

---

# EduQate: Generating Adaptive Curricula through RMABs in Education Settings

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        There has been significant interest in the development of personalized and adaptive  
2        educational tools that cater to a student’s individual learning progress. A crucial  
3        aspect in developing such tools is in exploring how mastery can be achieved across  
4        a diverse yet related range of content in an efficient manner. While Reinforcement  
5        Learning and Multi-armed Bandits have shown promise in educational settings,  
6        existing works often assume the independence of learning content, neglecting  
7        the prevalent interdependencies between such content. In response, we introduce  
8        *Education Network Restless Multi-armed Bandits* (EdNetRMABs), utilizing a  
9        network to represent the relationships between interdependent arms. Subsequently,  
10       we propose *EduQate*, a method employing interdependency-aware Q-learning to  
11       make informed decisions on arm selection at each time step. We establish the  
12       optimality guarantee of EduQate and demonstrate its efficacy compared to baseline  
13       policies, using students modeled from both synthetic and real-world data.

## 14    1 Introduction

15    The COVID-19 pandemic has accelerated the adoption of educational technologies, especially  
16    on eLearning platforms. Despite abundant data and advancements in modeling student learning,  
17    effectively capturing the learning process with interdependent content remains a significant challenge  
18    [9]. The conventional rules-based approach to creating personalized learning curricula is impractical  
19    due to its labor-intensive nature and need for expert knowledge. Machine learning-based systems offer  
20    a scalable alternative, automatically generating personalized content to optimize learning [22, 24].

21    One possible approach to model the learning process is the Restless Multi-Armed Bandits (RMAB,  
22    [26]), where a teacher agent selects a subset of arms (concepts) to teach each round. However,  
23    RMAB’s assumption that arms are independent is unrealistic in educational settings. For example,  
24    solving a math question on the area of a triangle requires knowledge of algebra, arithmetic, and  
25    geometry. Practicing this question should enhance proficiency in all three areas. Models that ignore  
26    such interdependencies may inaccurately predict knowledge levels by assuming each exercise impacts  
27    only a single area.

28    In response to this challenge, we introduce an interdependency-aware RMAB model to the education  
29    setting. We posit that by acknowledging and modeling the learning dynamics of interdependent  
30    content, both teachers and algorithms can strategically leverage overlapping utility to foster mastery  
31    over a broader range of topics within a curriculum. We advocate for RMABs as a fitting model for  
32    this context, as the inherent dynamics of such a model align closely with the learning process.

33    In this study, our objective is to derive a teacher policy that effectively recommends educational  
34    content to students, accounting for interdependencies among the content to enhance overall utility  
35    (that characterizes understanding and retention of content). Our contributions are as follows:

- 36 1. We introduce Restless Multi-armed Bandits for Education (EdNetRMABs), enabling the  
37 modeling of learning processes with interdependent educational content.
- 38 2. We propose EduQate, a Whittle index-based heuristic algorithm that uses Q-learning to  
39 compute an inter-dependency-aware teacher policy. Unlike previous methods, EduQate does  
40 not require knowledge of the transition matrix to compute an optimal policy.
- 41 3. We provide a theoretical analysis of EduQate, demonstrating guarantees of optimality.
- 42 4. We present empirical results on simulated students and real-world datasets, showing the  
43 effectiveness of EduQate over other teacher policies.

## 44 **2 Related Work and Preliminaries**

### 45 **2.1 Restless Multi-Armed Bandits**

46 The selection of the right time and manner for limited interventions is a problem of great practical im-  
47 portance across various domains, including health intervention [17, 5], anti-poaching operations [20],  
48 education [13, 6, 2], etc. These problems share a common characteristic of having multiple arms  
49 in a Multi-armed Bandit (MAB) problem, representing entities such as patients, regions of a forest,  
50 or students’ mastery of concepts. These arms evolve in an uncertain manner, and interventions are  
51 required to guide them from "bad" states to "good" states. The inherent challenge lies in the limited  
52 number of interventions, dictated by the limited resources (e.g., public health workers, the number of  
53 student interactions). RMAB, a generalization of MAB, offers an ideal model for representing the  
54 aforementioned problems of interest. RMAB allows non-active bandits to also undergo the Markovian  
55 state transition, effectively capturing uncertainty in arm state transitions (reflecting uncertain state  
56 evolution), actions (representing interventions), and budget constraints (illustrating limited resources).

57 RMABs and the associated Markov Decision Processes (MDP) for each arm offer a valuable model for  
58 representing the learning process. Firstly, leveraging the MDPs associated with each arm provides the  
59 flexibility to adopt nuanced modeling of learning content, accommodating different learning curves  
60 for various content based on students’ strengths and weaknesses. Secondly, the transition probabilities  
61 serve as a useful mechanism to model forgetting (through state decay due to passivity or negligence)  
62 and learning (through state transitions to the positive state from repeated practice). Considering  
63 these aspects, RMABs prove to be a beneficial framework for personalizing and generating adaptive  
64 curricula across a diverse range of students.

65 In general, computing the optimal policy for a given set of restless arms in RMABs is recognized as a  
66 PSPACE-hard problem [18]. The Whittle index [26] provides an approach with a tractable solution  
67 that is provably optimal, especially when each arm is indexable. However, proving indexability can  
68 be challenging and often requires specification of the problem’s structure, such as the optimality of  
69 threshold policies [17, 16]. Moreover, much of the research on Whittle Index policies has focused  
70 on two-action settings or requires prior knowledge of the transition matrix of the RMABs. Meeting  
71 these conditions proves challenging in the educational context, where diverse students interact with  
72 educational systems, each possessing different prior knowledge and distinct learning curves for  
73 various topics.

74 WIQL [5], on the other hand, employs a Q-learning-based method to estimate the Whittle Index and  
75 has demonstrated provable optimality without requiring prior knowledge of the transition matrix. We  
76 utilize WIQL as a baseline method in our subsequent experiments.

77 In a recent investigation by [12], RMABs were explored within a network framework, requiring the  
78 agent to manage a budget while allocating a high-cost, high-benefit resource to one arm to “unlock”  
79 potential lower-cost, intermediate-benefit resources for the arm’s neighbors. The network effects  
80 emphasized in their work are triggered by an intentional, active action, enabling the agent to choose  
81 to propagate positive externalities to a selected arm’s neighbors within budget constraints. In contrast,  
82 our study delves into scenarios where network effects are indirect results of an active action, and the  
83 agent lacks direct control over such effects. Thus, the challenge lies in accurately modeling these  
84 network effects and leveraging them when beneficial.

## 85 2.2 Reinforcement Learning in Education

86 In the realm of education, numerous researchers have explored optimizing the sequencing of in-  
87 structional activities and content, assuming that optimal sequencing can significantly impact student  
88 learning. RL is a natural approach for making sequential decisions under uncertainty [1]. While RL  
89 has seen success in various educational applications, effectively sequencing interdependent content in  
90 a personalized and adaptive manner has yielded mixed or insignificant results compared to baseline  
91 teacher policies [11, 21, 8]. In general, these RL works focus on data-driven methods using student  
92 activity logs to estimate students' knowledge states and progress, assuming that the interdependencies  
93 between learning content are encapsulated in students' learning histories [9, 3, 19]. In contrast, our  
94 work focuses on modelling these interdependencies directly.

95 Of particular relevance are factored MDPs applied to skill acquisition introduced by [11]. While fac-  
96 tored MDPs account for interdependencies amongst skills, decentralized policy learning is infeasible  
97 as policies must consider the joint state space. Our work leverages the advantage of decentralized  
98 policy learning provided by RMABs and introduces a novel decentralized learning approach that  
99 exploits interdependencies between arms.

100 Complementary to RL methods in education is the utilization of knowledge graphs to uncover  
101 relationships between learning content [9]. Existing research primarily focuses on establishing these  
102 relationships through data-driven methods (e.g. [7, 23]) often leveraging student-activity logs. In this  
103 work, we complement such research by presenting an approach where bandit methods can effectively  
104 operate with knowledge graphs derived by such methods.

## 105 3 Model

106 In this section, we introduce the Restless Multi-Armed Bandits for Education (EdNetRMABs). It  
107 is important to note that while we specifically apply EdNetRMABs to the education setting, the  
108 framework can be seamlessly translated to other scenarios where modeling the effects of active  
109 actions within a network is critical. For ease of access, a table of notations is provided in Table 2.

110 In education, a teacher recommends learning content, or items, to maximize student education, often  
111 with content from online platforms. Items are grouped by topics, such as "Geometry," where exposure  
112 to one piece of content can enhance knowledge across others in the same group. This cumulative  
113 learning effect which we refer to as "network effects", implies that exposure to an item is likely  
114 to positively impact the student's success on items within the same group. A successful teacher  
115 accurately estimates a student's knowledge state over repeated interactions, leveraging these network  
116 effects to promote both breadth and depth of understanding through recommendations.

