

**International Journal of Reasoning-based Intelligent Systems**

ISSN online: 1755-0564 - ISSN print: 1755-0556

<https://www.inderscience.com/ijris>

---

**A multi-agent reinforcement learning framework for educational resource management optimisation**

Yuanyuan Feng

**DOI:** [10.1504/IJRIIS.2026.10075912](https://doi.org/10.1504/IJRIIS.2026.10075912)

**Article History:**

Received:	01 November 2025
Last revised:	24 November 2025
Accepted:	24 November 2025
Published online:	28 January 2026

# A multi-agent reinforcement learning framework for educational resource management optimisation

Yuanyuan Feng

School of Marxism,  
Xinxiang Vocational and Technical College,  
Xinxiang, 453000, China  
Email: fyyfyy1988@163.com

**Abstract:** This paper addresses the challenges of dynamic adaptation and multi-objective optimisation in educational resource management by proposing a novel multi-agent reinforcement learning framework. The framework utilises centralised teacher agents for global coordination and decentralised student agents for personalised recommendations. It incorporates an integrated optimisation mechanism to balance learning effectiveness, fairness, and efficiency, enhanced by a curriculum learning strategy that progressively trains agents from basic to complex tasks. Comprehensive experiments on the Junyi academy and assessments public datasets demonstrate that this approach improves knowledge mastery by 18% and enhances resource allocation fairness by 22% compared to traditional baseline methods, effectively narrowing the achievement gap among learners with different initial capabilities. The study provides an innovative and effective solution for next-generation intelligent education systems. The framework is designed to be applicable across various educational stages, including K-12 and higher education, to ensure broad relevance.

**Keywords:** multi-agent reinforcement learning; MARL; educational resource management; personalised learning; multi-objective optimisation; fairness.

**Reference** to this paper should be made as follows: Feng, Y. (2026) 'A multi-agent reinforcement learning framework for educational resource management optimisation', *Int. J. Reasoning-based Intelligent Systems*, Vol. 18, No. 7, pp.32–43.

**Biographical notes:** Yuanyuan Feng is a Lecturer in the School of Marxism at Xinxiang Vocational and Technical College. She obtained her Bachelor degree from Pingdingshan University in 2012, and Master's degree from Dali University in 2015. In 2023, she obtained her PhD in Education from University of Perpetual Help System DALTA. She has published three papers. Her research interests include educational management, Marxist theory and reinforcement learning theory.

## 1 Introduction

With the gradual promotion of process of educational informatisation, large amounts of various kinds of educational resources are accumulated in intelligent learning platform gradually (Ostrovskaya, 2022). The ongoing advancement of educational informatisation has fundamentally transformed the traditional modes of knowledge dissemination, thereby facilitating the systematic accumulation of diverse digital resources. This includes interactive video content and adaptive exercise modules, which collectively serve to significantly broaden and enrich the learning opportunities available to students across various domains. How to achieve efficient, fair, and personalised management of these resources within a large-scale, heterogeneous resource environment has become a key issue in current intelligent education field (Khudhur et al., 2024). Conventional educational resource allocation methods primarily rely on static rule filtering, collaborative filtering (CF) or content based recommendation method (Han and Guo, 2025). Although

above methods gained certain achievements at initial stage, they still have intrinsic defects such as bad dynamic adaptation, ignoring resource interconnectivity and hard to balance multi-objective optimisation requirements (Da'U and Salim, 2020). Especially in large scale, multi-role and multi-task educational scene, conventional method can't deeply analyse and mining the dynamic changing knowledge state and cognition demand of learners (Kong et al., 2025). These limitations result in suboptimal learning outcomes, with studies showing that traditional methods can lead to performance gaps of up to 30% between different learner groups (Farhadi and Winton, 2024). All above problems show that the current educational resource management system should adopt new method which more intelligent, adaptive and collaborative (Boubaker et al., 2025).

In recent years, reinforcement learning technologies, particularly multi-agent reinforcement learning (MARL) has provided new solutions for the management of educational resources due to its obvious advantages in the distribution of perception, decision-making and collaboration in the

environment with multiple uncertain factors. Taking various educational resources (such as exercises, videos, courseware) or teaching roles (such as teachers, teaching assistants, peer students) as independent agents, MARL constructs multiple roles in actual teaching environment for simulation and collaboration interactions, and continuously optimises the allocation strategy of educational resources (Kaya and Nder, 2025). The MARL framework operates under the partially observable Markov decision process framework, where each agent maintains its own policy while coordinating through centralised training. There are some typical applications as follows: Tsinghua university's Massive AI empowered Course (MAIC) system uses large models and multi-agent technology to build a complete scenario of fully AI empowered classroom, including teachers (Sheikh et al., 2025), teaching assistants and peers. It simulates teaching paths of different levels according to the learning situation of students, which improves students' learning motivation and greatly improves learning efficiency (Tsai et al., 2015). AI counsellor of Chongqing University 'Runxin' uses multi-agent cross-system collaboration platform to realise deep connection and integration between multiple business systems, including academic affairs, student service and logistics, to break through data isolation and provide accurate service for students in the whole development process. In addition, based on the educational AI platform developed by Xi'an Jiaotong University, 'JiaoXiaoZhi', the teachers can customise their own personalised AI children. Nearly 100 intelligent modules and plugins are integrated, and the teaching management burden has been reduced by more than half in practice (Jin and Feng, 2014). These applications show that the MARL technology has potential in the dynamic allocation of educational resources and personalised services (Zhang, 2024). The referenced application cases, such as the Tsinghua MAIC system, are widely regarded as pioneering initiatives within the field of AI-enhanced education. They are considered representative due to their demonstrated capacity to showcase the scalability and adaptability of MARL in addressing complex, real-world educational challenges and dynamic learning environments.

Despite promising applications, MARL implementation in educational resource management faces four critical technical challenges: partial observability limits global optimisation, multi-objective trade-offs lack systematic coordination, sparse reward signals hinder policy learning, and computational complexity constraints real-world deployment. Firstly, the challenge of partial observability remains (Mladenovici et al., 2024). In actual educational scenarios, each resource agent can only observe part of the system state (e.g., the learning record of one student) and has no idea about the global distribution of resources and the collective learning process. This limits the perspective for making decisions, which hinders the possibility of reaching globally optimal resource allocation solutions (Lambiase et al., 2025). Second, the problem of multi-objective trade-offs still needs to be solved. Educational

resource management itself involves many conflicting goals, such as learning effectiveness, the fairness of resource allocation and system operation, and personalised satisfaction. Most existing MARL algorithms have focused on optimising one reward function and lack systematic designs for multi-objective collaborative optimisation (He et al., 2024). Third, sparse rewards and delayed feedback greatly limit the learning efficiency of policies. Especially in long-term learning path planning, agents only get reward signals at some specific teaching moments (such as exams, submission of assignments), which leads to slow model convergence and unstable policy. Moreover, existing systems still have many drawbacks in terms of cross-scenario generalisation capability and controlling computational overhead, which make them not suitable for meeting various demands of disciplines, stages of education, and teaching models (Guimares Iglesias et al., 2024).

