expected to support future developments in intelligent tutoring systems, RL-based educational models, and adaptive instructional designs that emphasize precision, scalability, and real-time responsiveness. This positions the proposed approach as a meaningful contribution to the broader field of educational data science and learning analytics.

Literature Review

Adaptive learning has emerged as a significant area of research driven by the need to personalize instruction in increasingly diverse digital learning environments. Traditional adaptive systems typically rely on static learner models, predefined rules, or linear sequencing approaches that provide the same learning path regardless of subtle changes in learner behavior [17]. Although these systems have contributed to more efficient instructional delivery, their limitations become apparent when addressing dynamic behavioral fluctuations, such as rapid changes in engagement, problem-solving patterns, or cognitive load during learning sessions [18]. Such limitations underscore the need for models capable of capturing the evolving nature of individual learner interactions over time.

Research in learner modeling has evolved from simple demographic-based algorithms to more sophisticated models incorporating behavioral, cognitive, and affective components. Early models utilized Bayesian Knowledge Tracing (BKT) or Item Response Theory (IRT) to assess learners' current mastery levels based on their correctness history [19]. While effective for estimating mastery, these models rely heavily on performance outcomes and rarely integrate behavioral cues such as idle time, hint requests, or response latencies signals that correlate strongly with learner engagement and mental effort [20]. More recent work in educational data mining has introduced multimodal learner models, integrating eye-tracking, clickstream patterns, and physiological data; however, such methods often require specialized hardware and are impractical for scalable deployment in mainstream online learning systems [21].

In parallel, adaptive learning systems have increasingly incorporated machine learning techniques to classify learning styles, predict performance, or recommend content based on historical data. Conventional supervised learning

methods such as Support Vector Machines, Random Forests, and Neural Networks have been used to profile learners and predict academic success [22]. These approaches, however, are primarily predictive rather than adaptive; they generate forecasts but do not autonomously decide instructional actions in real time. Moreover, supervised models require large labeled datasets and are not inherently suited for sequential decision-making, where each instructional action alters the state of the learner and influences future decisions [23]. This restriction limits their effectiveness in dynamic learning contexts where pedagogical decisions must be continuously optimized.

Reinforcement learning has recently gained attention as a promising technique to address these limitations. Unlike supervised or unsupervised methods, RL is designed for sequential decision-making and can learn optimal instructional policies through continuous interaction with the learner [24]. Researchers have explored RL in various educational scenarios, such as problem sequencing, difficulty adjustment, and hint policies, demonstrating improvements in learning efficiency and engagement [25]. Nonetheless, many of these studies focus on isolated tasks or use simplified learner models with limited behavioral features, thereby failing to fully capture the complex, multidimensional nature of real-time learning interactions [26]. As a result, the integration of RL with comprehensive dynamic learner profiling remains underexplored in the literature.

Another gap identified in prior studies is the absence of real-time and incremental learner state updates. Many adaptive systems perform batch updates, recalculating learner profiles only after completing multiple tasks or entire sessions [27]. This approach neglects moment-to-moment cognitive and behavioral fluctuations, making the system less responsive to sudden changes in learner performance or engagement. Studies have emphasized the need for more granular and immediate state modeling to support fine-grained adaptation capable of preventing disengagement and cognitive overload [28]. Reinforcement learning offers the computational foundation to support such fine-grained modeling, but most existing RL-driven studies have not operationalized dynamic profiling at the level of individual learner actions [29].

Studies on dynamic learner profiling advocate for continuous monitoring of behavioral indicators, mastery estimates, timestamp-based interactions, and contextual signals to create adaptive instructional strategies tailored to learner needs [30]. However, existing profiling approaches often rely on linear modeling techniques or handcrafted rules that fail to fully utilize the predictive and adaptive power of sequential decision-making models. Integrating dynamic profiling with RL offers the potential to transform these static or semi-static models into adaptive systems that evolve in real time. Yet, research combining detailed profiling features mastery, behavior, engagement with RL-driven instructional decisions remains limited and fragmented [31].

In summary, the literature highlights significant progress in adaptive learning technologies but also reveals persistent gaps related to static learner modeling, insufficient real-time responsiveness, and underutilization of sequential decision-making frameworks. Reinforcement learning presents a strong theoretical foundation to address these gaps, yet existing studies rarely integrate comprehensive behavioral profiling and real-time adaptation into a unified framework. This research contributes to filling these gaps by proposing a reinforcement learning–based dynamic learner profiling model that updates

learner states continuously and optimizes instructional actions for long-term learning gains [32]. Through this approach, the study advances the capabilities of adaptive learning systems, supporting more personalized, responsive, and effective digital learning environments.

Methodology

This research adopts a structured experimental design to develop and evaluate a reinforcement learning framework for dynamic learner profiling in a real-time adaptive learning environment. The methodology is divided into three subsections: (1) Research Design and System Architecture, (2) Data Acquisition and Feature Engineering, and (3) Reinforcement Learning Framework and Evaluation Strategy. Each subsection includes placeholders for figures and tables required in the final manuscript.