### 117 3.1 EdNetRMABs

118 The RMAB model tasks an agent with selecting  $k$  arms from  $N$  arms, constrained by a limit on the  
119 number of arms that can be pulled at each time step. The objective is to find a policy that maximizes  
120 the total expected discounted reward, assuming that the state of each arm evolves independently  
121 according to an underlying MDP.

122 The EdNetRMABs model extends RMABs by allowing for active actions to propagate to other arms  
123 dependent on the current arm when it is being pulled, thus relaxing the assumption of independent  
124 arms. This is operationalized by organising the arms in a network, and pulling of an arm results in  
125 changes for its neighbors, or members in the same group.

126 When applied to education setting, the EdNetRMABs is formalized as follows:

127 **Arms** Each arm, denoted as  $i \in 1, \dots, N$ , signifies an item. In the context of this networked  
128 environment, each arm belongs to a group  $\phi \in \{1, \dots, L\}$  representing the overarching topic that  
129 encompasses related items. It's important to note that arm membership is not mutually exclusive,  
130 allowing arms to be part of multiple groups. This flexibility enables a more nuanced modeling of  
131 interdependencies among educational content. For instance, a question involving the calculation of  
132 the area of a triangle may span both arithmetic and geometry groups.

133 **State space** In this framework, each arm possesses a binary latent state, denoted as  $s_i \in \{0, 1\}$ ,  
 134 where “0” represents an “unlearned” state, and “1” indicates a “learned” state. Considering all arms  
 135 collectively, these states serve as a representation of the student’s overall knowledge state. In the  
 136 current work, it is assumed that the states of all arms are fully observable, providing a comprehensive  
 137 model of the student’s understanding of the various educational concepts.

138 **Action space** To capture the network effects associated with arm pulls, we depart from the conven-  
 139 tional RMAB framework with a binary action space  $A = \{0, 1\}$  by introducing a pseudo-action. In  
 140 this modified setup, the action space is extended to  $A = \{0, 1, 2\}$ , where actions 0 and 2 represent  
 141 “no-pull” and “pull”, as commonly used in bandit literature. Notably, in EdNetRMABs, a third action  
 142 1 is introduced to simulate the network effects resulting from pulling another arm within the same  
 143 group. It is important to clarify that agents do not directly engage with action 1 but we employ it  
 144 solely for modeling network effects, hence the term “pseudo-action”.

145 **Transition function** For a given arm  $i$ , let  $P_{s,s'}^{a,i}$  represent the probability of the arm transitioning  
 146 from state  $s$  to  $s'$  under action  $a$ . It’s noteworthy that, in typical real-world educational settings, the  
 147 actual transition functions governing the states of the arms are often unknown and, even for the same  
 148 concept, may vary among students due to differences in prior knowledge. To address this challenge,  
 149 we adopt model-free approaches in this study, devising methods to compute the teacher policy without  
 150 relying on explicit knowledge of these transition functions. In the following experiments, we maintain  
 151 the assumption of non-zero transition probabilities, and enforce constraints that are aligned with the  
 152 current domain [17]: (i) The arms are more likely to stay in the positive state than change to the  
 153 negative state:  $P_{0,1}^0 < P_{1,1}^0$ ,  $P_{0,1}^1 < P_{1,1}^1$  and  $P_{0,1}^2 < P_{1,1}^2$ ; (ii) The arm tends to improve the latent  
 154 state if more efforts is spent on that arm, i.e., it is active or semi-active:  $P_{0,1}^0 < P_{0,1}^1 < P_{0,1}^2$  and  
 155  $P_{1,1}^0 < P_{1,1}^1 < P_{1,1}^2$ .

156 With the formalization of the EdNetRMABs model provided, we now apply it to an educational  
 157 context. In this scenario, the agent assumes the role of a teacher and takes actions during each time  
 158 step  $t \in \{1, \dots, T\}$ . Specifically, at each time step, the teacher recommends an item for the student to  
 159 study. We represent the vector of actions taken by the teacher at time step  $t$  as  $\mathbf{a}^t \in \{0, 1, 2\}^N$ . Here,  
 160 arm  $i$  is considered to be active at time  $t$  if  $a^t(i) = 2$  and passive when  $a^t(i) = 0$ . When arm  $i$  is  
 161 pulled, the set of arms that share the same group membership as arm  $i$ , denoted as  $\phi_i^-$  under goes  
 162 the pseudo-action, represented as  $a^t(j) = 1$  for all  $j \in \phi_i^-$ . In our framework, the teacher agent  
 163 acts on exactly one arm per time step to simulate the real-world constraint that the teacher can only  
 164 recommend one concept to students ( $\sum_i I_{a^t(i)=2} = 1, \forall t$ ). Subsequent to taking action, the teacher  
 165 receives  $\mathbf{s}^t \in \{0, 1\}^N$ , a vector reflecting the state of all arms, and reward  $r_t = \sum_{i=1}^N s^t(i)$ . The  
 166 vector  $\mathbf{s}^t$  represents the overall knowledge state of the student. The teacher agent’s goal, therefore, is  
 167 to maximize the long term rewards, either discounted or averaged.

## 168 4 EduQate

169 Q-learning [25] is a popular reinforcement learning method that enables an agent to learn optimal  
 170 actions in an environment by iteratively updating its estimate of state-action value,  $Q(s, a)$ , based on  
 171 the rewards it receives. At each time step  $t$ , the agent takes an action  $a$  using its current estimate of  $Q$   
 172 values and current state  $s$ , thus received a reward of  $r(s)$  and new state  $s'$ . We provide an abridged  
 173 introduction to Q-learning in the Appendix F.

174 Expanding upon Q-learning, we introduce *EduQate*, a tailored Q-learning approach designed for  
 175 learning Whittle-index policies in EdNetRMABs. In the interaction with the environment, the agent  
 176 chooses a single item, represented by arm  $i$ , to recommend to the student. In this context, the agent  
 177 possesses knowledge of the group membership  $\phi_i$  of the selected arm and observes the rewards  
 178 generated by activating arm  $i$  and semi-activating arms in  $\phi_i^-$ . *EduQate* utilizes this interaction to  
 179 learn the Q-values for all arms and actions.

180 To adapt Q-learning to EdNetRMABs, we propose leveraging the learned Q-values to select the arm  
 181 with the highest estimate of the Whittle index, defined as:

---

**Algorithm 1** Q-Learning for EdNetRMABs (EduQate)

---

**Input:** Number of arms  $N$   
Initialize  $Q_i(s, a) \leftarrow 0$  and  $\lambda_i(s) \leftarrow 0$  for each state  $s \in S$  and each action  $a \in \{0, 1, 2\}$ , for each arm  $i \in 1, \dots, N$ .  
Initialize replay buffer  $D$  with capacity  $C$ .  
**for**  $t$  in  $1, \dots, T$  **do**  
     $\epsilon \leftarrow \frac{N}{N+t}$   
    With probability  $\epsilon$ , select one arm uniformly at random. Otherwise, select arm with highest Whittle Index,  $i = \arg \max_i \lambda_i$ .  
    **for** arm  $n$  in  $1, \dots, N$  **do**  
        **if**  $n \neq i$  **then**  
            Set arm  $n$  to passive,  $a_n^t = 0$   
        **else**  
            Set arm  $n$  to active,  $a_n^t = 2$   
            **for**  $j \in \phi_i^-$  **do**  
                Set arms in same group as  $i$  to semi-active,  $a_j^t = 1$   
            **end for**  
        **end if**  
    **end for**  
    Execute actions  $\mathbf{a}^t$  and observe reward  $r^t$  and next state  $s^{t+1}$  for all arms  
    Store experience  $(s^t, \mathbf{a}^t, \mathbf{r}^t, s^{t+1})$  in replay buffer  $D$ .  
    Sample minibatch  $B$  of Experience from replay buffer  $D$ .  
    **for** Experience in minibatch  $B$  **do**  
        Update  $Q_n(s, a)$  using Q-learning update in Equation 11.  
        Compute  $\lambda_n$  using Equation 1  
    **end for**  
**end for**

---

$$\lambda_i = Q(s_i, a_i = 2) - Q(s_i, a_i = 0) + \sum_{j \in \phi_i^-} (Q(s_j, a_j = 1) - Q(s_j, a_j = 0)) \quad (1)$$

182 Here,  $\lambda_i$  is the Whittle Index estimate for arm  $i$ . In essence, the Whittle Index of arm  $i$  is computed as  
183 the linear combination of the value associated with taking action on arm  $i$  over passivity and the value  
184 of associated with semi-actively engaging with members from same group, compared to passivity.

185 To improve the convergence of Q-learning, we incorporate Experience Replay [15]. This involves  
186 saving the teacher algorithm’s previous experiences in a replay buffer and drawing mini-batches  
187 of samples from this buffer during updates to enhance convergence. In Section 4.1, we prove that  
188 EduQate will converge to the optimal policy. However, in practice, we may not have enough episodes  
189 to fully train EduQate. Therefore, we propose Experience Replay to mitigate the cold-start problem  
190 common in RL applications, a common problem where initial student interactions with sub-optimal  
191 teachers can lead to poor learning experiences [3].