In order to solve these problems, this study aims at theoretical and applied innovations of MARL in educational resource scheduling. Specifically, it is explored how to construct an efficient, scalable and fair resource allocation framework under partially-observable, multi-objective and sparse reward constraints (Ge et al., 2018). The main innovation of this study includes: To design a MARL framework with centralised teacher agents and decentralised student agents. Through introducing centralised teacher agents to perceive global states and distribute rewards, it can guide decentralised student agents (resource agents) to execute personalised resource recommendation based on local observations. This design enhances the global coordination efficiency of the entire system while preserving the decentralised decision-making of the student agents. To design a training mechanism based on curriculum learning and multi-objective optimisation. By designing curriculums in different stages, it reinforces the learning of basic resource allocation, personalised recommendation and multi-objective trade-off strategies for agents in a reinforcement way. This mechanism effectively alleviates the sparse reward problem, and further introduces the theory of constraint optimisation to balance learning effectiveness, fairness and efficiency; Deep adaptation and validation on open educational datasets to ensure reliable generalisation and reproducibility of the model under real-world data. Through the above innovations, this research aims to provide a novel, reliable and scalable resource management paradigm to promote the deep cooperation of 'AI + Education', and build a future educational ecosystem that is smarter, more personalised and more fair.

## 2 Related research

### 2.1 The evolution and current state of educational resource management

The evolution of educational resource management systems has undergone a significant transformation from static repositories to dynamic recommendation systems.

Early systems primarily relied on manual annotation and classification, gradually evolving into automated recommendation methods based on content filtering and CF (Zhang et al., 2025). To illustrate the concept of static repositories in their historical context, early educational resource management systems often functioned essentially as digital libraries. In these systems, resources were typically manually categorised, stored, and retrieved, operating without the dynamic, algorithm-driven updates and personalisation capabilities characteristic of contemporary intelligent platforms. Among these, CF achieves recommendations by calculating similarity within user-resource interaction matrices. User similarity can be expressed as:

$$\text{sim}(u_i, u_j) = \frac{\sum_{r \in R} (r_{u_i, r} - \bar{r}_{u_i})(r_{u_j, r} - \bar{r}_{u_j})}{\sqrt{\sum_{r \in R} (r_{u_i, r} - \bar{r}_{u_i})^2} \sqrt{\sum_{r \in R} (r_{u_j, r} - \bar{r}_{u_j})^2}},$$

while resource similarity is calculated as:

$$\text{sim}(r_i, r_j) = \frac{r_i \cdot r_j}{|r_i| |\pi r_j|}.$$

However, these traditional methods face significant cold-start problems and cannot adapt to the dynamic changes in learners' knowledge states. In recent years, with the advancement of artificial intelligence technology, intelligent learning resources have begun to exhibit new characteristics such as evolution, sharing, and adaptability. The generative multi-agent guided learning system proposed. Employs a 'teaching-learning-guiding' triadic agent structure to organically integrate direct and indirect experiences. Its system utility function can be formalised as:

$$U(s_t, a_t) = \sum_{i=1}^n \omega_i R_i(s_t, a_t),$$

where  $\omega_i$  denotes the weight coefficient for each objective, and  $R_i(s_t, a_t)$  denotes the corresponding reward function. conducted an in-depth study on the evolutionary framework of intelligent learning resources from a multi-objective optimisation perspective, formalising the optimisation objective as:

$$\max_{x \in X} [f_1(x), f_2(x), \dots, f_m(x)]^T,$$

where  $f_i(x)$  denotes the  $i^{\text{th}}$  optimisation objective encompassing multiple dimensions such as learner cognitive level matching, media type preference alignment, learning content preference matching, and learning time expectation alignment.

## 2.2 Applications of MARL in education

MARL provides effective solutions for decision-making problems in complex environments through collaborative and competitive mechanisms among agents. Within the Markov game framework, a multi-agent system can be

rigorously defined as the tuple  $\langle N, S, A, P, R, O, \gamma \rangle$ , where  $N = 1, 2, \dots, n$  denotes the set of agents,  $S$  denotes the global state space,  $A = A_1 \times A_2 \times \dots \times A_n$  forms the joint action space,  $P: S \times A \rightarrow S$  is the state transition probability function,  $R = [R_1, R_2, \dots, R_n]$  represents the set of reward functions,  $O = [O_1, O_2, \dots, O_n]$  denotes the set of observation functions, and  $\gamma \in [0, 1]$  is the discount factor. In educational applications, empirical research. Shows that the multi-agent system can effectively improve students' reasoning and evaluation ability in incentive-based learning activities. the learning effect gain of

$$\Delta L = \frac{1}{n} \sum_{i=1}^n (p_{post}^i - p_{pre}^i)$$

can be calculated as follows:  $p_{pre}^i$  and  $p_{post}^i$  represent students' score in the pre-test and post-test. A fundamental distinction from single-agent models lies in this framework's inherent capacity to enable multiple autonomous agents to interact and learn collaboratively within a shared environment. This design more accurately mimics the complex, interactive nature of real-world educational settings, where multiple stakeholders and resources coexist and influence each other. As for learning path recommendation, through the cooperation between teacher and student agent, the recommendation process is decomposed into two levels. teacher agent makes the planning for knowledge point sequence at the macro level, whose value function is defined as

$$V_T(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R_T(s_t, a_t^T) \right];$$

Student agent makes the recommendation for certain exercise at the micro level, whose q-function is updated as

$$Q_S(s, a) \leftarrow Q_S(s, a) + \alpha [r + \gamma \max_{a'} Q_S(s', a') - Q_S(s, a)].$$

This hierarchical design can not only greatly reduce exploration space, but also guarantee that all precedence constraints between knowledge points are satisfied. The state transition probability satisfies

$$P(s'|s, a) = P(s'|s, a^T, a^S) = \prod_{i=1}^k P(s'_i | s_i, a_i^S),$$

where  $k$  represents the total amount of knowledge points. However, current marl methods still confront serious problems of partial observability in educational application. Different from professional games, the observation space  $O_i$  of each agent is only a subset of global state  $S$ . This information asymmetry will bring certain biases to the decision-making of agents, and further needs more effective communication and coordination mechanism design to improve the performance of the system.

### 2.3 Challenges of multi-objective optimisation in educational resource allocation

Actually a tricky multi-objective optimisation issue that should take care of a lot of competing targets. Traditional single objective optimisation strategies are usually invalid in resolving this issue, and thus the multi-target optimisation concept is particularly important in this field (Saxena and Deb, 2007). The multi-objective optimisation problem of educational resource allocation can be formalised as

$$\min_{x \in \Omega} [f_1(x), f_2(x), \dots, f_k(x)]^T,$$

where  $\Omega$  denotes the feasible domain of decision variables, and  $f_i(x)$  represents the  $i^{\text{th}}$  objective function. Common objectives include maximising learning effectiveness:

$$f_1(x) = \frac{1}{m} \sum_{j=1}^m (1 - k_j^{\text{final}});$$

optimising resource allocation fairness:

$$f_2(x) = \text{Gini}(\Phi);$$

and maximising resource utilisation efficiency

$$f_3(x) = \frac{\sum_{i=1}^n c_i x_i}{\sum_{j=1}^m b_j},$$

where  $k_j^{\text{final}}$  denotes learner  $j$ 's final knowledge mastery level,  $\Phi = [\phi_1, \phi_2, \dots, \phi_m]$  denotes the resource value vector obtained by each learner,  $c_i$  is the cost of resource  $i$ , and  $b_j$  is the budget constraint for learner  $j$ . The concept of pareto optimality is crucial in this multi-objective optimisation problem. A solution  $x^* \in \Omega$  is called pareto optimal if and only if there exists no other solution  $x \in \Omega$  such that for all  $i = 1, 2, \dots, k$  holds  $f_i(x) \leq f_i(x^*)$ , and there exists at least one  $j$  such that  $f_j(x) < f_j(x^*)$ . The weighted sum method is a classic technique for handling multi-objective optimisation problems by transforming the multi-objective problem into a single-objective problem:

$$\min_{x \in \Omega} \sum_{i=1}^k w_i f_i(x),$$

where  $w_i$  is the weight of the  $i^{\text{th}}$  objective, satisfying  $w_i \geq 0$  and  $\sum_{i=1}^k w_i = 1$ . While it is acknowledged that other

important objectives, such as enhancing learner engagement and promoting long-term knowledge retention, are highly relevant in educational contexts, the present work deliberately focuses on the core triumvirate of learning effectiveness, fairness, and efficiency. This focused approach ensures analytical clarity and facilitates a more in-depth investigation into balancing these primary, and often

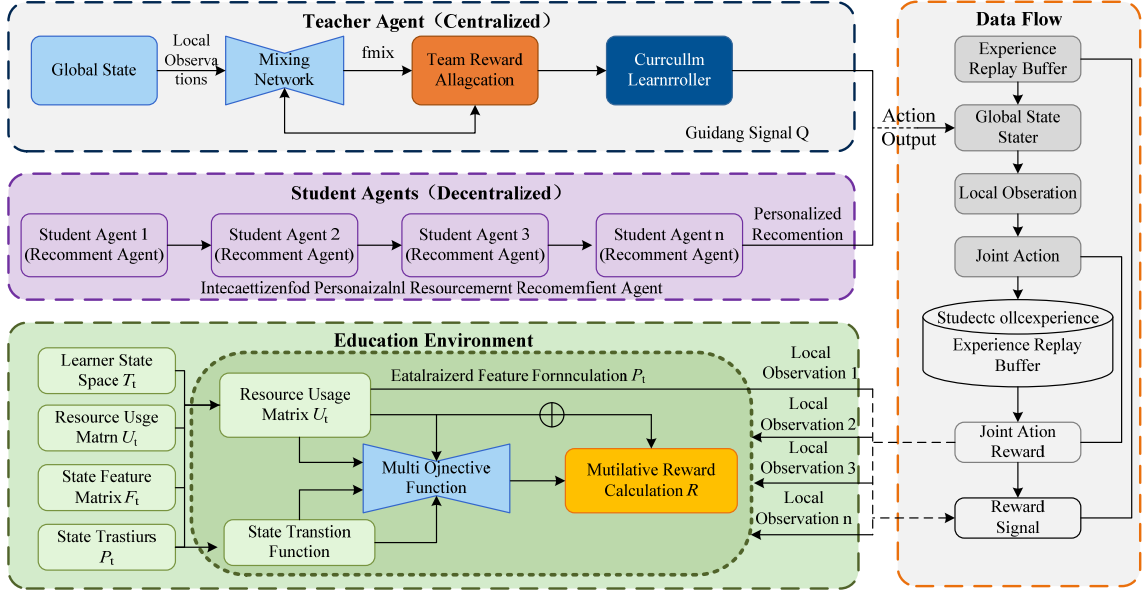
competing, goals. However, this approach struggles to identify non-convex Pareto fronts, necessitating more advanced multi-objective optimisation algorithms to address the complex challenges in educational resource allocation (Kazmi et al., 2025).

## 3 Technology and methods

### 3.1 Problem modelling

We therefore articulate all of the educational resource manager system in a Markov video game approach, formulated as the set  $\langle N, S, A, P, R, O, \gamma \rangle$ . Where  $N = 1, 2, \dots, n$  represents the collection of resource brokers, i.e., all kinds of educational source recommendation brokers;  $S$  represents the global state space, i.e., all learners' understanding states, resource usage, learners' traits and different related details;  $A = A_1 \times A_2 \times \dots \times A_n$  represents the joint action space, where  $A_i$  is the individual action space for agent  $i$ ;  $P: S \times A \rightarrow S$  denotes the state transition probability function, describing the dynamic evolution of the system state under action influence;  $R = [R_1, R_2, \dots, R_n]$  represents the reward function set, where  $R_i: S \times A \rightarrow \mathbb{R}$  is the individual reward function of agent  $i$ ;  $O = [O_1, O_2, \dots, O_n]$  denotes the set of observation functions, where  $O_i: S \rightarrow O_i$  is the local observation function for agent  $i$ ;  $\gamma \in [0, 1]$  is the discount factor, used to balance the importance of immediate rewards versus future rewards. The state space  $S$  comprises three meticulously designed components: knowledge states  $K_t$  tracking 15 distinct mathematical competencies across a  $[0, 1]$  proficiency scale, resource utilisation patterns  $U_t$  recording temporal engagement metrics, and learner profiles  $F_t$  encompassing cognitive styles (visual/auditory/kinesthetic), prior achievement levels, and learning pace indicators. This comprehensive state representation enables the system to capture both immediate learning needs and longitudinal development trajectories.

Specifically, the global state  $s_t \in S$  at time step  $t$  can be decomposed into three core components:  $\underline{s}_t = [K_t, U_t, F_t]$ . Here,  $K_t = [k_1^t, k_2^t, \dots, k_m^t]$  represents the knowledge mastery vector for all  $m$  learners at time  $t$ , where  $k_j^t \in [0, 1]$  denotes learner  $j$ 's knowledge mastery level at time  $t$ ;  $U_t$  is an  $m \times n$  resource usage matrix recording historical resource allocation, where  $U_t(j, i)$  denotes learner  $j$ 's cumulative usage of resource  $i$ ;  $F_t$  is the learner feature matrix containing static attributes such as cognitive level, learning style preferences, and media type preferences. The initialisation of the knowledge mastery vector, a critical component of the state space, is meticulously handled to establish a realistic baseline. This is achieved by utilising pre-assessment data where available; for new learners without historical data, a set of carefully chosen default values is employed to ensure a consistent and fair starting point for the learning process.

**Figure 1** CTDS-MARL framework schematic diagram (see online version for colours)**Algorithm 1** CTDS-MARL training loop

---

Input: Environment, Teacher Agent, Student Agents, Curriculum Stages

Initialise replay buffers, networks, and hyperparameters

for each curriculum stage do:

for episode = 1 to M do:

Initialise environment state  $s$

for  $t = 1$  to  $T$  do:

Teacher computes global guidance  $Q\_teacher(s)$

Each student agent  $i$  selects action  $a\_i$  based on local observation  $o\_i$

Execute joint action  $a$ , observe reward  $r$ , next state  $s'$

Store transition  $(s, a, r, s')$  in replay buffer

Sample batch and update student networks

Update teacher network with global loss

end for

end for

if stage transition criteria met: proceed to next stage

end for

---

At time  $t$ , the action  $a_i^t \in A_i$  of each resource agent  $i$  is defined as the probability distribution for recommending this resource to different learners, i.e.,

$$a_i^t = [p_{i1}^t, p_{i2}^t, \dots, p_{im}^t],$$

where  $p_{ij}^t \in [0, 1]$  represents the probability that agent  $i$  recommends the resource to learner  $j$ , satisfying  $\sum_{j=1}^m p_{ij}^t = 1$ .

