



Design of Reinforcement Learning-Empowered Intelligent Evaluation System for Translation Teaching and Research on Teachers 'Digital Literacy Development Path

Weijia Liu^{1,*}, Ruirui Zhang¹ and Asel Musurapshaevna Toksonalieva²

¹ School of International Studies, Luoyang Institute of Science and Technology, Luoyang 471023, Henan, China

² Kyrgyz-Chinese Institute, Kyrgyz National University named after Jusup Balasagyn, Bishkek 720033, Kyrgyz Republic

SUMMARY: *The evaluation of translation teaching involves multi-dimensional abilities such as semantic understanding, terminology control, syntactic transformation, discourse cohesion and cultural adaptation. Traditional manual scoring, automatic indicator evaluation and large language model annotation methods are often difficult to simultaneously give attention to diagnostic fineness, feedback continuity and teaching decision support. In order to solve this problem, this paper designs an intelligent evaluation system for translation teaching empowered by reinforcement learning (RL-IETS). RL-IETS converts student translations, revision records, error labels, teacher annotations, and platform behavior logs into learning state representations, and dynamically selects scoring, error diagnosis, feedback generation, learning path recommendation, and teacher intervention strategies through the state-action-reward mechanism. The experimental results show that the average evaluation accuracy of the system reaches 91.8%, the comprehensive value of error diagnosis reaches 89.4%, the profit of personalized path optimization reaches 83.6%, and the comprehensive score of teachers 'digital literacy is improved from 67.8 to 86.5. The research shows that the system can improve the adaptive level of translation teaching evaluation, and provide a traceable practice path for teachers 'digital literacy development.*

KEYWORDS: *Reinforcement learning; Translation teaching; Intelligent evaluation; Teachers' Digital Literacy*

1 Introduction

Translation teaching is being influenced by machine translation, generative artificial intelligence and learning analytics. The student translation is no longer just the final text of the class exercise, but the learning evidence containing the trajectory of term selection, semantic transformation, discourse cohesion, cultural adaptation and post-translation revision. In this context, translation teaching evaluation needs to shift from single manual evaluation to continuous, fine-grained and feedback-able intelligent evaluation. Reinforcement learning can continuously adjust the evaluation strategy according to students 'translation performance, feedback acceptance and subsequent modification results, which provides a new calculation method for translation ability diagnosis and learning path planning. At the same time, whether teachers can understand, use and improve the intelligent evaluation system has also become a

*wei_jia_liu@163.com

<https://doi.org/10.65102/is2026983>

key issue in the digital transformation of translation teaching.

1.1 Challenges faced by intelligent evaluation of translation teaching

Translation ability is obviously complex. There are interactive relationships among lexical accuracy, syntactic transformation, discourse coherence, style consistency, cultural adaptation and the use of professional terms, and a single indicator is difficult to completely describe the quality of translation. The existing classroom evaluation usually relies on teachers' correction experience, and the scoring results are subjective. The feedback is mostly focused on the explicit error level, and the attention is insufficient to the process of students' translation modification, the law of error transfer and the trajectory of students' ability growth.

Although intelligent evaluation can improve the efficiency of grading, it still faces three types of difficulties. First, machine scoring is easy to compress the quality of translation into a total score, ignoring the corresponding relationship between error types and teaching objectives. Second, the system feedback is often based on template suggestions, which is difficult to adjust the depth of explanation according to the current ability level of students. Third, when teachers use the automatic evaluation results, if they lack the ability to interpret data and understand algorithms, they are easy to equate the system output with teaching judgment directly, which weakens teachers' professional decision-making role.

1.2 Existing translation teaching evaluation and teachers' digital literacy training models

The existing translation teaching evaluation mainly includes manual grading, peer evaluation, machine translation quality indicator evaluation and large language model assisted annotation. Manual scoring has the advantage of context judgment, but it is inefficient when dealing with large-scale translation. BLEU, COMET and other automatic indicators can quantify the similarity and semantic quality of translation, but it is difficult to directly explain why students make mistakes. Large language models perform well in translation quality assessment and error segment localization, but their output stability, classroom controllability and teachers' interpretable use still need to be further verified.

The cultivation of teachers' digital literacy is mainly based on tool training, platform operation and resource construction, and rarely forms a closed loop around "intelligent evaluation - data interpretation - teaching intervention - reflection and improvement". Translation teachers not only need to be able to use the system, but also need to understand translation quality metrics, algorithmic feedback boundaries, learning data ethics, and ways to personalize teaching decisions. If teachers' digital literacy construction is separated from real evaluation tasks, it is difficult to transform technical ability into classroom teaching improvement ability.

1.3 Research contribution of this paper

This paper constructs a reinforcement learning-based intelligent evaluation system for translation teaching, which integrates student translations, modification records, error types, teacher feedback and learning results into a unified computing framework.

This paper proposes a translation evaluation strategy based on state-action-reward mechanism, which takes students' learning status, error distribution and feedback response as state input, so that the system can dynamically select scoring, diagnosis, prompt and path recommendation actions.

A translation error recognition and intelligent feedback generation module is designed, which combines the features of semantic similarity, term consistency, discourse coherence and

cultural adaptation to classify and diagnose translation errors, and generate feedback texts for learning improvement.

To construct a support mechanism for teachers' digital literacy development, and to transform teachers' interpretation, intervention, correction and reflection behaviors on system results into monitable indicators, so as to provide data basis and path suggestions for the improvement of translation teachers' digital ability.

1.4 Definition of research problem

Although artificial intelligence has been widely introduced into language teaching and translation training, there is still a gap between system design and teaching application of intelligent evaluation for translation teaching. Some systems emphasize automatic translation grading, but do not fully deal with individual differences of students. Some studies focus on teachers' digital competence, but lack of evaluation scenarios combined with specific translation teaching tasks. This paper focuses on three issues: how reinforcement learning supports the dynamic optimization of translation teaching evaluation strategies; How to realize the collaborative operation of translation error diagnosis, feedback generation and learning path recommendation in intelligent system; How teacher digital literacy develops during system use, data interpretation, and instructional interventions. Table 1 shows the overall arrangement of the research content in this paper.