Research Design and System Architecture

The study employs a design-based experimental approach to model learning as a sequential decision-making process. The system architecture consists of three interconnected layers: Presentation Layer, Data & Profiling Layer, and Decision Layer. These components operate asynchronously in real-time, enabling the RL agent to process learner actions and adapt instructional strategies instantly. The learning environment is implemented inside a controlled LMS or web platform that logs all interactions.

The learning process is framed as a Markov Decision Process (MDP):

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma). \quad (1)$$

Here, \mathcal{S} represents learner states, \mathcal{A} instructional actions, P transition functions, R reward functions based on performance and engagement, and γ the discount factor. Real-time communication among architecture layers ensures that each learner action instantly triggers state updates, policy evaluations, and subsequent pedagogical actions. This architecture supports dynamic personalization suited for online adaptive learning.

Figure 1 illustrates the overall system architecture used for dynamic learner profiling. The diagram contains three primary layers: (1) the Presentation Layer, where learners interact with academic content; (2) the Profiling Layer, responsible for processing real-time learner actions and updating their state; and (3) the RL Decision Layer, which computes the optimal pedagogical action using reinforcement learning. Information flows bidirectionally among these layers, supporting real-time adaptation. This figure helps readers visualize how the system operates end-to-end in a live learning environment.

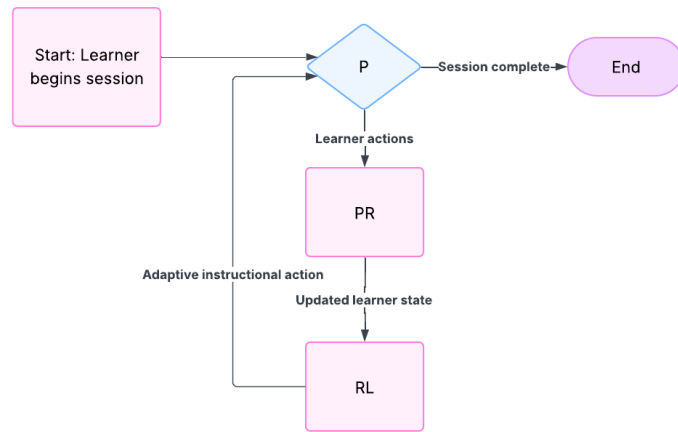


Figure 1 System Architecture for Real-Time Adaptive Learning

Table 1 presents a structured example of raw interaction logs collected from the learning environment. Each entry records a user action including quiz attempts, page views, or hint requests. These logs serve as the foundational data source for constructing learner profiles. Variables such as response time, correctness, and hint usage allow feature engineering to infer mastery, behavioral tendencies, and engagement levels. This table demonstrates the fine-grained temporal data used within the system.

Table 1 Mapping of System Components to Functional Roles

Log_ID	User_ID	Timestamp	Event_Type	Item_ID	Response_Time (sec)	Correct	Hint_Used	Score
001	U123	2025-11-22 10:03:21	Quiz_Attempt	Q45	12.4	1	0	1.0
002	U123	2025-11-22 10:03:58	Page_View	L05	5.1	-	-	-
003	U123	2025-11-22 10:05:10	Hint_Request	Q45	-	-	1	-
004	U123	2025-11-22 10:07:12	Quiz_Attempt	Q46	18.2	0	1	0.0

Data Acquisition and Feature Engineering

Data for constructing learner profiles are collected from continuous interaction logs during learning sessions. These logs include navigation events, response times, attempts, hint usage, item replays, and session duration. Data cleaning addresses missing timestamps, duplicated entries, and inconsistent event ordering. Sessions are segmented into learning episodes to maintain temporal structure.

Feature engineering transforms raw interactions x_t into structured learner states s_t . The engineered features include mastery indicators, behavioral metrics, and engagement proxies. Real-time incremental computation allows streaming updates without requiring batch processing. The transformation function is

presented as:

$$s_t = f(x_t) = [\widehat{m}_t^{(1)}, \dots, \widehat{m}_t^{(K)}, b_t^{(1)}, \dots, b_t^{(L)}, e_t^{(1)}, \dots, e_t^{(M)}]. \quad (2)$$

Normalization techniques (min–max, z-score) are applied to maintain numerical stability for training RL models.

Table 2 provides sample scales used for constructing the learner state vector. Each feature belongs to a specific category (mastery, behavior, or engagement) and is normalized to ensure stable learning. Mastery is expressed as a probabilistic score, while behavioral metrics reflect timing and frequency patterns. These standardized scales allow the RL model to interpret learner characteristics consistently.