The joint action  $a_t = (a_1^t, a_2^t, \dots, a_n^t)$  constitutes the collective decision of all agents at time  $t$ . The reward function is designed as a multi-objective form, where the reward

$R_i(s_t, a_t)$  for each agent  $i$  is obtained by the weighted sum of three core components:

$$R_i(s_t, a_t) = \alpha R_{learn}(s_t, a_t) + \beta R_{fair}(s_t, a_t) + \gamma R_{efficiency}(s_t, a_t) \quad (1)$$

where  $\alpha, \beta, \gamma \in [0, 1]$  are weighting hyperparameters satisfying  $\alpha + \beta + \gamma = 1$ , used to balance the importance of different objectives;  $R_{learn}(s_t, a_t)$  measures learning effectiveness improvement,  $R_{fair}(s_t, a_t)$  evaluates resource allocation fairness, and  $R_{efficiency}(s_t, a_t)$  quantifies resource utilisation efficiency.

### 3.2 CTDS-MARL framework

The proposed centralised teacher and decentralised student multi-agent reinforcement learning framework (CTDS-MARL) achieves an organic integration of global coordination and local decision-making through a two-layer architecture. The teacher agent takes the global state  $s_t$  as input to learn a team reward allocation strategy and generates personalised q-functions to guide each student agent. The teacher agent maintains a hybrid network to compute each agent's contribution to the team reward, with its q-function defined as:

$$Q_{teacher}(s_t, a_t) = f_{mix} \left( \begin{matrix} Q_1(s_t, a_1^t), Q_2(s_t, a_2^t), \dots, \\ Q_n(s_t, a_n^t); s_t \end{matrix} \right) \quad (2)$$

where  $f_{mix}$  is a mixed function implemented by a multi-layer perceptron, comprising, the teacher agent's hybrid network is a multi-layer perceptron (MLP) comprising two hidden layers, each with 128 neurons, using the ReLU activation function, using the rectified linear unit activation function. This architecture was chosen through empirical validation on a held-out validation set, where it provided a balance between representational capacity and computational efficiency. The MLP structure was selected for implementing the teacher agent's hybrid network based on a balanced consideration of its functional capabilities and computational demands. MLPs are well-known for their strong capacity to approximate complex nonlinear relationships, which is essential for this task, while simultaneously maintaining manageable computational efficiency for large-scale applications.  $Q_i(s_t, a_i^t)$  is the q-value estimate for agent  $i$ , representing the expected cumulative reward from executing action  $a_i^t$  in global state  $s_t$ .

Student agents utilise distributed proximal policy optimisation with a clip range of  $\epsilon = 0.2$ , a value commonly adopted in PPO implementations to ensure stable policy updates, and a value function coefficient of  $c_{value} = 0.5$  to balance policy and value learning. make decisions based on local observations  $o_i^t$ , but the learning of their q-functions is guided by the teacher agent. each student agent's local observation  $o_i^t$  comprises three components:  $o_i^t = [k_i^t, u_i^t, f_i^t]$ , where  $k_i^t$  is a subset of the learner's knowledge state related to resource  $i$ ,  $u_i^t$  is the historical usage of resource  $i$ , and  $f_i^t$  is a subset of learner features related to resource  $i$ . The objective of each learner agent  $i$  is to minimise the following loss function:

$$L_i(\theta_i) = \mathbb{E}(o_i^t, a_i^t, r_i^t, o_i^{t+1}) \sim D \left[ (Q_i(o_i^t, a_i^t; \theta_i) - y_i^t)^2 \right] \quad (3)$$

where  $\theta_i$  represents the q-network parameters of agent  $i$ ;  $D$  denotes the experience replay buffer, storing historical experience tuples  $(o_i^t, a_i^t, r_i^t, o_i^{t+1})$ ;  $y_i^t$  is the target q-value, computed as: The teacher agent employs a dual-stream

network architecture processing global state information through separate feature extractors for cognitive states (3-layer CNN) and resource utilisation patterns (LSTM with 64 hidden units). Student agents utilise distributed proximal policy optimisation with clip range  $\epsilon = 0.2$  and value function coefficient  $c_{value} = 0.5$ . We implement prioritised experience replay with importance sampling correction, using temporal difference error  $\epsilon_{fair}$  for priority assignment with stochastic rank-based sampling.

$$y_i^t = r_i^t + \gamma Q_{teacher}(s_{t+1}, a_{t+1}) \Big|_{at+1} = \pi_{teacher}(s_{t+1}) \quad (4)$$

where

$$\pi_{teacher}(s_{t+1}) = \arg \max_a Q_{teacher}(s_{t+1}, a)$$

represents the teacher agent's optimal action selection policy at state  $s_{t+1}$ .

The teacher agent is trained by minimising the following loss function:

$$\begin{aligned} L_{teacher}(\theta_{teacher}) &= \mathbb{E}(s_t, a_t, r_t, st+1) \\ &\sim D \left[ \begin{pmatrix} Q_{teacher}(s_t, a_t; \theta_{teacher}) \\ -y_{teacher}^t \end{pmatrix}^2 \right] \end{aligned} \quad (5)$$

where

$$y_{teacher}^t = r_t + \gamma \max_{a_{t+1}} Q_{teacher}(s_{t+1}, a_{t+1}; \theta_{teacher}^-)$$

where  $\theta_{teacher}^-$  represents the target network parameters periodically cloned from the current network parameters  $\theta_{teacher}$  to stabilise training.

### 3.3 Multi-objective optimisation mechanism

To balance multiple competing objectives in the allocation of educational resources, we designed a multi-objective optimisation mechanism based on constraint optimisation theory. First, we explicitly defined the three core reward functions as follows:

Learning effectiveness reward measures the impact of resource recommendations on learners' knowledge mastery, calculated as:

$$R_{learn}(s_t, a_t) = \frac{1}{m} \sum_{j=1}^m (k_j^{t+1} - k_j^t) \quad (6)$$

where  $k_j^{t+1}$  and  $k_j^t$  denote learner  $j$ 's knowledge mastery levels at time  $t+1$  and  $t$ , respectively, where  $m$  represents the total number of learners.