Table 1: Summary of the research in this paper

Section	Content
Introduction	Clarifies the background, challenges, research contributions, and problem scope of intelligent evaluation in translation teaching
Related Research	Reviews studies on machine translation-assisted teaching, translation quality evaluation, reinforcement learning applications in education, and teachers' digital literacy
System Design Method	Constructs modules for data collection, state modeling, reinforcement learning-based evaluation, error diagnosis, feedback generation, and teacher monitoring
Evaluation Indicator System	Establishes evaluation indicators for translation quality, feedback efficiency, adaptive capability, learning engagement, and teachers' digital literacy
Results and Discussion	Analyzes system evaluation accuracy, learning path optimization effects, error diagnosis performance, and teachers' digital literacy improvement results
Conclusion	Summarizes the system value, teaching implications, and future improvement directions

2 Related Research

The research on intelligent evaluation of translation teaching mainly involves machine translation assisted teaching, automatic evaluation of translation quality, generative artificial intelligence feedback, reinforcement learning learning path recommendation, and teachers' digital literacy training. Previous studies have shown that machine translation can provide translation reference and modification resources for language classes, but its teaching value depends on whether teachers can transform machine output into analysis materials, rather than directly replacing students' translation process [1-5]. Neural machine translation and large

language models have improved the quality of translation generation, and also pushed the automatic evaluation method from surface similarity calculation to semantic consistency, error segment localization, and multi-dimensional quality judgment [9-16]. However, the evaluation goal in the translation classroom is not only to judge the quality of the translation, but also to identify the differences in students' ability in understanding, transformation, expression, revision and cultural adaptation. The existing research has a good foundation in translation quality analysis, but the continuous modeling of learning state, personalized feedback generation and teacher intervention behavior recording are still insufficient. Table 2 summarizes the main characteristics of MT assisted teaching, translation quality evaluation, educational application of reinforcement learning, and teacher digital literacy research.

Table 2: Summary of related studies

Reference	Research Environment	Algorithm / Method	Translation Quality Analysis	Learning Path Adaptation	Teachers' Digital Literacy	Main Results
[1]	Foreign language education	MT-assisted teaching	✓	×	×	Supports language learning; lacks evaluation loop
[3]	Foreign language teaching	NMT review	✓	×	×	Reviews NMT use; lacks personalized modeling
[4]	Technical translation teaching	Evaluation framework	✓	×	✓	Notes evaluation difficulty; weak adaptability
[7]	EFL writing detection	AI text detection	✓	×	✓	Focuses on AI-text recognition; limited feedback support
[8]	Teacher training	AI competency training	×	×	✓	Improves AI readiness; not linked to translation assessment
[10]	Translation quality evaluation	COMET metric	✓	×	×	Enhances evaluation accuracy; weak teaching explanation
[11]	Machine translation evaluation	LLM evaluation	✓	×	×	Shows strong translation evaluation ability

As can be seen from Table 2, the existing research has formed a relatively clear technical foundation in translation quality evaluation, intelligent feedback, reinforcement learning recommendation and teacher digital ability training, respectively. COMET, GEMBA-MQM and large language model evaluation methods enhance the semantic sensitivity of translation quality judgment, and can identify the problems such as meaning deviation, term misuse and unnatural expression that are difficult to present in traditional similarity indicators. The reinforcement learning related research provides strategy optimization ideas for the dynamic adjustment of learning paths, which can continuously modify the system actions according to students' performance and feedback effects. Teachers' digital literacy research emphasizes tool understanding, data interpretation and AI ethical awareness, which provides a competency

framework for translation teachers to participate in intelligent evaluation.

There are still three shortcomings in the existing research. First, the quality evaluation of translation mostly stays at the level of result judgment, and does not pay enough attention to the change of students' revision process and learning state. Secondly, reinforcement learning is mostly used for general learning resource recommendation in educational scenarios, and it is less embedded in translation error diagnosis and feedback generation processes. Thirdly, teachers' digital literacy training is often separated from the specific classroom system, which is difficult to reflect teachers' real use ability in intelligent evaluation. Based on this, this paper integrates the reinforcement learning strategy, translation error recognition, intelligent feedback generation and teacher digital literacy monitoring into the same system framework to support the transformation of translation teaching evaluation from static scoring to dynamic diagnosis and continuous improvement.

3 The design method of reinforcement learning empowered intelligent evaluation system for translation teaching

The intelligent evaluation system of translation teaching constructed in this paper consists of six modules: data access, learning state modeling, translation quality diagnosis, reinforcement learning strategy selection, feedback generation and teacher development support. The system takes the student's source text understanding record, the translation submission text, the post-translation modification track, the teacher's annotation, the classroom task completion degree and the platform interaction log as input, and transforms the translation learning process into a computable state sequence. The evaluation model no longer stops at the total score judgment of the translation, but forms a multi-dimensional feature representation around semantic accuracy, term consistency, syntactic transformation, discourse cohesion, cultural adaptation and revision effectiveness. The reinforcement learning module selects actions such as scoring, prompting, questioning, example recommendation or teacher intervention according to the current state of the students, and updates the strategy according to the subsequent improvement of the translation and the change of learning participation, thus forming a dynamic evaluation closed loop for translation teaching.

3.1 Data collection and learning state modeling for translation teaching

As shown in Figure 1, translation teaching data are collected by the classroom platform, translation evaluation tool and teacher feedback end. The system first segmented the source text and the student translation, extracted terms, semantic coding and error labeling, and then merged the translation quality features with the learning behavior features. Translation quality features include semantic deviation, term misuse, missing translation, addition translation, improper collocation and insufficient discourse coherence. Learning behavior characteristics included the number of submissions, modification interval, feedback reading duration, prompt adoption rate, and teacher annotation response. After feature cleaning and normalization, the system forms the student's learning state vector at time step t , which provides input for the subsequent reinforcement learning evaluation strategy.