192 The pseudo-code is provided in Algorithm 1. Similar to WIQL [5], we employ a  $\epsilon$ -decay policy that  
193 facilitates exploration and learning in the early steps, and proceeds to exploit the learned Q-values in  
194 later stages.

#### 195 4.1 Analysis of EduQate

196 In this section, we analyze EduQate closely, and show that EduQate does not alter the optimality  
197 guarantees of Q-learning under the constraint that  $k = 1$  (Theorem 1). Our method relies on the  
198 assumption that teachers are limited to assign 1 item to the student at each time step. Theorem 2  
199 analyzes EduQate under the conditions that  $k > 1$ . Since our setting involves the semi-active actions,  
200 we should compute Equation 1. To reiterate,  $\phi_i$  here refers to the group that arm  $i$  belongs to, and  
201  $\phi_i^-$  is the same group but does not include arm  $i$ . If arm  $i$  is selected, then all the remaining arms in  
202 group  $\phi_i^-$  should be semi-active.

203 **Theorem 1** *Choosing the top arm with the largest  $\lambda$  value in Equation 1 is equivalent to maximizing*  
 204 *the cumulative long-term reward.*

205 *Proof.* According to the approach, we select the arm according to the  $\lambda$  value. Assume arm  $i$  has  
 206 the highest  $\lambda$  value, then for any arm  $j$  where  $j \neq i$ , we have

$$\lambda_i \geq \lambda_j \quad (2)$$

207 According to the definition of  $\lambda$  in Equation 1, we move the negative part to the other side, and the  
 208 left side becomes:

$$Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + Q(s_j, a_j = 0) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 0))$$

209 and the right side is similar. There are three cases:

210 • arm  $i$  and arm  $j$  are not connected, and group  $\phi_i$  and  $\phi_j$  has no overlap, i.e.,  $\phi_i \cap \phi_j = \emptyset$ . We add  
 211  $\sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0)$  on both sides. This denotes the addition of  $Q(s_z, a_z = 0)$  for all arm  $z$   
 212 that are not included in the set of  $\phi_i$  or  $\phi_j$ . We have the left side:

$$\begin{aligned} & Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + Q(s_j, a_j = 0) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 0)) + \sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) \\ &= Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + \sum_{j \notin \phi_i} (Q(s_j, a_j = 0)) \\ &= Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_i) \end{aligned} \quad (3)$$

213 Similarly, we do the same for the right side and thus, the equation 2 becomes

$$Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_i) \geq Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_j)$$

214 • arm  $i$  and arm  $j$  are not connected, but group  $\phi_i$  and  $\phi_j$  has overlap, i.e.,  $\phi_i \cap \phi_j \neq \emptyset$ . In this case,  
 215 we add  $\sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) - \sum_{z \in \phi_i \cap \phi_j} Q(s_z, a_z = 0)$  on both sides.

216 • arm  $i$  and arm  $j$  are connected, and group  $\phi_i$  and  $\phi_j$  has overlap, i.e.,  $\phi_i \cap \phi_j \neq \emptyset$ , and  $\{i, j\} \subset \phi_i \cap$   
 217  $\phi_j$ . This case is similar to the previous one, we add  $\sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) - \sum_{z \in \phi_i \cap \phi_j} Q(s_z, a_z =$   
 218  $0)$  on both sides.

219 The detailed proof is provided in Appendix B. □

220 Thus when  $k = 1$ , selecting the top arm according to the  $\lambda$  value is equivalent to maximizing the  
 221 cumulative long-term reward, and is guaranteed to be optimal.

222 **Theorem 2** *When  $k > 1$ , selecting the  $k$  arms is a NP-hard problem. The non-asymptotic tight*  
 223 *upper bound and non-asymptotic tight lower bound for getting the optimal solution are  $o(C(n, k))$*   
 224 *and  $\omega(N)$ , respectively.*

225 *Proof Sketch.* This problem can be considered as a variant of the knapsack problem. If we disregard  
 226 the influence of the shared neighbor nodes for two selected arms, then selecting arm  $i$  will not  
 227 influence the future selection of arm  $j$ . In such instances, the problem of selecting the  $k$  arms is  
 228 simplified to the traditional 0/1 knapsack problem, a classic NP-hard problem. Therefore, when  
 229 considering the effect of shared neighbor nodes for two selected arms, this problem is at least as  
 230 challenging as the 0/1 knapsack problem. □

231 When  $k > 1$ , it is difficult to compute the optimal solution, we provide a heuristic greedy algorithm  
 232 with the complexity of  $O(\frac{(2N-k)*k}{2})$  in Section C in the appendix.

## 233 5 Experiment

234 In this section, we demonstrate the effectiveness of EduQate against benchmark algorithms on  
 235 synthetic students and students derived from a real-world dataset, the Junyi Dataset and the OLI  
 236 Statics dataset. All experiments are run on CPU only. In our experiments, we compare EduQate with  
 237 the following policies:

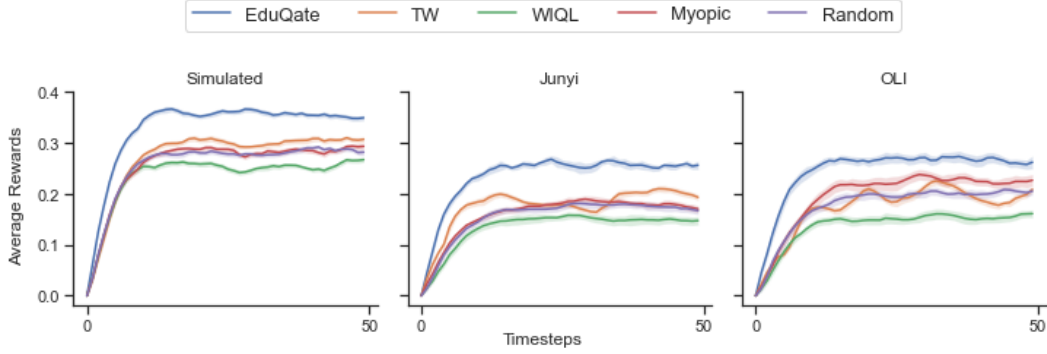


Figure 1: Average rewards for the respective algorithms on 3 datasets, averaged across 30 runs. Shaded regions represent standard error.

- 238 • **Threshold Whittle (TW)**: This algorithm, proposed by [17], utilizes an efficient closed-form  
 239 approach to compute the Whittle index, considering only the pull action as active. It operates under  
 240 the assumption that transition probabilities are known and stands as the state-of-the-art in RMABs.
- 241 • **WIQL**: This algorithm employs a Q-learning-based Whittle Index approach [5]. It learns Q-values  
 242 using the pull action as the only active strategy and calculates the Whittle Index based on the  
 243 acquired Q-values.
- 244 • **Myopic**: This strategy disregards the impact of the current action on future rewards, concentrating  
 245 solely on predicted immediate rewards. It selects the arm that maximizes the expected reward at  
 246 the immediate time step.
- 247 • **Random**: This strategy randomly selects arms with uniform probability, irrespective of the under-  
 248 lying state.

249 Inspired by work in healthcare settings [12, 14], we compare the policies by the *Intervention Benefit*  
 250 (*IB*), as shown in the following equation:

$$IB_{Random, EQ}(\pi) = \frac{\mathbb{E}_{\pi}(R(\cdot)) - \mathbb{E}_{Random}(R(\cdot))}{\mathbb{E}_{EQ}(R(\cdot)) - \mathbb{E}_{Random}(R(\cdot))} \quad (4)$$

251 where *EQ* represents EduQate, and *Random* represents a policy where the arms are selected at random.  
 252 Prior work in educational settings has demonstrated that random policies can yield robust learning  
 253 outcomes through spaced repetition [9, 10]. Therefore, to establish efficacy, successful algorithms  
 254 must demonstrate superiority over random policies. Our chosen metric, *IB*, effectively compares  
 255 the extent to which a challenger algorithm  $\pi$  outperforms a random policy in comparison to our  
 256 algorithm.

## 257 5.1 Experiment setup

258 In all experiments, we commence by initializing all arms in state 0 and permit the teacher algorithms  
 259 to engage with the student for a total of 50 actions, pulling exactly 1 arm (i.e.  $k = 1$ ) at each time step.  
 260 Following the completion of these actions, the episode concludes, and the student state is reset. This  
 261 process is iterated across 800 episodes, for a total of 30 seeds. The datasets used in our experiment  
 262 are described below:

263 **Synthetic dataset.** Given the domain-motivated constraints on the transition functions highlighted  
 264 in Section 3.1, we create a simulator based on  $N = 50$ ,  $S \in \{0, 1\}$ ,  $N_{\text{topics}} = 20$ . We randomly  
 265 assign arms to topic groups, and allow arms to be assigned to be more than one topic. Under this  
 266 method, number of arms under each group may not be equal. For each trial, a new transition matrix  
 267 is generated to simulate distinct student scenarios.