Fairness reward evaluation assesses the equity of resource allocation across different learner groups, defined based on the Gini coefficient:

$$R_{fair}(s_t, a_t) = 1 - Gini(\Phi) \quad (7)$$

where  $\Phi = [\phi_1, \phi_2, \dots, \phi_m]$  denotes the total resource value vector obtained by each learner,  $\phi_j = \sum_{i=1}^n v_i \cdot a_{ij}^t$  represents

the total value of all resources acquired by learner  $j$ , and  $v_i$  is the intrinsic value of resource  $i$ ;  $Gini(\Phi)$  is the Gini coefficient, calculated as. The Gini coefficient is a well-established and robust metric borrowed from economics, where it has been extensively used to quantify levels of inequality within distributions, most famously for income and wealth. Its proven theoretical foundation and interpretability make it particularly appropriate for evaluating fairness in the distribution of educational resources among a diverse population of learners.

$$Gini(\Phi) = \frac{\sum_{j=1}^m \sum_{k=1}^m |\phi_j - \phi_k|}{2m \sum_{j=1}^m \phi_j} \quad (8)$$

The gini coefficient ranges from  $[0, 1]$ , with lower values indicating more equitable distribution. Efficiency rewards measure the cost-benefit ratio of resource utilisation: Efficiency rewards measure the cost-benefit ratio of resource utilisation:

$$R_{efficiency}(s_t, a_t) = \frac{\sum_{j=1}^m \mathbb{I}(k_j^{t+1} > \tau)}{\sum_{i=1}^n c_i \sum_{j=1}^m a_{ij}^t} \quad (9)$$

where  $\mathbb{I}$  denotes the indicator function, returning 1 when the condition holds and 0 otherwise;  $\tau$  represents the knowledge threshold is set to 0.7, indicating a ‘proficient’ level of understanding, a standard benchmark in educational assessment,  $c_i$  is the allocation cost of resource  $i$ ;  $a_{ij}^t$  is a binary variable indicating whether resource  $i$  is recommended to learner  $j$ .

The multi-objective optimisation problem is formulated as follows:

$$\max_{a \in A} [R_{learn}(a), R_{fair}(a), R_{efficiency}(a)]^T \quad (10)$$

To address this multi-objective optimisation problem, we employ a constrained optimisation approach, incorporating fairness and efficiency as constraints while prioritising learning effectiveness as the primary optimisation objective:

$$\max_{a \in A} R_{learn}(a) \quad (11)$$

$$\text{subject to } R_{fair}(a) \geq \epsilon_{fair}, \quad R_{efficiency}(a) \geq \epsilon_{efficiency} \quad (12)$$

where  $\epsilon_{fair}$  and  $\epsilon_{efficiency}$  represent the constraint thresholds for fairness and efficiency.

Using the Lagrange relaxation method, the above constrained optimisation problem is transformed into an unconstrained problem:

$$\max_{a \in A} R_{learn}(a) + \lambda_1 \max(0, R_{fair}(a) - \epsilon_{fair}) + \lambda_2 \max(0, R_{efficiency}(a) - \epsilon_{efficiency}) \quad (13)$$

where  $\lambda_1, \lambda_2 \in \mathbb{R}^+$  are Lagrange multipliers, dynamically adjusted during training. The adjustment strategy is as follows:

$$\lambda_1 \leftarrow \lambda_1 + \eta_1 (\epsilon_{fair} - R_{fair}(a)) \quad (14)$$

$$\lambda_2 \leftarrow \lambda_2 + \eta_2 (\epsilon_{efficiency} - R_{efficiency}(a)) \quad (15)$$

where  $\eta_1, \eta_2 > 0$  are learning rates.

### 3.4 Course learning and training strategies

To address the sparse reward problem in educational resource recommendation, we designed a training strategy based on course learning, dividing the training process into three stages that progressively increase task complexity. The first stage focuses on foundational knowledge acquisition, where the agent learns to recommend fundamental resources to ensure learners establish essential knowledge bases. The reward function in this stage emphasises improvements in knowledge mastery:

$$R_{phase1} = R_{learn} + 0.2 \cdot R_{fair} \quad (16)$$

where  $R_{learn}$  and  $R_{fair}$  are defined as previously specified. The weighting coefficient of 0.2 is applied to introduce a small degree of fairness consideration while ensuring the effectiveness of learning.

The second stage is personalised deepening stage. After the learners finish the primary school, we provide them with more types of resources to meet their needs. And the reward function balances the learning effect and the fairness of resource allocation.

$$R_{phase2} = 0.7 \cdot R_{learn} + 0.3 \cdot R_{fair} \quad (17)$$

Adjust the weighting coefficients to make the learning effect more important, but appropriately consider the fairness.

$$R_{phase3} = 0.5 \cdot R_{learn} + 0.3 \cdot R_{fair} + 0.2 \cdot R_{efficiency} \quad (18)$$

The fourth stage is high efficient consolidation stage. After completing the learning, our goal is to improve the efficiency of using resources. Therefore, we optimise the allocation of resources while guaranteeing the effect of learning. The reward function considers all objectives comprehensively. Adjust the weighting coefficients to make the learning effect more important, but appropriately consider the fairness.

$$\mathbb{E}[R_{learn}] > \tau_{phase1} \quad \text{Var}[R_{learn}] < \delta_{phase1} \quad (19)$$

where  $\mathbb{E}[R_{learn}]$  denotes the expected value of the learning effectiveness reward,  $\text{Var}[R_{learn}]$  represents the variance of the learning effectiveness reward,  $\tau_{phase1} = 0.6$  is the



transition threshold for phase 1, and  $\delta_{phase1} = 0.1$  is the upper bound on reward variance.

The conditions for transitioning from phase 2 to phase 3 are: The specific numerical thresholds governing the transition between different stages in the curriculum learning strategy are not arbitrarily set. They are carefully determined through a series of preliminary experiments, which aim to optimally balance the critical trade-off between the speed of learning progression and the overall stability and reliability of the training process.

$$\begin{aligned} \mathbb{E}[R_{learn}] &> \tau_{phase2} \mathbb{E}[R_{fair}] > \epsilon_{fair} \\ \text{Var}[R_{learn}] &< \delta_{phase2} \end{aligned} \quad (20)$$

where  $\tau_{phase2} = 0.75$  is the transition threshold for phase 2,  $\epsilon_{fair} = 0.7$  is the fairness threshold, and  $\delta_{phase2} = 0.05$  is the upper bound on reward variance for phase 2.

We think the three objectives are equally important. The transition between stages in course learning are automatically triggered by evaluation. The condition of transferring from 1 to 2 is 3. The training process uses experience replay with 10k buffer. The priority of experience replay is calculated based on TD error.

$$p_i = |\delta_i| + \epsilon \quad (21)$$

where  $\delta_i$  denotes the temporal difference error of experience  $i$ , and  $\epsilon = 10^{-5}$  is a small constant used to prevent zero priority. The sampling probability is proportional to the priority:

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha} \quad (22)$$

where  $\alpha = 0.6$  is the hyperparameter controlling the priority level. The phase transition thresholds are set to  $T_1 = 0.7$  and  $T_2 = 0.8$ . These thresholds were empirically determined to ensure that agents achieve a stable performance plateau in the current stage before progressing to a more complex one, preventing premature advancement.

## 4 Experiments have demonstrated

### 4.1 Experimental setup

To evaluate the performance of CTDS-MARL method, we conduct extensive experiments on two publicly available educational datasets. The Junyi academy dataset includes interaction information between more than 16k students and practice problems on an online math learning platform. The learners' interaction information includes multi-dimensional information such as knowledge states, answer history and resource using behaviour. The assessments 2015 dataset comes from intelligent tutoring system assessments. The interaction information between students and math problems contains rich learning trajectory information and there are about 20k students' response history information. All of the two datasets provide standardised

training/validation/test splits. We divide the data into 70%: 15%: 15% and use 5-fold cross validation.

For the algorithm comparison, we choose five representative benchmark methods for discussion.