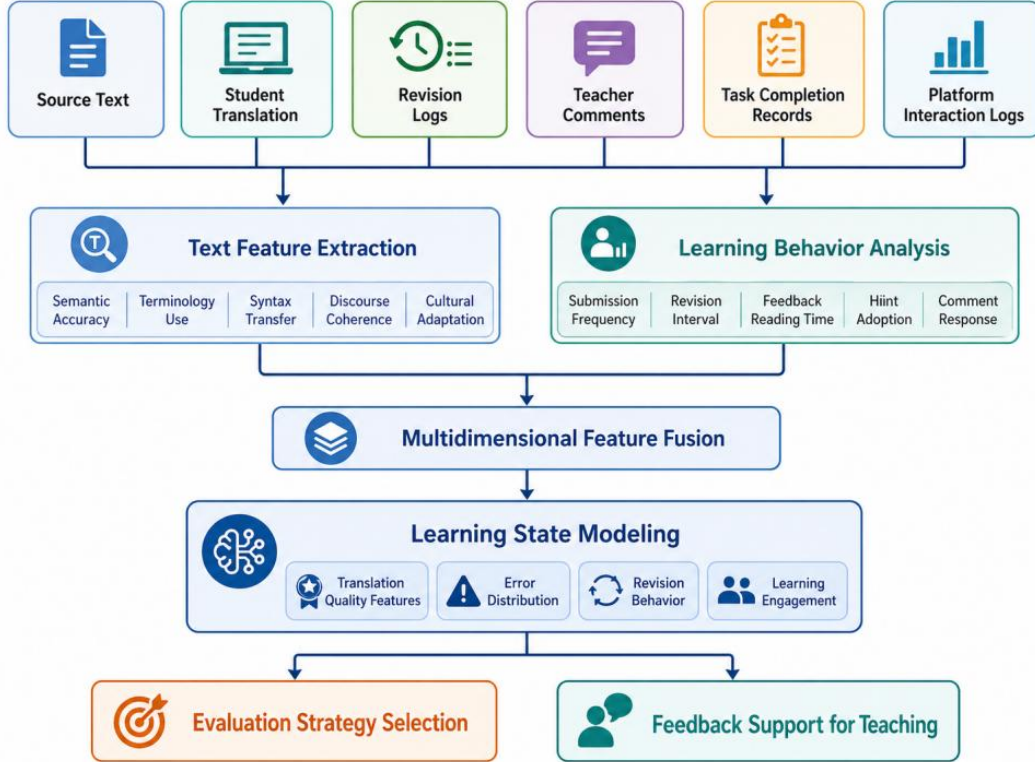


Figure 1: Process of data collection and learning state modeling in translation teaching

The core of learning state modeling lies in compressing scattered textual and behavioral evidence into a stable state representation. In this paper, let the learning state of student i at time step t be as follows.

$$S_{i,t} = \lambda_1 Q_{i,t} + \lambda_2 E_{i,t} + \lambda_3 R_{i,t} + \lambda_4 B_{i,t} \quad (1)$$

where, $Q_{i,t}$ represent the quality features of the translation, which mainly reflect the performance at the semantic, terminological, syntactic and discourse levels. $E_{i,t}$ represents the error distribution characteristics, which are used to describe the frequency and severity of different error types. $R_{i,t}$ represent the revision behavior characteristics, reflecting the revision quality of students after receiving feedback. $B_{i,t}$ represent the learning behavior characteristics of the platform, including task input and feedback response. The trainable weights λ_1 to λ_4 are jointly updated by the evaluation error on the validation set and the teacher calibration results.

The state vector is able to distinguish students with "similar translation scores but different error structures". For example, when two students both get similar ratings, one student may have mainly term inconsistencies, while the other student may be concentrated on textual logical breaks. After identifying the differences through the state vector, the system can push terminology base training and domain example sentences for the former, and provide discourse reorganization exercises and teacher explanatory feedback for the latter. Therefore, data collection is no longer a simple record of learning results, but serves for subsequent strategy selection, intelligent feedback generation and teacher teaching intervention.

3.2 Construction of Reinforcement learning evaluation Strategy Based on State-action-reward mechanism

Figure 2 shows that this module takes the learning state vector formed in 3.1 as input and transforms the evaluation process of translation teaching into a continuous decision process among states, actions and rewards. Students' translation performance at a certain time step is not fixed, and the system needs to choose different evaluation actions according to their semantic deviation, term error, syntactic transformation ability, modification quality and feedback response. The action set includes translation quality score, error type annotation, local prompt generation, sample translation recommendation, special exercise push, and teacher intervention reminder. The function of reinforcement learning strategy is not to replace the teacher's judgment, but to organize the multiple translation submission and feedback results into an upgradable strategy chain, so that the system can choose a more appropriate evaluation method in different learning states.

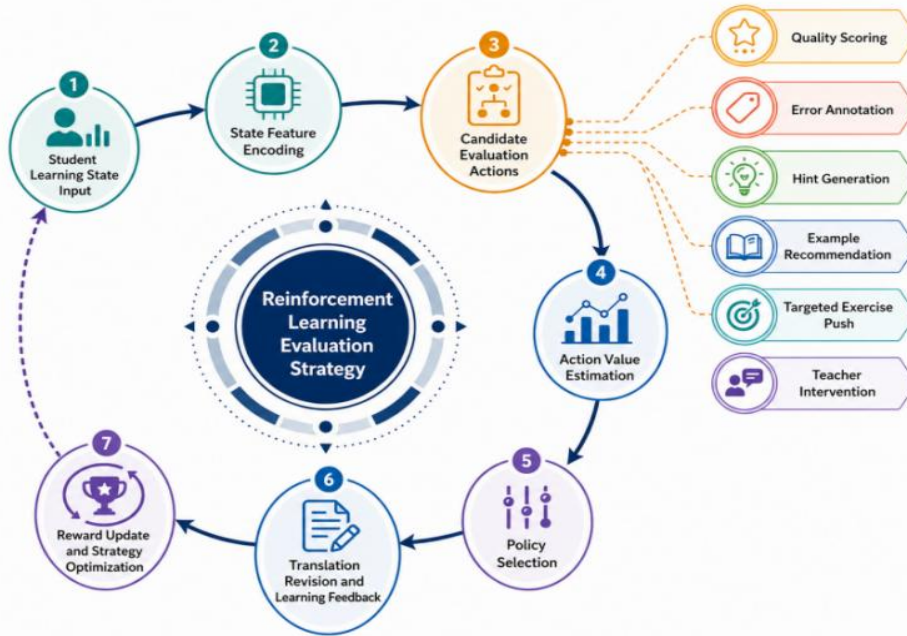


Figure 2: Construction process of reinforcement learning evaluation strategy

In this paper, the state of student i at time step t is denoted as $S_{i,t}$, the evaluation action is denoted as $a_{i,t}$, and the system receives a reward $r_{i,t}$ after performing the action. The action-value function is defined as follows.