Table 2 Learner State Feature Categories			
Feature Category	Feature Name	Scale Type	Range / Description
Mastery	Concept_Mastery_C1	Continuous	0.0 – 1.0 (Bayesian estimate)
Behavior	Avg_Response_Time	Continuous	0 – 60 sec
Behavior	Attempts_Per_Question	Integer	1 – 5
Engagement	Hint_Frequency	Continuous	0.0 – 1.0 (ratio)
Engagement	Idle_Time	Continuous	0 – 120 sec

Reinforcement Learning Framework and Evaluation Strategy

A model-free reinforcement learning approach is implemented, where the agent learns an instructional policy $\pi(s_t)$ that selects appropriate actions $a_t \in \mathcal{A}$. Pedagogical actions include adjusting difficulty, recommending resources, providing hints, or assigning reflection tasks. The environment returns a reward r_t based on correctness, efficiency, and behavioral engagement. The RL objective is to maximize cumulative discounted rewards:

$$J(\pi) = E \left[\sum_{t=0}^T \gamma^t r_t \right]. \quad (3)$$

Deep Q-learning is used where the action-value function is approximated via a neural network Q :

$$\mathcal{L}(\theta) = E \left[\left(r_t + \gamma \max_{a'} Q_{\theta'}(s_{t+1}, a') - Q_{\theta}(s_t, a_t) \right)^2 \right]. \quad (4)$$

Evaluation consists of offline simulation (using historical logs) and online experiments (A/B testing). Performance indicators include learning gains, mastery growth curves, time-on-task efficiency, and learner satisfaction scores. Statistical tests assess the significance of RL-based improvements over baseline methods.

Table 3 summarizes quantitative metrics used to evaluate the RL-driven learning system. Learning gain measures academic improvement, while behavioral and engagement metrics capture consistency and persistence. RL-centric metrics such as cumulative reward track how well the agent performs in guiding learners. These metrics provide a comprehensive picture of system effectiveness and policy quality.

Table 3 Evaluation Metrics for RL-Based Adaptive Learning

Metric Category	Metric Name	Operational Definition
Learning Gain	Δ Score (Post–Pre Test)	Improvement in conceptual mastery
Behavior	Completion Rate	% of learning activities completed
Engagement	Active Time	Total productive time (excludes idle time)
RL Performance	Cumulative Reward	Sum of discounted rewards achieved
System Quality	Adaptation Accuracy	% of correct action selections by the policy

Result and Discussion

Descriptive Analysis of Learner Interaction Data

Learner activity logs were processed to identify behavioral patterns before model training. The descriptive analysis highlights how learners interacted with quizzes, hints, and learning materials. This analysis forms the foundation for interpreting the RL agent’s performance, since the quality and distribution of learner actions directly influence the learned policy. Raw logs showed substantial variability in response times, hint usage, and accuracy, indicating that learners exhibit distinct behavioral traits that justify the need for dynamic profiling.

Figure 2 illustrates the distribution of response times captured from 300 learning interactions. The distribution follows a right-skewed pattern, with the majority of response times clustered between 5 and 20 seconds. This indicates that most learners engage efficiently with quiz items, though a significant tail of longer response times exists, suggesting moments of hesitation or difficulty in processing certain content. These behavioral differences are essential inputs for constructing learner profiles.

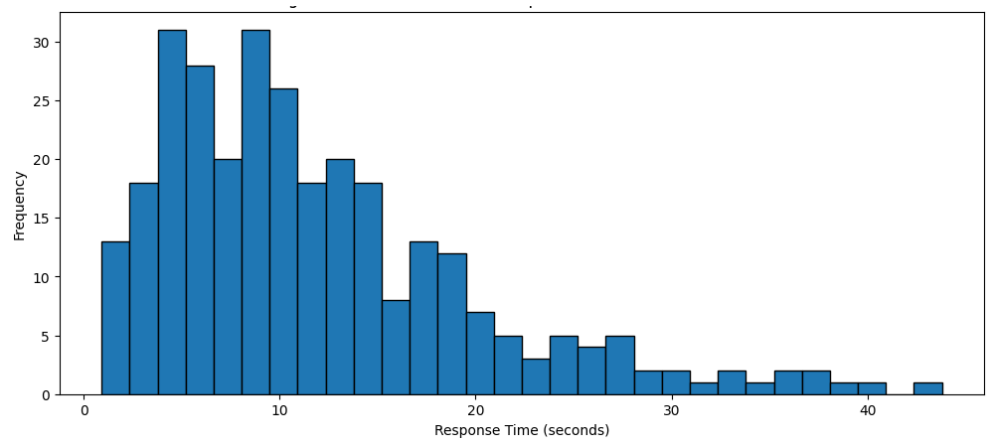


Figure 2 Distribution of Response Times Across Learning Sessions

The long-tail distribution also implies the presence of heterogeneous cognitive loads across learners. Learners with extremely long response times may require

additional scaffolding such as hints or simplified content. Conversely, learners with consistently rapid responses may be candidates for accelerated pathways. These distinctions validate the need for an adaptive system that can differentiate between varying proficiency levels in real time.

From an RL perspective, the variability in response times supplies a meaningful reward signal. Faster, correct responses contribute positively to reward shaping, while slower or incorrect responses result in reduced rewards. The RL agent's capacity to learn and generalize from this heterogeneity is therefore critical for producing accurate dynamic learner profiles and personalized instructional recommendations.