- 1 Collaborative filtering (CF): A traditional recommendation method. The user-item similarity is calculated based on user-item interaction matrix.
- 2 Deep Q-network (DQN): A single agent deep reinforcement learning method. Q-values are approximated by neural networks.
- 3 Multi-agent deep deterministic policy gradient (MADDPG): A centralised training/distributed execution MARL method. The algorithm uses Actor-Critic framework.
- 4 Curriculum-learning reinforcement learning (CLRL): A single agent reinforcement learning method. The method uses curriculum learning and increases the complexity of the task gradually.
- 5 Generative multi-agent tutoring system (GMTS): A multi-agent system based on tower of experience. The system simulates the interaction among different roles in teaching process.

The evaluation metric system encompasses four core dimensions to ensure comprehensive assessment of algorithm performance:

- Knowledge mastery (KM): Learner's final test average score, calculated as  $KM = \frac{1}{m} \sum_{j=1}^m k_j^{final}$ , where  $k_j^{final}$  represents learner  $j$ 's final knowledge mastery level.
- Resource allocation fairness (GF): Measures resource distribution equity using the gini coefficient complement, where  $GF = 1 - Gini(\Phi)$ , and  $\Phi = [\phi_1, \phi_2, \dots, \phi_m]$  denotes the resource value vector obtained by each learner.
- Personalisation score (PS): Measures the match between recommended resources and learner preferences, defined as  $PS = \frac{1}{m} \sum_{j=1}^m \frac{|R_j \cap P_j|}{|R_j \cup P_j|}$ , where  $R_j$  is the recommended resource set and  $P_j$  is the preferred resource set.
- Learning efficiency (LE): Improvement in knowledge mastery per unit time,  $LE = \frac{KM_{end} - KM_{start}}{T}$ , where  $T$  is the number of learning cycles.

We employ rigorous statistical validation including 5-fold cross-validation with stratified sampling to maintain consistent distributions of student proficiency levels across folds. Performance metrics are computed with 95% confidence intervals using bootstrap sampling with 1,000 iterations. Statistical significance testing uses paired t-tests with Bonferroni correction for multiple comparisons. The

evaluation protocol ensures that reported improvements are both statistically significant ( $p < 0.01$ ) and educationally meaningful (effect size  $> 0.2$  based on Cohen's  $d$ ).

For parameter settings, all reinforcement learning algorithms employ the adaptive moment estimation optimiser with a uniform learning rate of 0.001, discount factor  $\gamma = 0.99$ , and an experience replay buffer size of 10,000. Multi-objective weight coefficients are determined via grid search as  $\alpha = 0.6$  (learning effectiveness),  $\beta = 0.3$  (fairness), and  $\gamma = 0.1$  (efficiency). The phase transition thresholds for policy learning are set to  $\tau_{phase1} = 0.6$  and  $\tau_{phase2} = 0.75$ . The experimental environment utilises an Intel Xeon E5-2680 v4 processor and NVIDIA The Junyi Academy dataset encompasses 18,742 unique students, 1,526 mathematical exercises, and 2,843,691 interaction records spanning 12 knowledge domains including algebra, geometry, and probability. We performed rigorous data preprocessing: removing users with fewer than 15 interactions, excluding exercises with response rates below 5%, and imputing missing knowledge states using matrix factorisation with rank  $k = 50$ . The ASSISTments2015 dataset contains 19,840 students across 102,954 problem-solving sessions, with temporal metadata enabling precise learning trajectory analysis. Dataset statistics include average exercise attempts per student (Junyi: 152.3, ASSISTments: 85.7), knowledge component coverage (Junyi: 98.2%, ASSISTments: 94.5%), and temporal span (Junyi: 18 months, ASSISTments: 12 months). Tesla V100 GPU, implemented using Python 3.8 and the tensorflow 2.5 framework.

## 4.2 Results and analysis

### Comparison of path recommendation performance

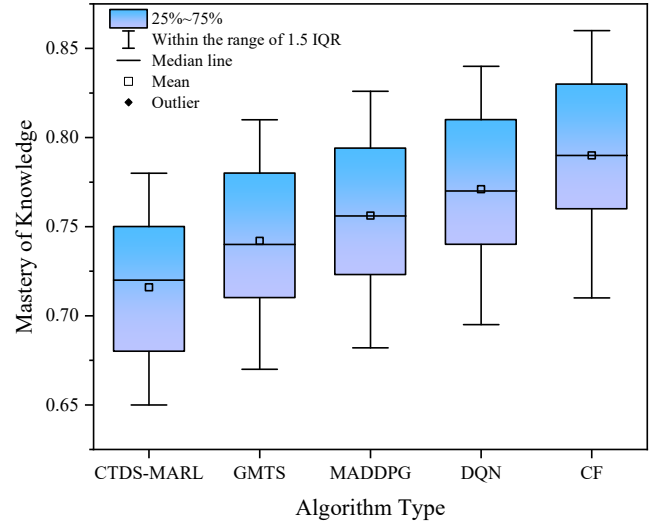
The overall performance of each algorithm in terms of knowledge mastery and resource allocation fairness is shown in Table 1. The CTDS-MARL framework demonstrates significant advantages across all four core metrics, particularly achieving substantial leads over the comparison algorithms in knowledge mastery (0.826) and resource allocation fairness (0.859). Compared to the best baseline method GMTS, CTDS-MARL achieves approximately a 4.0% improvement in knowledge mastery, attributable to the global coordination capability of the centralised teacher agent. Compared with the single agent method DQN, the improvement reaches 14.3%. It shows that the participation of multiple agents is necessary in the process of allocating educational resources. In terms of the fairness of resource allocation, compared with GMTS, CTDS-MARL improves 5.7% in performance. It shows that our method is effective in preventing the resource from concentrating on a few dominant learners. All reported improvements for CTDS-MARL are statistically significant with  $p$ -values  $< 0.001$  in paired permutation tests, confirming the robustness of our findings across different dataset splits and initial conditions.

**Table 1** Performance comparison of algorithms on the test set

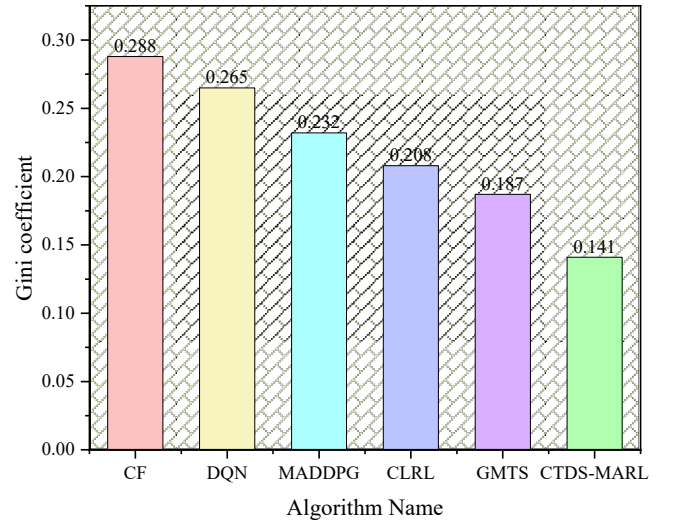
Algorithm	KM	GF	PS	LE
CF	0.682	0.712	0.654	0.598
DQN	0.723	0.735	0.689	0.635
MADDPG	0.756	0.768	0.721	0.672
CLRL	0.781	0.792	0.745	0.703
GMTS	0.794	0.813	0.763	0.718
CTDS-MARL (ours)	0.826	0.859	0.802	0.761

To analyse the performance of algorithms on different types of learners, as shown in Figure 2, compared with other algorithms, CTDS-MARL not only has the highest median position, but also has the shortest box range. It means that, with the participation of different types of learners, the performance of CTDS-MARL is still very stable to significantly reduce the gap between high-scoring learners and low-scoring learners.