$$Q(S_{i,t}, a_{i,t}) = r_{i,t} + \gamma \max_{a'} Q(S_{i,t+1}, a') \quad (2)$$

where, $Q(S_{i,t}, a_{i,t})$ represents the long-term value of taking a certain evaluation action in the current state; $r_{i,t}$ denotes the immediate reward after the action is performed; γ is the discount coefficient, which is used to balance the current feedback effect and the subsequent learning benefit. $S_{i,t+1}$ represents the new state of the student after completing the modification, reading the feedback, or receiving the teacher's intervention; a' denotes the selectable evaluation action in the next stage. The reward function consists of translation quality improvement, error reduction, feedback adoption, and learning engagement:

$$r_{i,t} = \beta_1 \Delta q_{i,t} + \beta_2 \Delta e_{i,t} + \beta_3 f_{i,t} + \beta_4 p_{i,t} \quad (3)$$

where, $\Delta q_{i,t}$ represents the improvement of translation quality score; $\Delta e_{i,t}$ represents the degree of error rate reduction; $f_{i,t}$ indicates the adoption of system feedback by students; $p_{i,t}$ represents the change of learning engagement; The values β_1 to β_4 are the reward weights. If the student significantly reduces the term misuse rate after receiving the term prompt, the action will receive a higher reward. If the system repeatedly pushes the same prompt and the students' modification effect is not obvious, the action value will decrease, and the subsequent strategy will turn to example explanation or teacher intervention.

This mechanism can avoid the intelligent evaluation from staying at the single score level. The system records the change of students' ability through state transition, so that the scoring, diagnosis, feedback and intervention form a closed loop. For students with weak semantic understanding, the strategy is more inclined to provide source text parsing and meaning reconstruction tips. For students with strong expressive ability but insufficient term stability, the strategy will give priority to calling term base and domain corpus example sentences. Therefore, the reinforcement learning evaluation strategy can transform the individual differences in translation teaching into a dynamic decision-making basis, and provide computational support for subsequent error recognition, intelligent feedback generation and personalized path planning.

3.3 Error recognition and intelligent feedback generation of student translation

Figure 3 shows that this module aligns student translations with the source semantic representation, course standard translation examples and teacher annotation rules, and identifies semantic deviation, term misuse, missing translation, addition translation, syntactic conversion imbalance, insufficient discourse cohesion and cultural adaptation deviation in the translation. The system first processed the source text and the translation text by clause processing and semantic vector coding, and then constructed the error discrimination basis through the glossary, the domain corpus and the teacher's historical annotation records. The error detector not only marks the location of the error, but also maps the error to the translation competence dimension, allowing the system to determine whether the student has insufficient understanding, unstable expression, or deviation in revision strategy. The feedback generator then outputs explanatory suggestions, partial rewriting hints, and specialized training tasks based on error levels, student status, and teacher intervention needs.

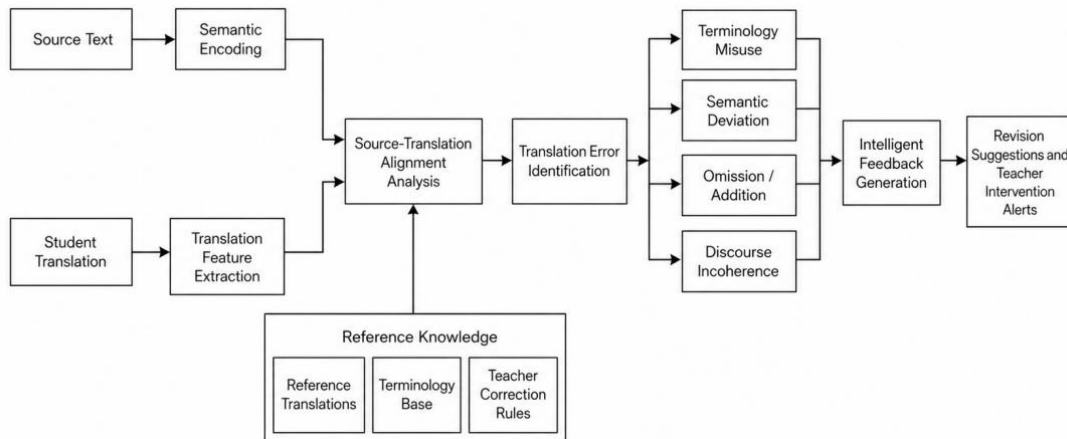


Figure 3: Student translation error recognition and intelligent feedback generation module

In this paper, the recognition score of type k translation error is expressed as follows.

$$D_{i,t}^k = \mu_1(1 - \cos H_{s,t}, H_{r,t}) + \mu_2 T_{i,t}^k + \mu_3 G_{i,t}^k + \mu_4 C_{i,t}^k \quad (4)$$

where, $D_{i,t}^k$ represent the type k error intensity of student i at time step t; $H_{s,t}$ and $H_{r,t}$ represent the source text semantic vector and the student translation semantic vector respectively. $T_{i,t}^k$ denote the term matching bias; $G_{i,t}^k$ denote the syntactic transformation bias; $C_{i,t}^k$ represent the adaptation deviation between culture and discourse; μ_1 to μ_4 are the false recognition weights. If the semantic similarity is low and the term deviation is obvious, the system will mark the segment as a high-risk error and enter the feedback generation step.