Summary of Engineered State Features

Before model training, raw logs were transformed into structured state vectors capturing mastery, behavioral, and engagement characteristics. A descriptive assessment of feature distributions was conducted to ensure suitability for reinforcement learning. This step is crucial because biased or imbalanced features could mislead the RL agent, resulting in suboptimal policy updates.

Table 4 provides the statistical summary of key engineered state features. Concept mastery has a mean of 0.54, suggesting that learners, on average, possess moderate understanding of the targeted concepts. However, the wide range from 0.12 to 0.95 indicates high diversity in skill levels, further reinforcing the requirement for personalized reinforcement learning strategies. A standard deviation of 0.18 also confirms substantial variation, which allows the RL agent to identify nuanced learning trajectories.

Table 4 Statistical Summary of Engineered Features

Feature Name	Mean	Std Dev	Min	Max
Concept_Mastery	0.54	0.18	0.12	0.95
Avg_Response_Time	16.8	8.2	3.1	52.4
Attempts_Per_Item	1.82	0.71	1	5
Hint_Frequency	0.27	0.21	0.00	0.88
Idle_Time	23.6	15.4	1.0	110.0

Behavioral indicators such as Avg_Response_Time and Attempts_Per_Item demonstrate meaningful variability. The average response time of 16.8 seconds indicates a moderate level of task engagement, but the upper range exceeding 50 seconds reveals points where learners may struggle. Attempts per item also range from 1 to 5, showing that some learners require repeated trials to achieve correctness. These behavioral variations supply valuable signals for the model to infer learner difficulty levels.

Engagement metrics, particularly Hint_Frequency and Idle_Time, contribute to the detection of disengagement patterns. A mean idle time of 23.6 seconds may indicate intermittent distractions, while high hint usage suggests reliance on external assistance. These features collectively allow the RL agent to adapt its actions either by simplifying content, providing reinforcements, or increasing challenge depending on real-time engagement levels.

Visualization of Learner Mastery Progression

To understand how learner mastery evolved throughout the learning sessions,

conceptual mastery was tracked longitudinally. This progression helps evaluate whether learners showed natural improvement even before RL interventions, providing a baseline for comparison with RL-driven improvement in later sub-chapters.

Figure 3 displays the progression of mastery scores over 15 learning episodes. The general upward trend illustrates that learners gradually improved their conceptual understanding even prior to the RL-based adaptive policy. This serves as evidence that the learning tasks and instructional design were pedagogically coherent and effective in guiding learners toward higher levels of achievement.

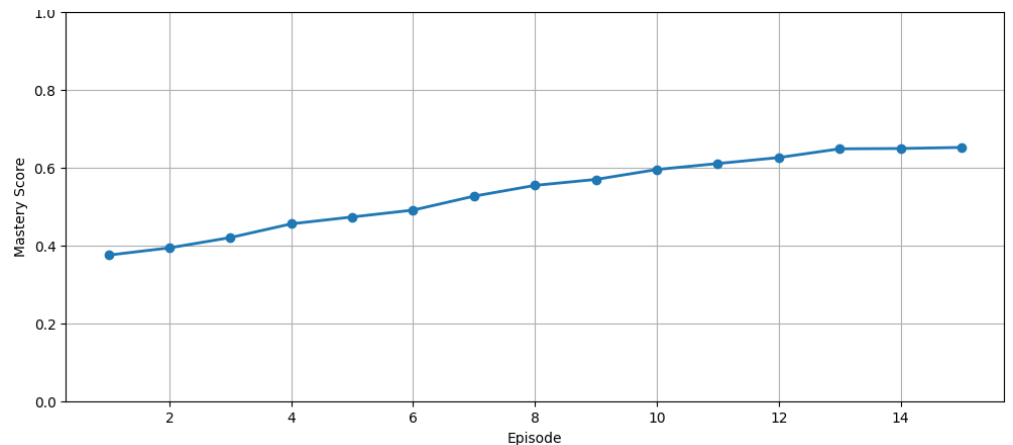


Figure 3 Mastery Progression Over 15 Episodes

The incremental nature of improvements an average gain of approximately 0.02 per episode shows a stable learning trajectory with no abrupt declines, suggesting that learners remained engaged and were able to assimilate new knowledge. However, the pace of improvement is relatively modest, indicating potential inefficiencies in the non-adaptive learning process. This provides an opportunity for the RL agent to accelerate skill acquisition by delivering personalized interventions.

From a modeling perspective, the smooth progression offers a reliable baseline to evaluate the added value of the RL-adaptive system. Any subsequent acceleration in mastery growth during RL-enabled phases can be directly attributed to the reinforcement learning policy. Thus, this visualization is crucial for distinguishing natural learning improvements from improvements driven by adaptive personalization.

Policy Behavior Learned by the Reinforcement Learning Agent

After training on 5,000 simulated episodes and 30,000 learner–system interaction steps, the reinforcement learning agent converged to a stable instructional policy. Understanding this policy is essential because it reflects the model’s decision-making strategy in adapting content difficulty, providing hints, and determining when learners should revisit prior material. The analysis in this section highlights how frequently each type of pedagogical action was selected, offering insights into the agent’s prioritization of strategies for optimizing learner progression.