268 **Junyi dataset.** The Junyi dataset [7] is an extensive dataset collected from the Junyi Academy <sup>1</sup>,  
 269 an eLearning platform established in 2012 on the basis of the open-source code released by Khan

<sup>1</sup><http://www.junyiacademy.org/>

Table 1: Comparison of policies on synthetic, Junyi, and OLI datasets.  $\mathbb{E}[R]$  represents the average reward obtained in the final episode of training. Statistic after  $\pm$  represents standard error across 30 trials.

Policy	Synthetic		Junyi		OLI	
	$\mathbb{E}[IB](\%) \pm$	$\mathbb{E}[R] \pm$	$\mathbb{E}[IB](\%) \pm$	$\mathbb{E}[R] \pm$	$\mathbb{E}[IB](\%) \pm$	$\mathbb{E}[R] \pm$
Random	-	$26.84 \pm 0.46$	-	$15.82 \pm 0.34$	-	$18.46 \pm 0.35$
WIQL	$-49.03 \pm 15.07$	$24.60 \pm 0.43$	$-26.77 \pm 7.39$	$14.01 \pm 0.97$	$-60.20 \pm 19.38$	$14.33 \pm 0.42$
Myopic	$-3.44 \pm 5.81$	$27.07 \pm 0.52$	$10.74 \pm 3.13$	$16.86 \pm 0.356$	$39.92 \pm 12.00$	$20.51 \pm 0.48$
TW	$37.21 \pm 17.02$	$28.50 \pm 0.47$	$31.284 \pm 2.65$	$15.819 \pm 0.34$	$0.20 \pm 9.27$	$18.07 \pm 0.21$
<b>EduQate</b>	<b>100.0</b>	<b><math>34.33 \pm 0.49</math></b>	<b>100.0</b>	<b><math>24.53 \pm 0.31</math></b>	<b>100.0</b>	<b><math>25.47 \pm 0.47</math></b>

270 Academy. In this dataset, there are nearly 26 million student-exercise interactions across 250 000  
 271 students in its mathematics curriculum. For this experiment, we selected the top 100 exercises with  
 272 the most student interactions to create our student models. Using our method to generate groups, the  
 273 resultant EdNetRMAB has  $N = 100$  and  $N_{topics} = 21$ .

274 **OLI Statics dataset.** The OLI Statics dataset [4] comprises student interactions with an online  
 275 Engineering Statics course<sup>2</sup>. In this dataset, each item is assigned one or more Knowledge Compo-  
 276 nents (KCs) based on the related topics. After filtering for the top 100 items with the most student  
 277 interactions, the resultant EdNetRMAB includes  $N = 100$  items and  $N_{topics} = 76$  distinct topics.

## 278 5.2 Creating student models

279 In this section, we outline the procedure for generating student models aimed at simulating the  
 280 learning process. To clarify, a student model in this context is defined as a set of transition matrices  
 281 for all items. These matrices are employed with EdNetRMABs to simulate the learning dynamics.

282 We employ various strategies to model transitions within the RMAB framework. Active transitions  
 283 are determined by assessing the average success rate on a question before and after a learning  
 284 intervention. Passive transitions are influenced by difficulty ratings, with more challenging questions  
 285 more prone to rapid forgetting. Semi-active transitions, on the other hand, are computed as proportion  
 286 of active transition, guided by similarity scores. Here, we provide an outline and the full details can  
 287 be found in Appendix D.

288 **Active Transitions.** We use data on students’ correct response rate after interacting with an item to  
 289 create the transition matrix for action 2, based on the change in correctness rates before and after a  
 290 learning intervention.

291 **Passive Transitions.** To construct passive transitions for items, we use relative difficulty scores to  
 292 determine transitions based on difficulty levels. We assume that higher difficulty correlates with a  
 293 greater likelihood of forgetting, resulting in higher failure rates. Specifically, higher difficulty values  
 294 correspond to higher  $P_{1,0}^0$  values, indicating a greater likelihood of forgetting. The transition matrix  
 295 for the passive action  $a = 0$  is then randomly generated, with values influenced by difficulty levels.

296 **Semi-active Transitions.** To derive semi-active transitions, we use similarity scores between exercises  
 297 from the Junyi dataset. We first normalize these scores to the range  $[0, 1]$ . Then, for any chosen arm,  
 298 we compute its transition matrix under the semi-active action  $a = 1$  as a proportion of its active  
 299 action transitions,  $P_{0,1}^1 = \sigma(P_{0,1}^2)$ , where  $\sigma$  signifies the similarity proportion.

300 The arm’s transition matrix for the semi-active action varies due to different similarity scores between  
 301 pairs in the same group. To address this, we use the average similarity score to determine the  
 302 proportion. Since the OLI dataset does not contain similarity ratings, we assume a constant similarity  
 303 rating of  $\sigma = 0.8$  for all pairs.

## 304 6 Results

305 The experimental results for the synthetic, Junyi, and OLI datasets are shown in Table 1. We report  
 306 the average intervention benefit  $IB$  and final episode rewards from thirty independent runs for five  
 307 algorithms: EduQate, TW, WIQL, Myopic, and Random. EduQate consistently outperforms the other  
 308 policies across all datasets, demonstrating higher intervention benefits and average rewards.

<sup>2</sup><https://oli.cmu.edu/courses/engineering-statics-open-free/>



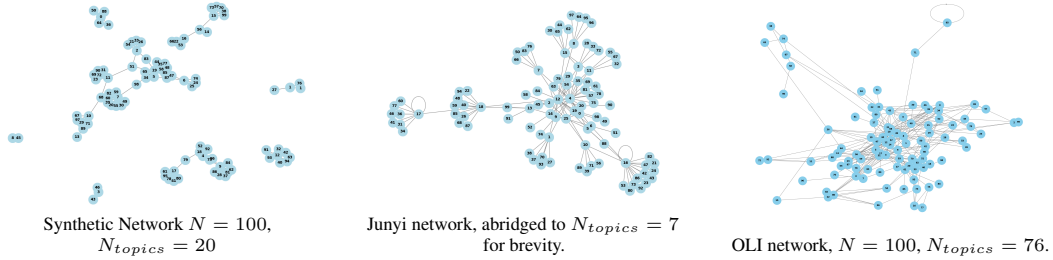


Figure 2: This visualization compares network complexities from our experiments. The synthetic dataset (left) shows simpler, isolated groups, while the real-world datasets (Junyi, middle; OLI, right) displays more intricate and interconnected relationships amongst items.

309 In terms of  $IB$ , we note that all challenger policies do not exceed 50%, indicating two key points.  
 310 First, as noted in prior works [9], our results confirm that random policies in educational settings are  
 311 robust and difficult to surpass, even when algorithms are equipped with knowledge of the learning  
 312 dynamics. Second, our interdependency-aware EduQate performs well over random policies and  
 313 other algorithms, highlighting the importance of considering network effects and interdependencies  
 314 in EdNetRMABs.

315 Notably, WIQL, which relies solely on Q-learning for active and passive actions, performs worse  
 316 than a random policy, likely due to misattributing positive network effects to passive actions. Despite  
 317 having access to the transition matrix, TW does not perform as well as the interdependency-aware  
 318 EduQate. While it has demonstrated effectiveness in traditional RMABs, TW weaknesses become  
 319 evident in the current setting, where pulling an arm has wider implications to other arms. Overall,  
 320 EduQate has demonstrated robust and effective performance in maximizing rewards across different  
 321 datasets. Figure 1 shows the average rewards obtained in the final episode for each algorithm.

322 Figure 2 provides visualizations of the networks generated from synthetic students and mined from  
 323 real-world datasets. The synthetic dataset produces networks with distinct isolated groups, contrasting  
 324 with the more intricate and interconnected networks from the Junyi and OLI datasets, reflecting  
 325 real-world complexities. Despite these differing topologies and levels of interdependency, EduQate  
 326 performs well under all network setups. In Appendix E.1, we explore the effects of different network  
 327 topologies by varying the number of topics while limiting the membership of each item. We find that  
 328 as network interdependencies are reduced, the network effects diminish, and such EdNetRMABs  
 329 can be approximated to traditional RMABs with independent arms. Under these conditions, our  
 330 algorithm does not perform as well as other baseline policies.

331 Finally, an ablation study detailed in Appendix E.2 examines the effectiveness of the replay buffer in  
 332 EduQate. The study shows that the replay buffer helps overcome the cold-start problem, where initial  
 333 learning episodes provide sub-optimal experiences for students [3].

## 334 7 Conclusion and Limitations

335 In this paper, we introduced EdNetRMABs to the education setting, a variant of MAB designed to  
 336 model interdependencies in educational content. We also proposed EduQate, a novel Whittle-based  
 337 learning algorithm tailored for EdNetRMABs. Unlike other Whittle-based algorithms, EduQate com-  
 338 puts an optimal policy without requiring knowledge of the transition matrix, while still accounting  
 339 for the network effects of pulling an arm. We demonstrated the guaranteed optimality of a policy  
 340 trained under EduQate and showcased its effectiveness on synthetic and real-world datasets, each  
 341 with its own characteristic.