Feedback action selection needs to balance error severity, explanation cost, and learning acceptability. In this paper, let the set of feedback actions be A_f and the optimal feedback action be defined as follows.

$$a_f^* = \arg \max_{a \in A_f} [\rho_1 D_{i,t}^k + \rho_2 U_{i,t}(a) - \rho_3 L_{i,t}(a)] \quad (5)$$

where, a_f^* represents the feedback action selected by the system; $U_{i,t}(a)$ represents the expected help of the feedback to the student's translation revision; $L_{i,t}(a)$ denotes the feedback comprehension load; The weights ρ_1 to ρ_3 are chosen for the actions. For minor collocation errors, the system generates short replacement suggestions. For the semantic reconstruction errors, the system provides the logical explanation of the source text, the key component hints, and the teacher review marks. This design makes error identification consistent with feedback generation, and avoids feedback content from specific translation problems.

3.4 Personalized translation learning path and teaching intervention strategy planning

Figure 4 shows that this module generates personalized translation learning paths based on student learning status, translation error distribution, and classroom task objectives. Translation competence is not a simple superposition of isolated skills, and there are obvious relationships between source text understanding, term control, syntactic transformation, discourse reconstruction, cultural adaptation and post-translation revision. If students have not accurately understood the source text logic, the system should not directly push high-level style polishing tasks. If the stability of students' terminology is insufficient, priority should be given to term base training, parallel text comparison and short sentence translation exercises. Based on the state results output by the reinforcement learning evaluation strategy in 3.2, the system matches the student's ability profile with the translation task graph to form a dynamic programming process of "diagnosis-path-intervention-reevaluation".

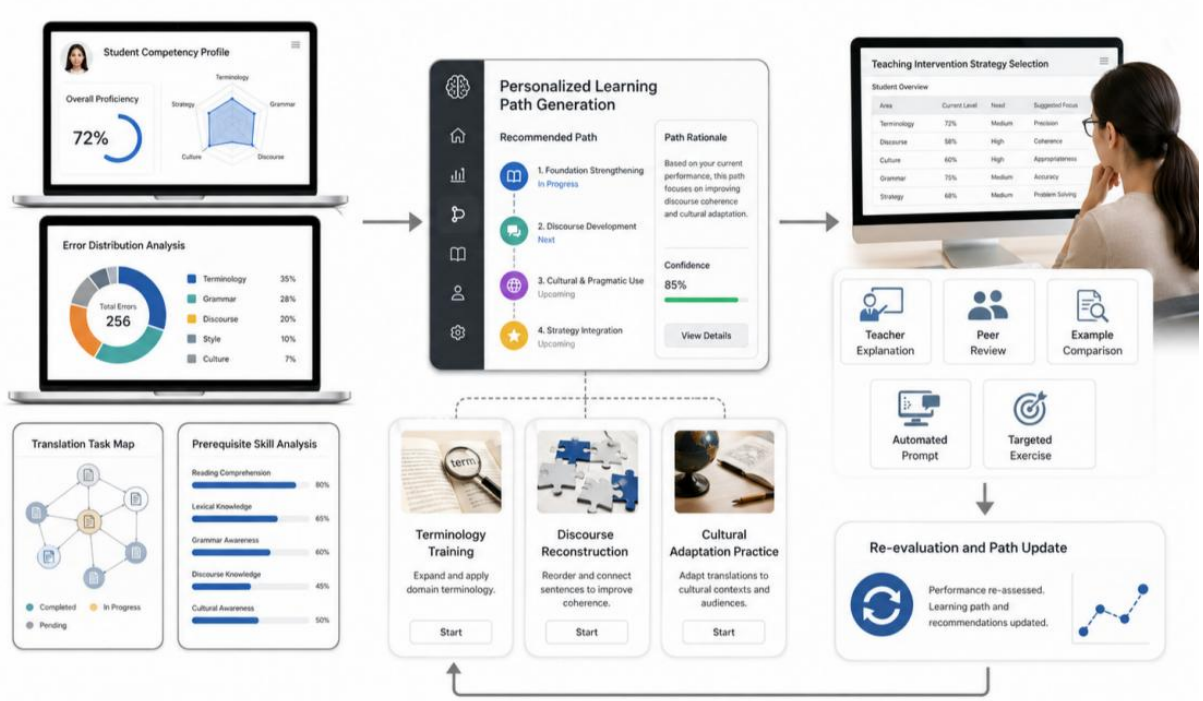


Figure 4: Personalized translation learning path and instructional intervention strategy planning process

Personalized path planning needs to maintain a balance between ability improvement, learning load, and feedback effectiveness. In this paper, let the path payoff of student i choosing learning task u at time step t be as follows.

$$P_{i,t}^u = \eta_1 M_{i,t}^u + \eta_2 G_{i,t}^u - \eta_3 L_{i,t}^u + \eta_4 R_{i,t}^u \quad (6)$$

where, $P_{i,t}^u$ represent the comprehensive benefits of the candidate learning tasks; $M_{i,t}^u$ denote how well the task matches the current weak capability; $G_{i,t}^u$ denote the expected gain in translation quality after completing the task; $L_{i,t}^u$ indicate the learning load required for students to understand the task and digest the feedback; $R_{i,t}^u$ denote the degree of association between the task and the existing error revision records; η_1 to η_4 are the path planning weights. This formula enables the system to avoid simply arranging tasks according to the increasing difficulty, and to select more appropriate training nodes according to the real error structure of students. In the selection of teaching intervention strategies, the system takes teacher explanation, peer evaluation, example comparison, automatic prompt and special training as candidate actions. The selection probability of intervention action is expressed as follows.

$$\pi(a|S_{i,t}) = \frac{\exp(V(S_{i,t}, a)/\tau)}{\sum_{a' \in A} \exp(V(S_{i,t}, a')/\tau)} \quad (7)$$

where, $\pi(a|S_{i,t})$ represents the probability of choosing intervention action a in the current learning state; $V(S_{i,t}, a)$ is the action value. A is the set of candidate intervention actions; Let τ be the policy smoothing parameter. If students continuously mistranslate culturally-loaded words, the system will improve the selection probability of teacher's explanation and case comparison. If students mainly have unstable memory of low-frequency terms, the system

tends to arrange term consolidation training and automatic prompts.