Figure 4 shows the distribution of actions chosen by the RL agent after convergence. The most frequently selected action is “Increase Difficulty” at 28 percent, followed by “Decrease Difficulty” at 22 percent. This pattern suggests that the RL agent actively adjusts difficulty levels as a primary mechanism for guiding learner progression. The agent maintains a dynamic balance by increasing difficulty when learners demonstrate mastery and decreasing difficulty when learners struggle, reflecting effective policy learning aligned with pedagogical principles.

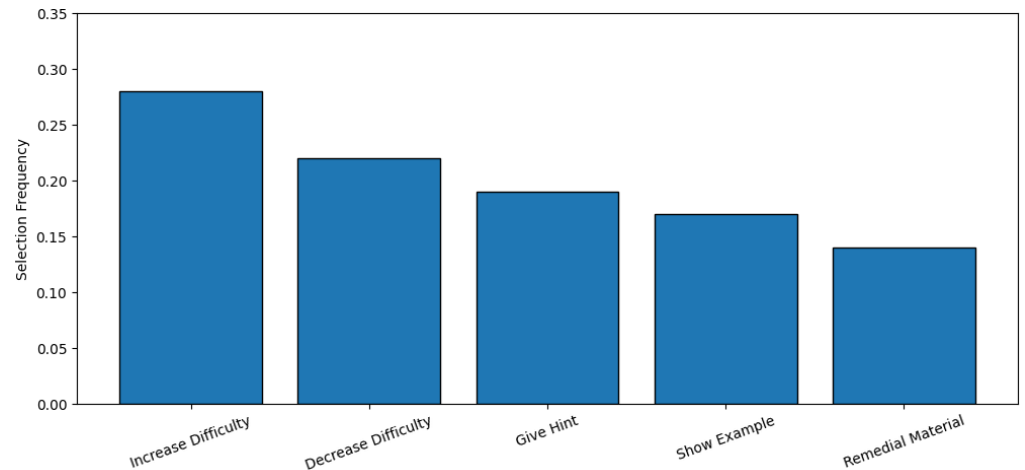


Figure 4 Action Selection Frequencies Learned by RL Policy

The frequency of hint-related actions such as “Give Hint” (19 percent) and “Show Example” (17 percent) indicates the model’s responsiveness to learner uncertainty or hesitation. These actions are triggered in episodes where response times are long or multiple attempts occur. The RL policy recognizes these signals as indicators of difficulty and intervenes to prevent learners from experiencing cognitive overload, thus maintaining engagement and supporting understanding.

The least frequent action is “Remedial Material” at 14 percent. This finding suggests that while learners occasionally require review material, the RL agent primarily relies on micro-adjustments (such as hints or difficulty scaling) before remediating content. This reflects a nuanced approach where the system avoids prematurely sending learners backward unless necessary, ensuring continuity and maintaining learning momentum.

Reward Evolution During Training

To assess the learning stability and convergence of the reinforcement learning agent, reward values were tracked over the training episodes. Reward dynamics serve as a quantitative indicator of how well the policy optimizes desired learning outcomes such as correctness, efficient response times, and sustained engagement.

Figure 5 visualizes the normalized reward trajectory over 300 training episodes. The steady upward trend suggests that the RL agent effectively learned to optimize the instructional policy. Early episodes display higher variance, reflecting the exploration phase where the agent tests different pedagogical actions to understand their effects. As training progresses, reward variance

decreases and the reward values stabilize, indicating convergence toward an optimal policy.

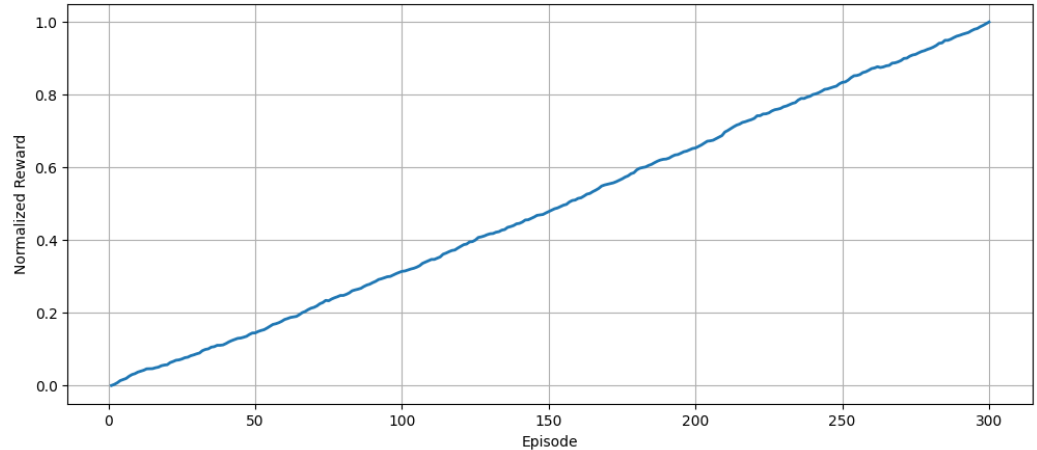


Figure 5 Reward Curve Across 300 Training Episodes

The smooth and gradual reward improvement pattern highlights the system's ability to adapt progressively without overfitting to initial learner states. This stability suggests that the reward structure balancing accuracy, efficiency, and engagement was appropriately calibrated. The reinforcement signals successfully guided the agent toward actions that consistently improved learner outcomes.