342 Our work assumes that student knowledge states are fully observable and available at all times, which  
 343 is a limitation. Despite this, we believe our work is significant and can inspire further research to  
 344 improve efficiencies in education. For future work, we aim to extend EduQate to handle partially  
 345 observable states and address the cold-start problem in education systems by minimizing the initial  
 346 exploratory phase.

347 **References**

- 348 [1] Richard C Atkinson. Ingredients for a theory of instruction. *American Psychologist*, 27(10):  
349 921, 1972.
- 350 [2] Aqil Zainal Azhar, Avi Segal, and Kobi Gal. Optimizing representations and policies for  
351 question sequencing using reinforcement learning. *International Educational Data Mining*  
352 *Society*, 2022.
- 353 [3] Jonathan Bassen, Bharathan Balaji, Michael Schaarschmidt, Candace Thille, Jay Painter, Dawn  
354 Zimmaro, Alex Games, Ethan Fast, and John C Mitchell. Reinforcement learning for the  
355 adaptive scheduling of educational activities. In *Proceedings of the 2020 CHI conference on*  
356 *human factors in computing systems*, pages 1–12, 2020.
- 357 [4] Norman Bier. Oli engineering statics - fall 2011 (114 students), 2011. URL [https://](https://pslcdatashop.web.cmu.edu/DatasetInfo?datasetId=590)  
358 [pslcdatashop.web.cmu.edu/DatasetInfo?datasetId=590](https://pslcdatashop.web.cmu.edu/DatasetInfo?datasetId=590).
- 359 [5] Arpita Biswas, Gaurav Aggarwal, Pradeep Varakantham, and Milind Tambe. Learn to intervene:  
360 An adaptive learning policy for restless bandits in application to preventive healthcare. *arXiv*  
361 *preprint arXiv:2105.07965*, 2021.
- 362 [6] Colton Botta, Avi Segal, and Kobi Gal. Sequencing educational content using diversity aware  
363 bandits. 2023.
- 364 [7] Haw-Shiuan Chang, Hwai-Jung Hsu, and Kuan-Ta Chen. Modeling exercise relationships in  
365 e-learning: A unified approach. In *EDM*, pages 532–535, 2015.
- 366 [8] Shayan Doroudi, Vincent Aleven, and Emma Brunskill. Robust evaluation matrix: Towards a  
367 more principled offline exploration of instructional policies. In *Proceedings of the fourth (2017)*  
368 *ACM conference on learning@ scale*, pages 3–12, 2017.
- 369 [9] Shayan Doroudi, Vincent Aleven, and Emma Brunskill. Where’s the reward? a review of rein-  
370 forcement learning for instructional sequencing. *International Journal of Artificial Intelligence*  
371 *in Education*, 29:568–620, 2019.
- 372 [10] Hermann Ebbinghaus. *Über das gedächtnis: untersuchungen zur experimentellen psychologie*.  
373 Duncker & Humblot, 1885.
- 374 [11] Derek Green, Thomas Walsh, Paul Cohen, and Yu-Han Chang. Learning a skill-teaching  
375 curriculum with dynamic bayes nets. In *Proceedings of the AAAI Conference on Artificial*  
376 *Intelligence*, volume 25, pages 1648–1654, 2011.
- 377 [12] Christine Herlihy and John P. Dickerson. Networked restless bandits with positive externalities,  
378 2022.
- 379 [13] Andrew S Lan and Richard G Baraniuk. A contextual bandits framework for personalized  
380 learning action selection. In *EDM*, pages 424–429, 2016.
- 381 [14] Dexun Li and Pradeep Varakantham. Avoiding starvation of arms in restless multi-armed bandits.  
382 In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent*  
383 *Systems*, pages 1303–1311, 2023.
- 384 [15] Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and  
385 teaching. *Machine learning*, 8:293–321, 1992.
- 386 [16] Keqin Liu and Qing Zhao. Indexability of restless bandit problems and optimality of whittle  
387 index for dynamic multichannel access. *IEEE Transactions on Information Theory*, 56(11):  
388 5547–5567, 2010.
- 389 [17] Aditya Mate, Jackson A Killian, Haifeng Xu, Andrew Perrault, and Milind Tambe. Collapsing  
390 bandits and their application to public health interventions. *arXiv preprint arXiv:2007.04432*,  
391 2020.

- 392 [18] Christos H Papadimitriou and John N Tsitsiklis. The complexity of optimal queueing network  
393 control. In *Proceedings of IEEE 9th Annual Conference on Structure in Complexity Theory*,  
394 pages 318–322. IEEE, 1994.
- 395 [19] Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J  
396 Guibas, and Jascha Sohl-Dickstein. Deep knowledge tracing. *Advances in neural information*  
397 *processing systems*, 28, 2015.
- 398 [20] Yundi Qian, Chao Zhang, Bhaskar Krishnamachari, and Milind Tambe. Restless poachers:  
399 Handling exploration-exploitation tradeoffs in security domains. In *Proceedings of the 2016*  
400 *International Conference on Autonomous Agents & Multiagent Systems*, pages 123–131, 2016.
- 401 [21] Avi Segal, Yossi Ben David, Joseph Jay Williams, Kobi Gal, and Yaar Shalom. Combining  
402 difficulty ranking with multi-armed bandits to sequence educational content. In *Artificial*  
403 *Intelligence in Education: 19th International Conference, AIED 2018, London, UK, June 27–30,*  
404 *2018, Proceedings, Part II 19*, pages 317–321. Springer, 2018.
- 405 [22] Shitian Shen, Markel Sanz Ausin, Behrooz Mostafavi, and Min Chi. Improving learning &  
406 reducing time: A constrained action-based reinforcement learning approach. In *Proceedings of*  
407 *the 26th conference on user modeling, adaptation and personalization*, pages 43–51, 2018.
- 408 [23] Anni Siren and Vassilios Tzerpos. Automatic learning path creation using oer: a systematic  
409 literature mapping. *IEEE Transactions on Learning Technologies*, 2022.
- 410 [24] Utkarsh Upadhyay, Abir De, and Manuel Gomez Rodriguez. Deep reinforcement learning of  
411 marked temporal point processes. *Advances in Neural Information Processing Systems*, 31,  
412 2018.
- 413 [25] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3):279–292, 1992.
- 414 [26] Peter Whittle. Restless bandits: Activity allocation in a changing world. *Journal of applied*  
415 *probability*, pages 287–298, 1988.

417 **A Table of Notations**

Table 2: Notations

Notation	Description
$N, N_{topics}$	$N$ : number of arms in EdNetRMABs; $N_{topics}$ : number of topic groups
$s_i^t$	$s_i^t$ : state of arm $i$ at time step $t$ . 1: learned, 0: unlearned.
$a_i^t$	$a_i^t$ : action of arm $i$ at time step $t$ . 0: passive action, 1: semi-active action, 2: active action.
$\mathbf{s}, \mathbf{a}$	$\mathbf{s}, \mathbf{a}$ : joint state vector and joint action vector of EdNetRMABs.
$\phi_i, \phi_i^-$	$\phi_i$ : the set of arms that includes the arm $i$ and its connected neighbors, $\phi_i^-$ : $\phi_i$ that exclude arm $i$ .
$P_{s,s'}^{i,a}$	$P_{s,s'}^{i,a}$ is the probability of transition from state $s$ to $s'$ when arm $i$ is taking action $a$ .
$Q_i(s_i, a_i)$	$Q_i(s_i, a_i)$ is the state-action value function for the arm $i$ when taking action $a_i$ with state $s_i$ .
$V_i(s_i)$	The value function for arm $i$ at the state $s_i$ .