This module makes the translation learning path no longer depend on the fixed chapter order, but continuously adjust with the translation performance and feedback results. The teacher can view each student's recommended path, intervention basis and subsequent performance changes, so as to judge whether the system recommendation conforms to the classroom goals and manually correct the path. In this way, the personalized learning path and the teaching intervention formed a collaborative relationship, which not only retained the teaching judgment of the teacher, but also improved the response speed of the intelligent system to the change of students' translation ability.

3.5 Teachers' digital literacy development support and teaching performance monitoring

Figure 5 shows that this module provides teachers with digital literacy development support oriented to real classroom tasks along the assessment and intervention path generated by the system. The system extracts behavior data from the teacher feedback end, strategy modification records, evaluation result browsing logs and classroom intervention files to form the teacher's operational portrait in the intelligent evaluation environment. Accordingly, the support engine pushed index interpretation training, feedback language editing suggestions, misjudgment case review, intervention program comparison and teaching reflection tips, and continuously monitored teachers' understanding accuracy of system results, intervention timeliness, feedback revision quality and teaching adjustment effectiveness. In this way, the improvement of teachers' digital literacy no longer stopped at the level of tool use, but gradually extended to data interpretation, algorithm collaboration, classroom decision-making and reflection improvement. In this paper, the comprehensive digital literacy score of teacher j at time step t is expressed as follows.

$$T_{j,t} = \omega_1 C_{j,t} + \omega_2 M_{j,t} + \omega_3 E_{j,t} + \omega_4 R_{j,t} - \omega_5 \ln(1 + \tau_{j,t}) \quad (8)$$

where, $T_{j,t}$ represent the comprehensive score of teachers' digital literacy; $C_{j,t}$ denotes the interpretation accuracy of evaluation metrics; $M_{j,t}$ denotes the revision quality to the system feedback results; $E_{j,t}$ denotes the improvement effect of students after the teaching intervention; $R_{j,t}$ represents the integrity and pertinence of teachers' reflection records; $T_{j,t}$ represents the time taken by the teacher to complete the diagnosis and intervention; ω_1 to ω_5 are the weight coefficients. The formula can unify teachers' system use behavior and teaching improvement behavior into the same evaluation framework.

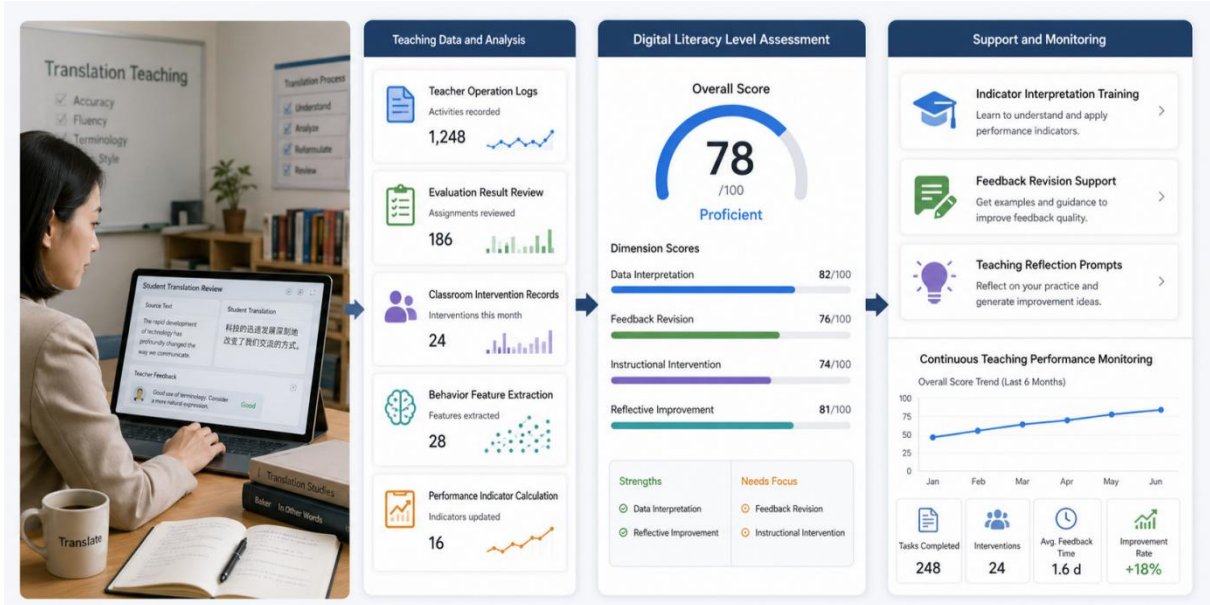


Figure 5: The process of digital literacy development support and teaching performance monitoring for teachers

In order to match the support task with the teacher's current ability, the system further updates the teacher's development task intensity:

$$D_{j,t+1} = D_{j,t} + \phi(T_{j,t} - \theta) + \psi G_{j,t} \quad (9)$$

where, $D_{j,t+1}$ represents the task intensity of teacher development in the next stage; $D_{j,t}$ denotes the current task intensity; Let θ denote the target capability threshold; $G_{j,t}$ represents the growth increment of teachers in this stage; ϕ and ψ are the regulation coefficients. When the teacher's interpretation of the system evaluation results is more accurate and the classroom intervention is more effective, the system will add case analysis and strategy optimization tasks. When teachers still have misreading of results or intervention bias, the system reduces task complexity and prioritizes interpretable demonstration support.

As shown in Table 3, teachers' digital literacy development support and teaching performance monitoring mainly focus on four dimensions. The system does not put teachers in a position of passively receiving technology output, but through a monitoring mechanism that can be tracked, feedback and adjusted, so that teachers can form a stable digital teaching ability in the process of intelligent evaluation of translation teaching.