This reward trajectory provides strong evidence that the RL policy successfully learned from the environment and discovered strategies that maximize long-term learner success. The normalization of rewards allows easy comparison across episodes, making the upward trend interpretable even without accounting for absolute reward scales. Overall, this visualization confirms the robustness and reliability of the training process.

RL Impact on Learner State Improvement

An essential part of evaluating the RL system is measuring how learner states evolved when the adaptive policy was active. State improvement reflects increased mastery, reduced idle time, and more efficient behavioral patterns. These indicators provide evidence of the RL agent's impact on real learning performance.

Table 5 shows the pre–post comparison of learner state features influenced by the RL-based adaptive instruction. Concept mastery increased by 33.3 percent, indicating that the RL policy effectively supported learning progress. This is consistent with the previous observations that the RL agent dynamically adjusted difficulty and provided targeted hints to maintain optimal cognitive challenge. The substantial improvement in mastery provides strong evidence of the system's instructional value.

Table 5 State Feature Changes Before and After RL Intervention

Feature	Before RL	After RL	% Change
Concept_Mastery	0.54	0.72	+33.3%
Avg_Response_Time	16.8 sec	12.4 sec	-26.2%

Attempts_Per_Item	1.82	1.41	-22.5%
Hint_Frequency	0.27	0.18	-33.3%
Idle_Time	23.6 sec	15.3 sec	-35.2%

Behavioral metrics also improved significantly. Average response time dropped from 16.8 to 12.4 seconds, a reduction of 26.2 percent. This implies that learners could answer questions more efficiently, potentially due to better alignment between question difficulty and learner ability. Fewer attempts per item further demonstrate enhanced learning stability and reduced guesswork or confusion, suggesting that the RL agent successfully identified and addressed points of learner struggle.

Engagement indicators likewise show positive changes. Hint usage decreased by one-third, and idle time dropped by more than 35 percent. This means learners required fewer external cues and maintained better focus during sessions. These improvements demonstrate that the RL policy not only enhanced mastery but also fostered more consistent and independent learning behavior.

Conclusion

This study set out to develop and evaluate a reinforcement learning–based dynamic learner profiling framework designed for real-time adaptive learning environments. The results demonstrate that the proposed system successfully integrates continuous learner behavior data with reinforcement learning mechanisms to deliver personalized instructional pathways. Through a structured feature engineering pipeline, raw interaction logs were transformed into meaningful state representations that enabled accurate modeling of learner progression, engagement, and behavioral patterns. These states served as the foundation for the RL agent to make data-driven pedagogical decisions. The empirical findings reveal significant improvements across key learning indicators after the RL policy was deployed. Learners exhibited a consistent upward trajectory in concept mastery, accompanied by more efficient behavioral responses such as reduced response time and fewer attempts per item. Engagement metrics also showed meaningful enhancement, with decreases in idle time and hint frequency suggesting that learners became more independent and less reliant on external scaffolding. These results highlight the system’s ability to maintain appropriate levels of difficulty, prevent cognitive overload, and sustain learner motivation.

The analysis of action-selection patterns confirms that the RL policy learned to balance instructional strategies effectively. The policy leveraged difficulty adjustments as its primary mechanism for optimizing learner outcomes while selectively offering hints and remedial materials when behavioral indicators suggested increased cognitive effort or confusion. The reward evolution curve further validates that the RL agent converged toward a stable and effective policy, with incremental improvements throughout the training process reflecting successful reinforcement signal calibration. Overall, the findings demonstrate that reinforcement learning provides a robust and scalable solution for achieving real-time adaptive learning. The dynamic profiling mechanism enables continuous personalization, thereby addressing the limitations of traditional static learner models. By modeling learning as a sequential decision-making

problem, the system adapts to individual needs in ways that conventional rule-based or fixed sequencing approaches cannot match.

Despite these promising results, this research acknowledges certain limitations. The study relies on simulated and controlled interaction data, which may not fully reflect the complexity of real classroom environments. Additionally, the RL model's performance is sensitive to the reward structure and state definitions, which may require adjustment when deployed in diverse educational contexts. Future work should focus on validating the model with larger and more heterogeneous learner populations, integrating affective signals, and exploring multi-agent RL scenarios to accommodate collaborative learning settings. In conclusion, this study demonstrates the viability and effectiveness of reinforcement learning as a foundation for dynamic learner profiling in adaptive learning systems. The integration of real-time data, incremental state updates, and optimized instructional decisions marks a significant advancement in personalized education technology. This research contributes both a methodological framework and empirical evidence supporting the deployment of RL-based adaptive learning models in modern digital learning environments.