418 **B Proof for the theorem**

419 We rewrite the theorem here for ease of explanation.

420 **Theorem 3** Choose top arms according to the  $\lambda$  value in Equation 1 is equivalent to maximize the  
421 cumulative long-term reward.422 *Proof.* According to the approach, we select the arm according to the  $\lambda$  value. Assume arm  $i$  has  
423 the highest  $\lambda$  value, then for any arm  $j$ , where  $i \neq j$ , we have

$$\begin{aligned}
& \lambda_i \geq \lambda_j \\
& Q(s_i, a_i = 1) - Q(s_i, a_i = 0) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1) - Q(s_i, a_i = 0)) \geq Q(s_j, a_j = 1) - Q(s_j, a_j = 0) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 1) - Q(s_j, a_j = 0)) \\
& Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + Q(s_j, a_j = 0) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 0)) \geq Q(s_j, a_j = 1) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 1)) + Q(s_i, a_i = 0) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 0))
\end{aligned} \tag{5}$$

424 There are two cases:

425 • **arm  $i$  and arm  $j$  are not connected, and group  $\phi_i$  and  $\phi_j$  has no overlap, i.e.,  $\phi_i \cap \phi_j = \emptyset$ .** We  
426 add  $\sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0)$  on both sides, we can have the left side:

$$\begin{aligned}
& Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + Q(s_j, a_j = 0) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 0)) + \sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) \\
& = Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + \sum_{j \notin \phi_i^-} (Q(s_j, a_j = 0)) \\
& = Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_i)
\end{aligned} \tag{6}$$

427 Similarly, the right side becomes

$$Q(s_j, a_j = 1) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 1)) + \sum_{i \notin \phi_j} (Q(s_i, a_i = 0)) = Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_j) \tag{7}$$

428 Thus, the equation 2 becomes

$$Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_i) \geq Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_j) \tag{8}$$

429 • **arm  $i$  and arm  $j$  are not connected, but group  $\phi_i$  and  $\phi_j$  has overlap, i.e.,  $\phi_i \cap \phi_j \neq \emptyset$ .** In this  
430 case, we add  $\sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) - \sum_{z \in \phi_i \cap \phi_j} Q(s_z, a_z = 0)$  on both sides, we can have the

431 left side:

$$\begin{aligned}
& Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + Q(s_j, a_j = 0) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 0)) + \sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) - \sum_{z \in \phi_i \cap \phi_j} Q(s_z, a_z = 0) \\
& = Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 0)) + \sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) - \sum_{z \in \phi_i \cap \phi_j} Q(s_z, a_z = 0) \\
& = Q(s_i, a_i = 1) + \sum_{i \in \phi_i^-} (Q(s_i, a_i = 1)) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 0)) \\
& = Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_i)
\end{aligned} \tag{9}$$

432 Similarly, the right side becomes

$$Q(s_j, a_j = 1) + \sum_{j \in \phi_j^-} (Q(s_j, a_j = 1)) + \sum_{i \notin \phi_j} (Q(s_i, a_i = 0)) = Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_j) \tag{10}$$

433 • **arm  $i$  and arm  $j$  are connected, and group  $\phi_i$  and  $\phi_j$  has overlap, i.e.,  $\phi_i \cap \phi_j \neq \emptyset$ , and**  
434  **$\{i, j\} \subset \phi_i \cap \phi_j$ .** This case is similar to the previous one, we add  $\sum_{z \notin \phi_i \wedge z \notin \phi_j} Q(s_z, a_z = 0) -$   
435  $\sum_{z \in \phi_i \cap \phi_j} Q(s_z, a_z = 0)$  on both sides, we can have the left side:  $Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_i)$  and the right side  
436  $Q(\mathbf{s}, \mathbf{a} = \mathbb{I}_j)$ .

437

□

438 We show that, using Theorem 1, selecting the top arms according to the  $\lambda$  value is guaranteed to  
439 maximize the cumulative long-term reward, thus proving it to be optimal.

440 However when it comes to the case where  $k > 1$ , selecting the top  $k$  arms according to the  $\lambda$  value  
441 is not guaranteed to be optimal. Let the  $\Phi$  denote the set of arms that are selected, i.e.,  $a_i = 2$  if  
442  $i \in \Phi$ . Because once the arm  $i$  is added to the selected arm set  $\Phi$ , the benefit of selecting arm  $j$  will  
443 also be influenced if the arm  $j$  has the shared connected neighbor arms with arm  $i$ , i.e.,  $\phi_i \cap \phi_j \neq \emptyset$ .  
444 To this end, finding the optimal solution is difficult, as we need to list all the possible solution sets.  
445 The non-asymptotic tight upper bound and non-asymptotic tight lower bound for getting the optimal  
446 solution are  $o(C(n, k))$  and  $\omega(N)$ , respectively.

447 We provide the proof for Theorem 2: *Proof.* When considering the influence of the shared neighbor  
448 nodes for two selected arms, then selecting arm  $i$  will influence the future benefit of selecting arm  
449  $j$  if arm  $i$  and arm  $j$  have the overlapped neighbor nodes, i.e.,  $\phi_i \cap \phi_j \neq \emptyset$ . This is because the  
450 calculation of  $\lambda_j$ , as some arms  $z \in \phi_i \cap \phi_j$  already receive the semi-active action  $a = 1$  due to the  
451 selection of arm  $i$ , the subsequent selection of arm  $j$  would not double introduce the benefit from  
452 those arms  $z$  who already included in  $\phi_i$ . However, if the top  $k$  arms ranked according to their  $\lambda$   
453 value do not have any overlaps in their connected neighbor nodes, i.e.,  $\phi_i \cap \phi_j = \emptyset$  for  $\forall i, j$ , where  
454 arm  $i$  and arm  $j$  are top  $k$  arms according to  $\lambda$  value. We can directly add those top  $k$  arms to the  
455 action set  $\Phi$ , and the solution is guaranteed to be optimal. Then we have the non-asymptotic tight  
456 lower bound for getting the optimal solution which is  $\omega(N)$ . Otherwise, if the top  $k$  arms ranked  
457 according to their  $\lambda$  value have any overlaps in their connected neighbor nodes, to get the optimal  
458 solutions, we need to list all possible combinations of the  $k$  arms, which have the  $C(n, k)$  cases, and  
459 computing the corresponding sum of the  $\lambda$  value. In this case, we can derive that the non-asymptotic  
460 tight upper bound for getting the optimal solution is  $o(C(n, k))$ . □

## 461 C Greedy algorithm when $k > 1$

462 When  $k > 1$ , it is difficult to compute the optimal solution as we might list all possible solutions, and  
463 the complexity is  $O(C(n, k))$ , Thus we provide a heuristic greedy algorithm to find the near-optimal  
464 solutions. The process to decide the selected arm set  $\Phi$  is as follows:

- 465 1. We first compute the independent  $\lambda$  value for each arm  $i$ , where  $i \in \{1, \dots, N\}$ , where  
466  $\lambda_i = Q(s_i, a_i = 1) - Q(s_i, a_i = 0) + \sum_{j \in \phi_i^-} (Q(s_i, a_i = 2) - Q(s_i, a_i = 0))$ ;
- 467 2. We add the arm with the top  $\lambda$  value to the set  $\Phi$ ;
- 468 3. We recompute the  $\lambda$  value for the each arm, note that we will remove  $Q(s_j, a_j)$  in the  $\lambda$   
469 equation if  $j \in \Phi$  or  $j \in \phi_j$  for  $\forall i \in \Phi$ ;

470 4. we add the arm with the top  $\lambda$  value to the set  $\Phi$ , and repeat the step 3 and 4 until we add  $k$   
471 arms to set  $\Phi$ .

472 The intuition of such a heuristic greedy algorithm is to add the arm that maximizes the marginal gain  
473 to the action. And the complexity for the greedy algorithm is  $O(\frac{(2N-k)*k}{2})$ .

## 474 **D Generating Student Models from Junyi and OLI Dataset**

475 In this section, we describe the features in Junyi and OLI dataset which we use in developing the  
476 transition matrices.

477 The datasets contain the following features which we use in various aspects to generate the student  
478 models and the network:

- 479 • **Topic & Knowledge Component Classification:** Items are classified into topics (Junyi) or  
480 KCs (OLI). This classification is employed to group items and establish the initial network.
- 481 • **Similarity:** The Junyi dataset offers expert ratings for exercise similarity, enabling a nuanced  
482 approach to form richer group memberships. High similarity scores group exercises together,  
483 irrespective of topic tags.
- 484 • **Difficulty:** The Junyi dataset provides expert ratings to determine the relative difficulty of  
485 exercise pairs. In the OLI dataset, we use the overall correct response rate as a measure of  
486 difficulty.
- 487 • **Rate of Correctness:** By analyzing student-exercise interactions, we calculate the frequency  
488 of correct answers for each question, offering insights into the improvement of knowledge  
489 over time.

### 490 **D.1 Active Transitions**

491 **Junyi Dataset** The Junyi dataset contains `earned_proficiency` feature which indicates if the  
492 student has achieved mastery of the topic based on Khan Academy's algorithm<sup>3</sup>. Thus, we take the  
493 number of attempts before `earned_proficiency=1` as  $P_{0,1}^2$ , and the errors made during mastery as  
494  $P_{1,0}^2$ .

495 **OLI Dataset** We possess records of students' accuracy on quiz questions after studying specific  
496 topics. To derive the transition matrix for the student with the corresponding action 2, we utilize the  
497 change in correctness rate before and after a learning intervention.

498 Given that proportion of correct attempts at time  $t$  as  $a^t$ , then  $a^{t+1} = P_{0,1}^2(1 - a^t) + P_{1,1}^2(a^t)$ . We  
499 use a linear regressor to estimate the respective  $P^2$ , constraining it to produce positive values and  
500 clipping the values to 0.99 when required.

### 501 **D.2 Passive Transitions**

502 To construct passive transitions for exercises, we utilize relative difficulty scores to determine  
503 transitions based on difficulty levels. We operate under the assumption that the difficulty of an  
504 exercise is linked to its likelihood of being forgotten, thereby resulting in a higher failure rate. More  
505 precisely, higher difficulty values of an exercise correspond to higher  $P_{1,0}^0$  values, indicating a greater  
506 likelihood of forgetting. The transition matrix for the passive action  $a = 0$  is then randomly generated,  
507 with the values influenced by the difficulty levels.