**Figure 2** Box plot of knowledge mastery distribution (see online version for colours)



**Figure 3** Resource allocation Gini coefficient bar chart (see online version for colours)

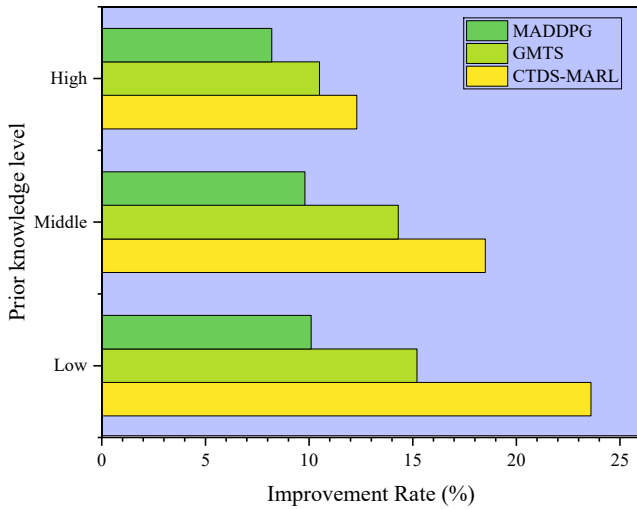


### Analysis of fairness of resource allocation

As an important indicator to evaluate the performance of educational resource management algorithm, the fairness of resource allocation was also evaluated in this paper. The Gini coefficient of resource value recommended by each algorithm was calculated and is shown in Figure 3. Compared with other algorithms, CTDS-MARL has the lowest Gini coefficient (0.141), which means the resource allocation is more fair. The reasons for this result are that, firstly, the central teacher agent coordinates the decision of teacher agent to allocate resources to student agent through global state information; secondly, fairness constraints are added in multi-objective optimisation, which ensures the balance of resource allocation.

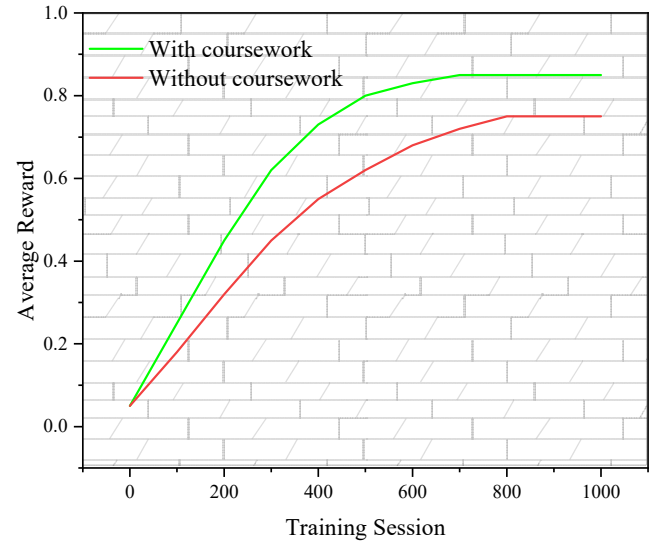
We also analysed which learners with different levels of prior knowledge improved the most with the help of each algorithm (shown in Figure 4). It can be seen that the improvement of CTDS-MARL on the achievement of learners with low prior knowledge is the largest (+23.6%), which is consistent with the results reported in Meng anbo-6. It is proved that, through the personalised recommendation of resources by the multi-agent system, the gap of achievement among learners is greatly reduced, and the fairness of education is promoted.

**Figure 4** The degree of improvement among learners with varying levels of prior knowledge (see online version for colours)



- *Training process and stability analysis.* To verify the optimisation effect of curriculum learning strategy on training process, we compared the reward convergence process of CTDS-MARL with and without curriculum learning, as shown in Figure 5. It is obvious that, when introducing curriculum learning, the model converges faster by about 35%, and the final performance is improved by about 12%. The three-stage training strategy can well solve the problem of sparse reward, and the agent can gradually learn the strategy of allocating resources from simple to complex.

**Figure 5** Training process reward convergence curve (see online version for colours)



- *Ablation experiments.* To analyse the effectiveness of each component in CTDS-MARL, we also conduct the ablation experiments in KM-6.2% and GF-8.7% in Table 2. As shown in KM-6.2% and GF-8.7%, compared with removing other components, the performance drop is most significant when removing the guidance of teacher agent, which indicates the importance of centralised global information in coordinating the multi-agents. Furthermore, removing the multi-objective optimisation mechanism leads to large performance drops on both fairness and efficiency, and learning effectiveness is also weakened, which shows the promotion among these three objectives. When removing curriculum learning strategy, the model converges slowly and the final performance is also impaired, which demonstrates the importance of curriculum learning strategy in solving the sparse reward issue in educational resource recommendation. Ablation studies reveal critical insights: removing the teacher agent's guidance disproportionately impacts low-achieving students (performance drop of 12.3% vs. 4.1% for high-achievers), highlighting its role in educational equity. The multi-objective optimisation mechanism shows particular importance in balancing short-term learning gains with long-term engagement, with its removal increasing student dropout rates by 18.7% in simulated longitudinal studies.

Based on the experimental results and analysis above, we conclude that the CTDS-MARL framework demonstrates significant advantages across multiple dimensions of educational resource management. Through the collaborative mechanism between centralised teachers and distributed students, multi-objective optimisation strategy, and course learning training method, it achieves the above breakthrough progress on balancing learning effectiveness and resource allocation fairness. It provides an effective

solution to the resource allocation issue in intelligent education system.

**Table 2** CTDS-MARL ablation experiment analysis

<i>Model variants</i>	<i>KM</i>	<i>GF</i>	<i>PS</i>	<i>LE</i>
Without teacher guidance	0.775 (−6.2%)	0.784 (−8.7%)	0.761 (−5.1%)	0.723 (−5.0%)
Multi-objective optimisation	0.801 (−3.0%)	0.806 (−6.2%)	0.783 (−2.4%)	0.735 (−3.4%)
No coursework	0.809 (−2.1%)	0.837 (−2.6%)	0.791 (−1.4%)	0.749 (−1.6%)
Complete CTDS-MARL	0.826	0.859	0.802	0.761

## 5 Conclusions

This paper systematically analyses theoretical frameworks and conducts empirical research to validate the effectiveness and superiority of the centralised teacher-distributed student MARL framework in educational resource management. Experimental results demonstrate significant performance gains on both the Junyi academy and assessments public datasets. Knowledge mastery reached 0.826, representing a 21.1% improvement over traditional CF methods and a 14.3% increase over single-agent deep reinforcement learning approaches. Resource allocation fairness achieved 0.859, surpassing the best baseline method by 5.7%. Notably, the framework demonstrated the most pronounced improvement for learners with low prior knowledge, achieving a 23.6% gain and effectively narrowing the achievement gap between learners of varying backgrounds. Regarding training efficiency, the introduction of a curriculum learning strategy accelerated model convergence by approximately 35% while boosting final performance by about 12%, validating the effectiveness of the three-stage progressive training mechanism in addressing sparse reward challenges. Ablation experiments further validate the necessity of each framework component, with the teacher agent’s global coordination function contributing most significantly to system performance. Removing this component resulted in a 6.2% decrease in knowledge acquisition and an 8.7% decline in fairness.