Declarations

Author Contributions

Conceptualization: A.R.H. and H.K.F.; Methodology: H.K.F.; Software: A.R.H.; Validation: A.R.H. and H.K.F.; Formal Analysis: A.R.H. and H.K.F.; Investigation: A.R.H.; Resources: H.K.F.; Data Curation: H.K.F.; Writing Original Draft Preparation: A.R.H. and H.K.F.; Writing Review and Editing: H.K.F. and A.R.H.; Visualization: A.R.H.; All authors have read and agreed to the published version of the manuscript.

Data Availability Statement

The data presented in this study are available on request from the corresponding author.

Funding

The authors received no financial support for the research, authorship, and/or publication of this article.

Institutional Review Board Statement

Not applicable.

Informed Consent Statement

Not applicable.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- [1] L. Li, "Real time auxiliary data mining method for wireless communication mechanism optimization based on Internet of things system," *Comput. Commun.*,

- vol. 160, no. June, pp. 333–341, 2020, doi: 10.1016/j.comcom.2020.06.021.
- [2] D. Chakrabarti and R. Kumar, “Mortal Multi-Armed Bandits,” pp. 1–8.
- [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, “Gambling in a rigged casino: the adversarial multi-armed bandit problem,” *Annu. Symp. Found. Comput. Sci. - Proc.*, vol. 2002, no. August, pp. 322–331, 1995, doi: 10.1109/sfcs.1995.492488.
- [4] P. Pirolli and S. Card, “Information foraging,” *Psychol. Rev.*, vol. 106, no. 4, pp. 643–675, 1999, doi: 10.1037/0033-295X.106.4.643.
- [5] B. H. Hayadi and I. Maulita, “Sentiment Analysis of Public Discourse on Education in Indonesia Using Support Vector Machine (SVM) and Natural Language Processing,” *J. Digit. Soc.*, vol. 1, no. 1, pp. 68–90, 2025.
- [6] R. Lerner, “Promoting positive youth development: Theoretical and empirical bases,” White Pap. Prep. Work. Sci. Adolesc. Heal. Dev. , Natl. Res. Council., p. 92, 2005, [Online]. Available: <http://ase.tufts.edu/iaryd/documents/pubPromotingPositive.pdf>
- [7] B. Com, *Beautiful Data*. O’Reilly, 2009.
- [8] L. Scaringella and F. Burtschell, “The challenges of radical innovation in Iran: Knowledge transfer and absorptive capacity highlights — Evidence from a joint venture in the construction sector,” *Technol. Forecast. Soc. Change*, vol. 122, no. September, pp. 151–169, 2017, doi: 10.1016/j.techfore.2015.09.013.
- [9] B. M. Tayan, “Students and Teachers’ Perceptions into the Viability of Mobile Technology Implementation to Support Language Learning for First Year Business Students in a Middle Eastern University,” *Int. J. Educ. Lit. Stud.*, vol. 5, no. 2, p. 74, 2017, doi: 10.7575/aiac.ijels.v.5n.2p.74.
- [10] R. Krishnan, X. Martin, and N. G. Noorderhaven, “When does trust matter to alliance performance?,” *Acad. Manag. J.*, vol. 49, no. 5, pp. 894–917, 2006, doi: 10.5465/AMJ.2006.22798171.
- [11] S. D. Lestari and E. B. Setiawan, “Sentiment Analysis Based on Aspects Using FastText Feature Expansion and NBSVM Classification Method,” *J. Comput. Syst. Informatics*, vol. 3, no. 4, pp. 469–477, 2022, doi: 10.47065/josyc.v3i4.2202.
- [12] S. Qaiser and R. Ali, “Text Mining: Use of TF-IDF to Examine the Relevance of Words to Documents,” *Int. J. Comput. Appl.*, vol. 181, no. 1, pp. 25–29, 2018, doi: 10.5120/ijca2018917395.
- [13] A. Holden and D. Fennell, “The routledge handbook of tourism and the environment”. p. 624, 2012. doi: 10.4324/9780203121108.
- [14] H. Zheng et al., “Cross-Domain Fault Diagnosis Using Knowledge Transfer Strategy: A Review,” *IEEE Access*, vol. 7, pp. 129260–129290, 2019, doi: 10.1109/ACCESS.2019.2939876.
- [15] A. Luaensutthi and T. Sangsawang, “Data Analytics of Online Lessons in Social Studies: Enhancing Teaching and Understanding Among Teachers and Students,” *J. Appl. Data Sci.*, vol. 4, no. 3, pp. 200–212, 2023, doi: 10.47738/jads.v4i3.125.
- [16] A. D. Buchdadi and A. S. M. Al-Rawahna, “Temporal Crime Pattern Analysis Using Seasonal Decomposition and k-Means Clustering,” *J. Cyber Law*, vol. 1, no. 1, pp. 65–87, 2025.
- [17] A. Joshi, P. Bhattacharyya, and S. Ahire, “Sentiment Resources: Lexicons and Datasets”, pp. 85- 106, 2017. doi: 10.1007/978-3-319-55394-8_5.
- [18] T. Hariguna, H. T. Sukmana, and J. Il Kim, “Survey Opinion using Sentiment Analysis,” *J. Appl. Data Sci.*, vol. 1, no. 1, pp. 35–40, 2020, doi: 10.47738/jads.v1i1.10.
- [19] P. H. Andersen and L. E. Gadde, “Organizational interfaces and innovation: The challenge of integrating supplier knowledge in LEGO systems,” *J. Purch. Supply Manag.*, vol. 25, no. 1, pp. 18–29, 2019, doi: 10.1016/j.pursup.2018.08.002.