### 508 **D.3 Semi-active Transitions**

509 To derive semi-active transitions, the Junyi dataset contains similarity scores between two distinct  
510 exercises, quantifying their similarity on a 9-point Likert scale. Once the transition matrices are  
511 computed under the active action  $a = 2$  for all arms, we proceed to calculate the transition matrix

---

<sup>3</sup><http://david-hu.com/2011/11/02/how-khan-academy-is-using-machine-learning-to-assess-student-mastery.html>

512 for the semi-active action  $a = 1$ . This involves normalizing the similarity scores to the range  $[0, 1]$ ,  
 513 denoted as  $\sigma$ . For any chosen arm/topic, we can then compute its neighbor’s transition matrix under  
 514 the semi-active action  $a = 1$  with  $P_{0,1}^1 = \sigma(P_{0,1}^2)$ , where  $\sigma$  signifies the similarity proportion. It is  
 515 worth noting that an arm’s transition matrix for the semi-active action varies due to different neighbors  
 516 being selected — different neighbors correspond to different similarity scores.

517 To address this, we can store the transition matrix of semi-active actions for different neighbor  
 518 selection scenarios, preserving the flexibility of our algorithm. In this work, for simplicity, we opt  
 519 not to distinguish the impact of different neighbors being selected. Instead, we calculate the average  
 520 similarity for all arms in a group average them, and use the resultant average as  $\sigma$ .

521 For the OLI Statics dataset, we use a constant value of  $\sigma = 0.8$  since there are no similarity scores  
 522 available.

## 523 E Additional Experiment Results and Discussion

### 524 E.1 Comparing Different Network Setups

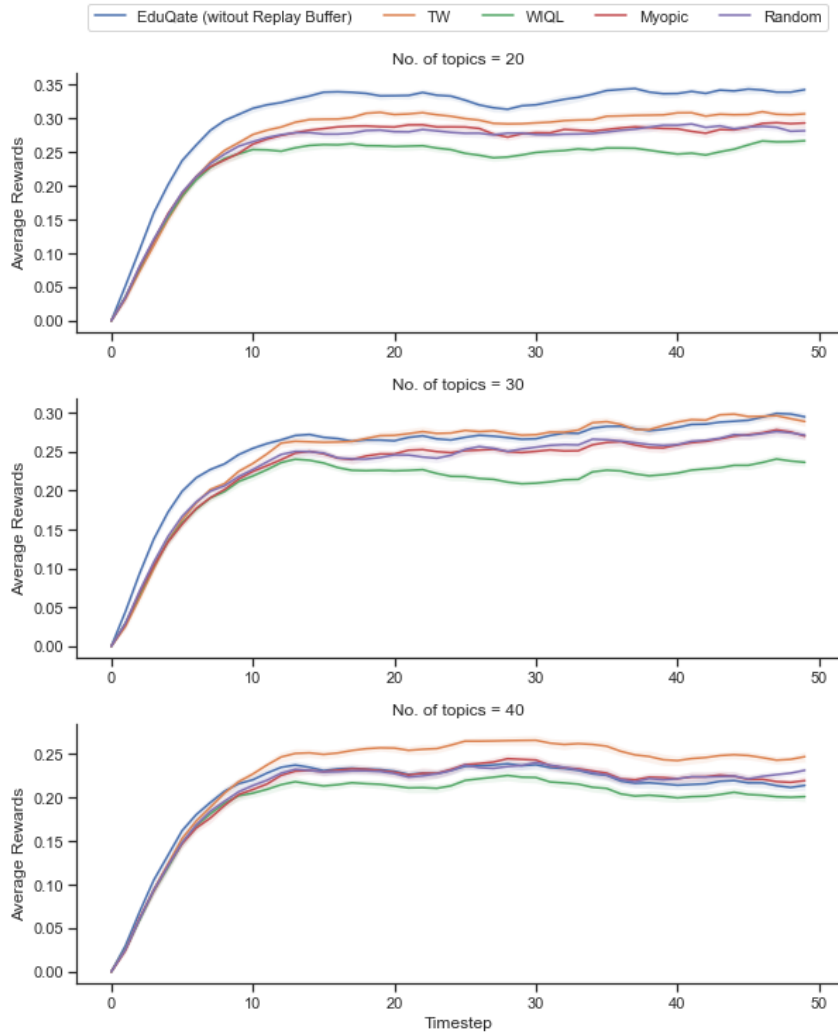


Figure 3: Average rewards for the respective algorithms, on the last episode of training. Note that as  $N_{topics}$  increase, the network effects are reduced, and most algorithms are not better than a random policy.

Table 3: Comparison of policies on synthetic dataset, with different network setups. Note that that as  $N_{topics}$  increase, the reliability of any algorithms decreases, as seen by the standard deviations of their average  $IB$ . EduQate- here refers to the EduQate algorithm without replay buffer.

$N_{topics}$	POLICY	$\mathbb{E}[IB] (\%) (\pm)$
20	WIQL	$-57.9 \pm 13.1$
	MYOPIC	$0.24 \pm 8.2$
	TW	$32.6 \pm 7.0$
	EDUQATE-	<b>100.0</b>
30	WIQL	$-292 \pm 1162$
	MYOPIC	$180 \pm 600$
	TW	$122 \pm 277$
	EDUQATE-	100
40	WIQL	$307 \pm 1069$
	MYOPIC	$212 \pm 526$
	TW	$4.34 \pm 1124$
	EDUQATE-	100

525 We present the results for different network setups in Table 3. We note that as the number of topics  
 526 approach the number of arms (i.e.  $N_{topics} = \{30, 40\}$ ), all algorithms perform in a highly unstable  
 527 manner, as reflected in the standard deviations presented. We emphasize here that the performance  
 528 of EduQate is dependent on the quality of the network it is working on, and tends to thrive in more  
 529 complex, yet realistic scenarios, such as the Junyi dataset presented in Figure 2. We present an  
 530 example of a graph generated when  $N_{topics} = 40$  in Figure 4, where we notice that many arms do  
 531 not belong to a group. Under this network, the EdNetRMAB can be approximated to a traditional  
 RMAB, where the arms are independent of each other.

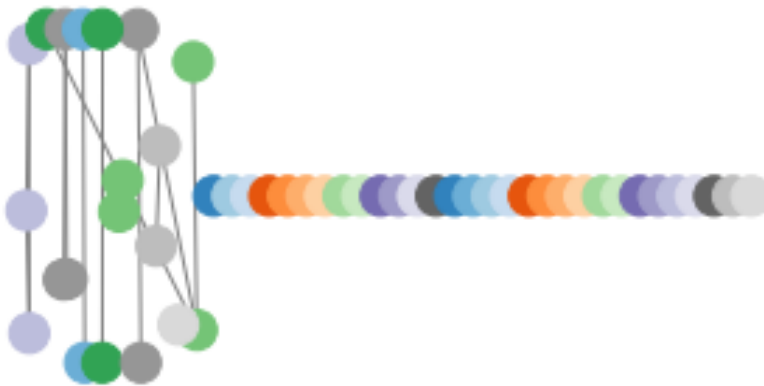


Figure 4: Synthetic network when  $N_{topics} = 40$ . Note that some arms are without group members, and do not receive benefits from networks. Node colors represent topic groups.



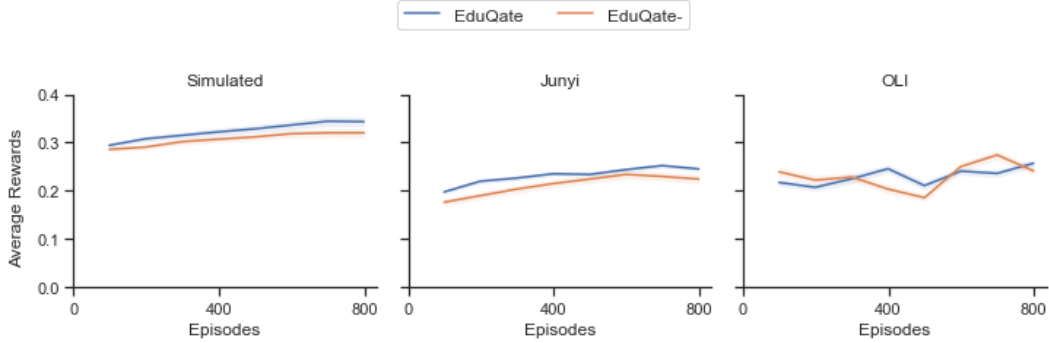


Figure 5: Average rewards across 800 episodes of training, across 30 seeds. EduQate- (orange) refers to the EduQate algorithm without replay buffer.

## 533 E.2 Ablation of Replay Buffer

Table 4: Comparison of EduQate with and without (EduQate-) Experience Replay Buffer policies across different datasets. Results reported are of the final episode of training.