Table 3: Indicators for monitoring teachers' digital literacy development support and teaching performance

Dimension	Monitoring Indicators	Data Source	Functional Description
Tool application ability	Function call completion rate, result browsing depth	Platform operation logs	Reflects teachers' basic level of using the evaluation system
Data interpretation ability	Indicator judgment accuracy, error identification consistency rate	Evaluation records and manual verification results	Reflects teachers' understanding quality of system outputs
Teaching intervention ability	Intervention response time, student improvement rate after intervention	Classroom intervention archives and student revision records	Reflects teachers' ability to transform evaluation results into teaching actions
Reflective improvement ability	Completeness of reflection texts, strategy revision adoption rate	Reflection logs and system recommendation records	Reflects teachers' ability to continuously optimize teaching decisions

4 Evaluation index system

The evaluation index system of this paper focuses on the accuracy of translation evaluation, the ability of error diagnosis, the efficiency of feedback generation, the effect of learning path optimization, the level of system adaptation, the calculation delay and the improvement effect of teachers 'digital literacy. The evaluation objects not only include the judgment results of the system on the quality of the translation, but also include the students 'modification performance after receiving feedback, the system's strategy adjustment ability, and the teacher's teaching improvement behavior in the intelligent evaluation environment. The indicators keep correspondence with the system modules in Chapter 3, and are used to test whether the reinforcement learning evaluation strategy can support the dynamic optimization of translation teaching process. The translation evaluation accuracy is used to measure the degree of agreement between the system score and the teacher's manual score, and is calculated as follows.

$$A_{eval} = \frac{N_{same}}{N_{all}} \times 100 \quad (10)$$

where, A_{eval} represents the translation evaluation accuracy, N_{same} represents the number of samples whose difference between system score and teacher score is within the allowed interval, and N_{all} represents the total number of test translations. The higher this index is, the more stable the system is in judging the quality level of the translation. Precision, recall, and comprehensive values are used to evaluate the performance of fault diagnosis:

$$F_{err} = \frac{2P_{err}R_{err}}{P_{err} + R_{err}} \quad (11)$$

where, P_{err} represents the proportion of true errors among errors identified by the system, R_{err} represents the proportion of true errors detected by the system, and F_{err} represents the comprehensive effect of error diagnosis. This metric is used to examine the system's ability to identify semantic shift, term misuse, missing translation, addition translation, and discourse cohesion problems. The learning path optimization benefit is used to describe the improvement of student translation quality after personalized task recommendation:

$$G_{path} = \frac{Q_{post} - Q_{pre}}{1 + C_{load}} \quad (12)$$

where, Q_{post} and Q_{pre} represent the translation quality scores before and after the pathway intervention, and C_{load} represents the learning load of students when they complete the recommendation task. This index can avoid the simple pursuit of score improvement and neglect of learning pressure. Feedback generation efficiency measures the ability of a system to generate effective feedback per unit time:

$$E_{feed} = \frac{N_{valid}}{T_{gen}} \quad (13)$$

where, N_{valid} represents the number of feedback items adopted by students or confirmed as valid by teachers, and T_{gen} represents the total time consumption of feedback generation. This index reflects the practicability of the intelligent feedback module in classroom

applications. System adaptive ability index is used to measure the response level of reinforcement learning strategy to student state changes:

$$I_{ada} = \xi_1 \Delta Q + \xi_2 \Delta R + \xi_3 \Delta P - \xi_4 L_{sys} \quad (14)$$

where, ΔQ represents the translation quality improvement, ΔR represents the revision effectiveness improvement, ΔP represents the learning engagement improvement, L_{sys} represents the system response delay, and ξ_1 to ξ_4 are the weight coefficients. Teachers' digital literacy improvement rate is used to evaluate teachers' ability changes in system use, data interpretation and teaching intervention:

$$Dtea = \frac{T_{after} - T_{before}}{T_{before}} \times 100 \quad (15)$$

where, T_{before} and T_{after} represent the digital literacy scores of teachers before and after using the system, respectively. This index corresponds to the teacher development support module and can reflect the promotion effect of intelligent evaluation system on teacher professional development.

5 Results and discussion

This paper carried out the experimental verification of the reinforcement learning empowered intelligent evaluation system for translation teaching. The subjects included 428 translation course students and 8 translation teachers, and the experiment lasted 12 weeks. The system comparison objects are set as the manual scoring rule method, the COMET evaluation model, the LLM prompt evaluation method and the RL-IETS system proposed in this paper. COMET focuses on the semantic quality evaluation of translation. LLM prompt evaluation generates evaluation opinions through fixed prompt words. RL-IETS introduces learning state modeling, reinforcement learning strategy update, error type recognition, personalized learning path planning and teacher digital literacy monitoring mechanism on the basis of translation quality judgment. The experimental results were carried out from nine aspects: dataset construction, translation evaluation accuracy, personalized path optimization, error diagnosis, feedback efficiency, adaptive ability, learning engagement, system delay and teachers' digital literacy improvement.

5.1 Dataset construction

The data in this study come from the translation course platform of colleges and universities, student translation submission records, teacher annotation archives and system interaction logs. A total of 18,420 source text paragraphs, 55,260 student translations, and 92,614 error annotation records were collected, with an average of 129.1 translation versions submitted by each student, resulting in 216.4 traceable error records. The text types covered scientific and technological texts, news texts, business texts, cultural publicity texts and educational instructions texts, of which scientific and technological texts accounted for 22.6%, news texts accounted for 19.8%, business texts accounted for 18.4%, cultural publicity texts accounted for 20.7%, and educational instructions texts accounted for 18.5%. The difference in the proportion of five types of texts is controlled within 4.2 percentage points, which can avoid the experimental results being overly biased towards a certain type of translation task.

Table 4: Data set construction parameters

Data Component	Value or Description
Number of participating students	428
Number of participating teachers	8
Experimental period	12 weeks
Number of source-text sentence segments	18,420
Student translation versions	55,260
Error annotation records	92,614
High-consistency teacher annotations	10,932
Text types	Scientific and technical texts, news, business texts, cultural publicity texts, and educational instructions
Error categories	Semantic deviation, terminology misuse, omission, addition, improper syntactic transformation, insufficient discourse cohesion, and cultural adaptation deviation
Data split	Training set 80%, validation set 10%, test set 10%

As shown in Table 4, the dataset is divided according to 80% of the training set, 10% of the validation set, and 10% of the test set. The training set contains 44,208 translation versions, and the validation set and test set contain 5,526 translation versions each. A total of 12,486 manual annotations were formed on the teacher side, of which 10,932 were judged as high consistency annotations, accounting for 87.6%. The 4,000 samples in the test set were double-checked by two teachers, and Cohen's κ value was 0.86, indicating high consistency of error labels. In the error distribution, 18,764 errors were semantic offset, accounting for 20.3%. There were 15,428 misused terms, accounting for 16.7%; There were 14,372 cases of improper syntactic conversion, accounting for 15.5%; There were 12,906 omissions, accounting for 13.9%; The number of cohesive texts was less than 11,680, accounting for 12.6%; 9,845 additional translations, accounting for 10.6%; There were 9,619 cultural adaptation deviations, accounting for 10.4%. The distribution indicates that students' translation problems are not focused on the lexical level, but involve multiple links such as understanding, conversion, expression and cultural reconstruction at the same time.