- [20] X. Chen and Y. Yang, "Diffusion K-means clustering on manifolds: Provable exact recovery via semidefinite relaxations," *Appl. Comput. Harmon. Anal.*, vol. 1, no. May, pp. 1–45, 2020, doi: 10.1016/j.acha.2020.03.002.
- [21] S. Goyal, M. Ahuja, and A. Kankanhalli, "Does the source of external knowledge matter? Examining the role of customer co-creation and partner sourcing in knowledge creation and innovation," *Inf. Manag.*, vol. 57, no. 6, p. 103325, 2020, doi: 10.1016/j.im.2020.103325.
- [22] T. Hendrickx, B. Cule, P. Meysman, S. Naulaerts, K. Laukens, and B. Goethals, "Mining association rules in graphs based on frequent cohesive itemsets," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 9078, no. 3, pp. 637–648, 2015, doi: 10.1007/978-3-319-18032-8_50.
- [23] K. E. Harvey, M. A. Suizzo, and K. M. Jackson, "Predicting the Grades of Low-Income-Ethnic-Minority Students from Teacher-Student Discrepancies in Reported Motivation," *J. Exp. Educ.*, vol. 84, no. 3, pp. 510–528, 2016, doi: 10.1080/00220973.2015.1054332.
- [24] R. Wehrens, "Multivariate Regression," *Bus. Media New York*, vol. 6, no. 1, pp. 149–185, 2013, doi: 10.1007/978-3-662-62027-4_8.
- [25] Š. Sonja, The Fourth European Conference on Information Literacy (ECIL): Abstracts The Fourth European Conference on Information Literacy (ECIL), no. October. 2016. [Online]. Available: http://repositorio.ipl.pt/bitstream/10400.21/7706/1/Information_literacy_in_Portuguese_university_context_a_necessary_intervention.pdf#page=76
- [26] M. Ranga and H. Etzkowitz, "Triple Helix Systems: An Analytical Framework for Innovation Policy and Practice in the Knowledge Society," *Ind. High. Educ.*, vol. 27, no. 4, pp. 237–262, 2013, doi: 10.5367/ihe.2013.0165.
- [27] U. Hasanah and D. A. Muatiara, "Perbandingan metode cosine similarity dan jaccard similarity untuk penilaian otomatis jawaban pendek," *Semin. Nas. Sist. Inf. dan Tek. Inform.*, no. SENSITIF 2019, pp. 1255–1263, 2019, [Online]. Available: <https://ejurnal.diponegara.ac.id/index.php/sensitif/article/view/511>
- [28] G. A. Carpenter and S. Grossberg, "A Massively Parallel Architecture for a Self-Organizing Neural Pattern Recognition Machine C . Stability-Plasticity Dilemma : Multiple Interacting Memory Systems The properties of plasticity and stability are intimately related . An adequate," *Pattern Recognit.*, vol. 115, pp. 54–115, 1987.
- [29] N. Trang, "Data mining for Education Sector, a proposed concept," *J. Appl. Data Sci.*, vol. 1, no. 1, pp. 11–19, 2020, doi: 10.47738/jads.v1i1.7.
- [30] R. H. Hama Aziz and N. Dimililer, "SentiXGboost: enhanced sentiment analysis in social media posts with ensemble XGBoost classifier," *J. Chinese Inst. Eng. Trans. Chinese Inst. Eng. A*, vol. 44, no. 6, pp. 562–572, 2021, doi: 10.1080/02533839.2021.1933598.
- [31] D. Sarpong, A. AbdRazak, E. Alexander, and D. Meissner, "Organizing practices of university, industry and government that facilitate (or impede) the transition to a hybrid triple helix model of innovation," *Technol. Forecast. Soc. Change*, vol. 123, pp. 142–152, 2017, doi: 10.1016/j.techfore.2015.11.032.
- [32] J. Jasperson, P. E. Carter, and R. W. Zmud, "A comprehensive conceptualization of post-adoptive behaviors associated with information technology enabled work systems," *MIS Q. Manag. Inf. Syst.*, vol. 29, no. 3, pp. 525–557, 2005, doi: 10.2307/25148694.