POLICY	$\mathbb{E}[IB] (\%) \pm$		
	SYNTHETIC	JUNYI	OLI
EDUQATE-	104.74 $\pm$ 32.56	76.90 $\pm$ 4.72	107.30 $\pm$ 11.77
EDUQATE	100.0	100.0	100.0
POLICY	$\mathbb{E}[R] \pm$		
	SYNTHETIC	JUNYI	OLI
EDUQATE-	32.032 $\pm$ 0.469	22.133 $\pm$ 0.544	25.16 $\pm$ 0.432
EDUQATE	34.331 $\pm$ 0.489	24.527 $\pm$ 0.314	25.468 $\pm$ 0.469

534 We investigate the importance of the Experience Replay buffer in EduQate, as shown in Figure 5 and  
 535 Table 4. For the Simulated and Junyi datasets, EduQate without Experience Replay (EduQate-) does  
 536 not achieve the performance levels of the full EduQate algorithm within 800 episodes, highlighting the  
 537 importance of methods that aid Q-learning convergence. In real-world applications, slow convergence  
 538 can result in students experiencing a curriculum similar to a random policy, leading to sub-optimal  
 539 learning experiences during the early stages. This issue is known as the cold-start problem [3].  
 540 Future work in EdNetRMABs should explore methods to overcome cold-start problems and improve  
 541 convergence in Q-learning-based methods.

## 542 F Q-Learning

543 Q-learning [25] is a popular reinforcement learning method that enables an agent to learn optimal  
 544 actions in an environment by iteratively updating its estimate of state-action value,  $Q(s, a)$ , based on  
 545 the rewards it receives. The objective, therefore, to learn  $Q^*(s, a)$  for each state-action pair of an  
 546 MDP, given by:

$$Q^*(s, a) = r(s) + \sum_{s' \in S} P(s, a, s') \cdot V^*(s')$$

547 where  $V^*(s')$  is the optimal expected value of a state, is given by:

$$V^*(s) = \max_{a \in A} (Q(s, a))$$

548 Q-learning estimates  $Q^*$  through repeated interactions with the environment. At each time step  $t$ ,  
 549 the agent takes an action  $a$  using its current estimate of  $Q$  values and current state  $s$ , thus received a  
 550 reward of  $r(s)$  and new state  $s'$ . Q-learning then updates the current estimate using the following:

$$\begin{aligned}
 Q_{new}(s, a) \leftarrow & (1 - \alpha) \cdot Q_{old}(s, a) \\
 & + \alpha \cdot (r(s) \\
 & + \gamma \cdot \max_{a \in A} Q_{old}(s', a))
 \end{aligned}
 \tag{11}$$

551 where  $\alpha \in [0, 1]$  is the learning rate that controls updates, and  $\gamma$  is the discount on future rewards  
 552 associated with the MDP.

## 553 G Experiment Details and Hyperparameters

Category	Parameter	Value
Replay buffer	buffer_size	10000
	batch_size	64
WIQL/EduQate	$\gamma$	0.95
	$\alpha$	0.1

Table 5: Hyperparameters for Replay Buffer and Q-learning

## 554 H NeurIPS Paper Checklist

### 555 1. Claims

556 Question: Do the main claims made in the abstract and introduction accurately reflect the  
557 paper's contributions and scope?

558 Answer: [Yes]

559 Justification: We summarize our contributions and provide the scope of the paper in the  
560 abstract and introduction.

561 Guidelines:

- 562 • The answer NA means that the abstract and introduction do not include the claims  
563 made in the paper.
- 564 • The abstract and/or introduction should clearly state the claims made, including the  
565 contributions made in the paper and important assumptions and limitations. A No or  
566 NA answer to this question will not be perceived well by the reviewers.
- 567 • The claims made should match theoretical and experimental results, and reflect how  
568 much the results can be expected to generalize to other settings.
- 569 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
570 are not attained by the paper.

### 571 2. Limitations

572 Question: Does the paper discuss the limitations of the work performed by the authors?

573 Answer: [Yes]

574 Justification: Limitations were discussed in the final section.

575 Guidelines:

- 576 • The answer NA means that the paper has no limitation while the answer No means that  
577 the paper has limitations, but those are not discussed in the paper.
- 578 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 579 • The paper should point out any strong assumptions and how robust the results are to  
580 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
581 model well-specification, asymptotic approximations only holding locally). The authors  
582 should reflect on how these assumptions might be violated in practice and what the  
583 implications would be.
- 584 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
585 only tested on a few datasets or with a few runs. In general, empirical results often  
586 depend on implicit assumptions, which should be articulated.
- 587 • The authors should reflect on the factors that influence the performance of the approach.  
588 For example, a facial recognition algorithm may perform poorly when image resolution  
589 is low or images are taken in low lighting. Or a speech-to-text system might not be  
590 used reliably to provide closed captions for online lectures because it fails to handle  
591 technical jargon.
- 592 • The authors should discuss the computational efficiency of the proposed algorithms  
593 and how they scale with dataset size.
- 594 • If applicable, the authors should discuss possible limitations of their approach to  
595 address problems of privacy and fairness.
- 596 • While the authors might fear that complete honesty about limitations might be used by  
597 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
598 limitations that aren't acknowledged in the paper. The authors should use their best  
599 judgment and recognize that individual actions in favor of transparency play an impor-  
600 tant role in developing norms that preserve the integrity of the community. Reviewers  
601 will be specifically instructed to not penalize honesty concerning limitations.

### 602 3. Theory Assumptions and Proofs

603 Question: For each theoretical result, does the paper provide the full set of assumptions and  
604 a complete (and correct) proof?

605 Answer: [Yes]

606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659

Justification: Proofs are provided in Appendix 4.1.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Experiment details are provided in both the main body and the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

**5. Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711

Answer: [Yes]

Justification: Code and the transition matrices are provided as supplementary materials.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Relevant details are provided in the main body, as well as the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: In our experiments, we report and display the standard error across all seeds.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
  - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
  - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 720 8. Experiments Compute Resources

721 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

724 Answer: [Yes]

725 Justification: The current paper only requires CPU-level of compute and is mentioned in the Experiment section.

727 Guidelines:

- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- The answer NA means that the paper does not include experiments.
  - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 736 9. Code Of Ethics

737 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

739 Answer: [Yes]

740 Justification: All datasets used were anonymized by the respective authors.

741 Guidelines:

- 742
- 743
- 744
- 745
- 746
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 747 10. Broader Impacts

748 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

750 Answer: [Yes]

751 Justification: The current work has positive implications for applied machine learning in education settings, and is discussed in the Introduction section. As far as we can see, we don't think there are negative impacts for education.

754 Guidelines:

- 755
- 756
- 757
- 758
- 759
- 760
- 761
- The answer NA means that there is no societal impact of the work performed.
  - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The current paper does not release any new assets.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Code [17] and datasets [7, 4] were appropriately cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- 814 • If this information is not available online, the authors are encouraged to reach out to  
815 the asset’s creators.

816 **13. New Assets**

817 Question: Are new assets introduced in the paper well documented and is the documentation  
818 provided alongside the assets?

819 Answer: [NA]

820 Justification: [NA]

821 Guidelines:

- 822 • The answer NA means that the paper does not release new assets.
- 823 • Researchers should communicate the details of the dataset/code/model as part of their  
824 submissions via structured templates. This includes details about training, license,  
825 limitations, etc.
- 826 • The paper should discuss whether and how consent was obtained from people whose  
827 asset is used.
- 828 • At submission time, remember to anonymize your assets (if applicable). You can either  
829 create an anonymized URL or include an anonymized zip file.

830 **14. Crowdsourcing and Research with Human Subjects**

831 Question: For crowdsourcing experiments and research with human subjects, does the paper  
832 include the full text of instructions given to participants and screenshots, if applicable, as  
833 well as details about compensation (if any)?

834 Answer: [NA]

835 Justification: [NA]

836 Guidelines:

- 837 • The answer NA means that the paper does not involve crowdsourcing nor research with  
838 human subjects.
- 839 • Including this information in the supplemental material is fine, but if the main contribu-  
840 tion of the paper involves human subjects, then as much detail as possible should be  
841 included in the main paper.
- 842 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
843 or other labor should be paid at least the minimum wage in the country of the data  
844 collector.

845 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
846 Subjects**

847 Question: Does the paper describe potential risks incurred by study participants, whether  
848 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
849 approvals (or an equivalent approval/review based on the requirements of your country or  
850 institution) were obtained?

851 Answer: [NA]

852 Justification: [NA]

853 Guidelines:

- 854 • The answer NA means that the paper does not involve crowdsourcing nor research with  
855 human subjects.
- 856 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
857 may be required for any human subjects research. If you obtained IRB approval, you  
858 should clearly state this in the paper.
- 859 • We recognize that the procedures for this may vary significantly between institutions  
860 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
861 guidelines for their institution.
- 862 • For initial submissions, do not include any information that would break anonymity (if  
863 applicable), such as the institution conducting the review.