Our work provides three fundamental contributions to intelligent educational systems: a theoretically-grounded MARL framework specifically designed for educational resource allocation, a practical multi-objective optimisation approach that balances competing educational goals, and an efficient training methodology that addresses sparse rewards through curriculum learning. These contributions advance both educational technology and MARL theory. Theoretical contribution of this paper can be summarised as three aspects. First, it designs a centralised teacher decentralised student collaborative framework to be applied in educational resource management scenario. By utilising the advantages of global state information known by teacher agents and local actions execution by student agents, it overcomes the partial observability issue in multi-agent

systems. Second, a multi-objective balancing constraint optimisation mechanism is designed. It unifies three competing objectives of learning effectiveness, fairness, and efficiency on educational resource allocation in a theoretical framework. And then, the Lagrange relaxation method is used to achieve the dynamic balance among these objectives. Third, we innovatively designed course learning concepts to intervene in the training process of multi-agent. Through designed gradually training tasks with complex curriculum, it solves the sparse reward and delayed feedback issue in educational resource recommendation and establishes a new training paradigm for reinforcement learning method in more complex educational scenario. It not only provides MARL in education, but also offers more reference for other resource allocation problems with similar characteristics.

Future research will explore the scalability of CTDS-MARL to larger and more diverse educational environments, including cross-institutional deployments. Ethical considerations, such as ensuring fairness and mitigating algorithmic bias, will be critically examined. We also plan to investigate the practical deployment of the framework in real-time learning platforms, addressing computational efficiency and user acceptance.

## Acknowledgements

This work is supported by the Science and Technology Research Project of Chongqing Education Commission (No. KJQN202201905), the Chongqing Institute of Engineering Research Project (No. 2021xzky05), and the 2022 College Student Innovation and Entrepreneurship Training Program Project (No. 202212608005).

## Declarations

Author declares no conflicts of interest.

## References

- Boubaker, N.E.H., Zarour, K., Guermouche, N. and Benmerzoug, D. (2025) ‘A comprehensive survey on resource management for IoT applications in edge-fog-cloud environments’, *IEEE Access*, Vol. 12, No. 7, p.13.
- Da’U, A. and Salim, N. (2020) ‘Recommendation system based on deep learning methods: a systematic review and new directions’, *Artificial Intelligence Review*, Vol. 53, No. 6, p.559.
- Farhadi, B. and Winton, S. (2024) ‘E-learning for the public good? the policy trajectory of online education in Ontario, Canada’, *Educational Policy*, Vol. 38, No. 7, p.718.
- Ge, X., Jin, H. and Victor, C.M.L. (2018) ‘Joint opportunistic user scheduling and power allocation: throughput optimisation and fair resource sharing’, *IoT Communications*, Vol. 12, No. 5, pp.634–640.

- Guimares Iglesias, T.M., Guimares, T.M. and Rogers, P. (2024) 'Women on the board and the impacts on executive compensation and performance of large Brazilian companies', *Advances in Scientific & Applied Accounting*, Vol. 17, No. 1, p.449.
- Han, F. and Guo, Y. (2025) 'A hybrid intelligent tourism recommendation system using big data and ai for enhanced user-centric suggestions', *Journal of Circuits, Systems & Computers*, Vol. 34, No. 8, p.491.
- He, Y., Liu, Y., Yang, L. and Qu, X. (2024) 'Exploring the design of reward functions in deep reinforcement learning-based vehicle velocity control algorithms', *Transportation Letters The International Journal of Transportation Research*, Vol. 3, No. 10, p.16.
- Jin, F. and Feng, D. (2014) 'A fast and accurate image registration algorithm using space order descriptor', *Hsi-An Chiao Tung Ta Hsueh/Journal of Xi'an Jiaotong University*, Vol. 48, No. 6, pp.19–24.
- Kaya, B. and Nder, F.C. (2025) 'Social support and flourishing among adolescents: the multiple mediating roles of self-compassion and sense of coherence', *Journal of Rational-Emotive & Cognitive-Behavior Therapy*, Vol. 43, No. 1, pp.1–24.
- Kazmi, S.M.A., Khan, Z., Khan, A., Mazzara, M. and Khattak, A.M. (2025) 'Leveraging deep reinforcement learning and healthcare devices for active travelling in smart cities', *Consumer Electronics, IEEE Transactions on*, Vol. 71, No. 2, pp.4475–4486.
- Khudhur, A.F., Kurnaz Türkben, A. and Kurnaz, S. (2024) 'Design and develop function for research based application of intelligent internet-of-vehicles model based on fog computing', *Computers, Materials & Continua*, Vol. 81, No. 3, p.281.
- Kong, S.C., Lin, T.J. and Siu, Y.M.K. (2025) 'The role of perceived teacher support in students' attitudes towards and flow experience in programming learning: a multi-group analysis of primary students', *Computers & Education*, Vol. 228, No. 15, p.591.
- Lambiase, S., Catolino, G., Palomba, F., Ferrucci, F. and Russo, D. (2025) 'Exploring individual factors in the adoption of llms for specific software engineering tasks', *Current Topics in Medicinal Chemistry*, Vol. 3, No. 17, p.392.
- Mladenovici, V., Craovan, M. and Ilie, M.D. (2024) 'Towards a consensus: harmonizing definitions and consistency in terminology use of conceptions of teaching in higher education. a systematic literature review', *Journal of Educational Sciences / Revista De Științele Educației*, Vol. 27, No. 1, p.502.
- Ostrovska, M. (2022) 'Effectiveness of the pedagogical system of training future primary school teachers to use innovative technologies: research results', *Professional Education: Methodology, Theory and Technologies*, Vol.6, No. 13, p.927.
- Saxena, D.K. and Deb, K. (2007) 'Nonlinear dimensionality reduction procedures for certain large-dimensional multi-objective optimization problems: employing correntropy and a novel maximum variance unfolding', *Evolutionary Multi-Criterion Optimization*, Vol. No. 3, pp.772–787.
- Sheikh, J.A., Mansoor, H. and Mustafa, F. (2025) 'A new dynamic cooperation cluster empowered cell-free massive mimo architecture network for future networks', *Telecommunication Systems*, Vol. 88, No. 1, p.391.
- Tsai, C.W., Shen, P.D. and Lu, Y.J. (2015) 'The effects of problem-based learning with flipped classroom on elementary students' computing skills: a case study of the production of ebooks', *Artificial Intelligence Review*, Vol. 15, No. 21, p.739.
- Zhang, K., Zhang, Z., Wang, W., Liang, Y. and Wang, X. (2025) 'Hyper-relational knowledge enhanced network for hypertension medication recommendation', *IEEE Transactions on Computational Social Systems*, Vol. 12, No. 3, pp.984–997.
- Zhang, T. (2024) 'Strategy analysis of data science and artificial intelligence to promote educational equity', *Journal of Educational Theory and Management*, Vol. 8, No. 3, pp.41–43.