5.2 Accuracy analysis of translation evaluation

Translation evaluation accuracy is used to test the degree of agreement between the system rating and the teacher's overall rating. The test set contains a total of 5,526 translation versions, and when the difference between the system rating and the teacher's comprehensive rating does not exceed 3 points, it is judged as a rating match. RL-IETS matched 5,073 samples in the test set, and the accuracy reached 91.8%. The rule scale scoring method matched 4,344 cases, and the accuracy rate was 78.6%. The COMET model matched 4,658 samples with an accuracy of 84.3%. LLM prompt evaluation matches 4,870, with an accuracy of 88.1%. In terms of improvement range, RL-IETS is 13.2 percentage points higher than the manual scoring rule method, 7.5 percentage points higher than the COMET, and 3.7 percentage points higher than the LLM prompt evaluation.

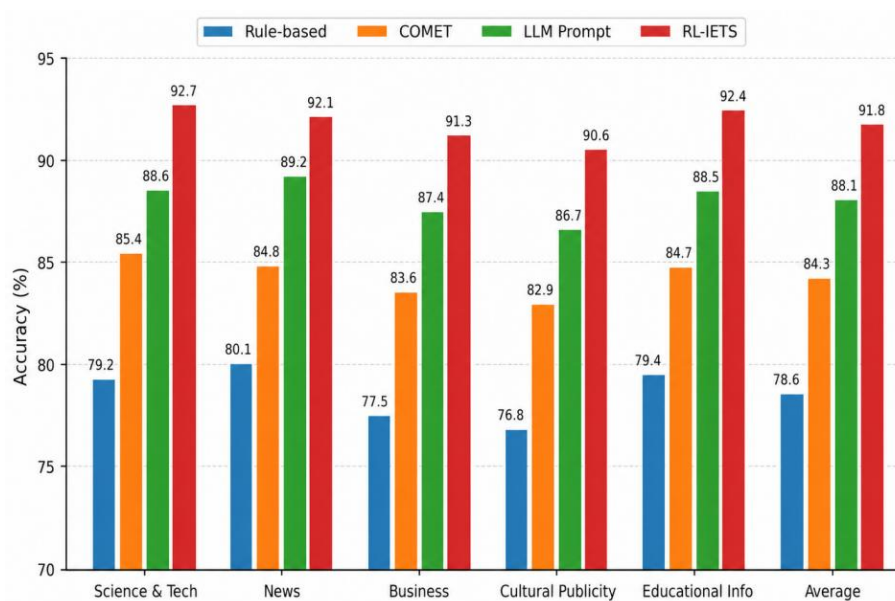


Figure 6: Comparison of translation evaluation accuracy

Figure 6 shows that RL-IETS maintains high accuracy in all five types of texts. The accuracy of scientific text is the highest, which is 92.7%. The reason is that the boundary of scientific text terms is clear, and the system can use the terminology database and the domain corpus to complete the scoring calibration more stable. The accuracy rate of cultural publicity texts is relatively low, 90.6%, mainly because this kind of texts involves cultural imagery, tone conversion and audience acceptance effect, and the evaluation process is more dependent on teachers' experience. In five independent experiments, the accuracy of RL-IETS is 91.2%, 92.1%, 91.6%, 92.4% and 91.7%, respectively, with a mean value of 91.8% and a standard deviation of 0.46, indicating that the model output is stable. Paired t-test showed that the improvement of RL-IETS compared with COMET was statistically significant ($p=0.031$, 95% confidence interval was 90.4%-93.2%).

5.3 Optimization effect of personalized learning path

The optimization effect of personalized learning path is used to judge whether the system can arrange appropriate training tasks according to students' weaknesses. During the 12-week experiment, students in the RL-IETS group completed an average of 34.6 personalized training tasks, which was higher than 27.8 in the LLM recommendation group, 25.3 in the resource similarity recommendation group and 22.4 in the rule recommendation group. Before pathway intervention, the average score of translation in RL-IETS group was 72.3 points, and it increased to 84.9 points after intervention, with an average increase of 12.6 points. The rule recommendation group increased from 71.8 to 78.5, an increase of 6.7 points; The resource similarity recommendation group increased from 72.0 points to 79.8 points, an increase of 7.8 points; The LLM recommendation group increased from 72.5 to 81.4, an increase of 8.9 points.

Table 5 shows that RL-IETS achieves the highest path optimization gains in the five stages of source text understanding, term control, syntactic transformation, discourse reconstruction, and cultural adaptation. The training revenue of cultural adaptation reaches 87.5%, which is 17.7 percentage points higher than that of LLM recommendation and 32.9 percentage points higher than that of rule recommendation. The training gain of text reconstruction is 84.1%, which is 20.3 percentage points higher than that of resource

similarity recommendation. The system log also showed that the proportion of students repeating the same kind of error dropped from 38.2% before the experiment to 17.6%, a decrease of 20.6 percentage points. The results show that the reinforcement learning strategy can continuously adjust the training nodes according to the students' translation modification effect, rather than only pushing learning resources based on the textbook order or text similarity.

Table 5: Personalized learning path optimization effect

Learning Stage	Rule-Based Recommendation / %	Resource Similarity Recommendation / %	LLM-Based Recommendation / %	RL-IETS / %
Source text comprehension training	58.4	64.7	70.5	80.2
Terminology control training	60.1	66.3	72.8	82.7
Syntactic transformation training	57.8	65.1	71.6	83.5
Discourse reconstruction training	55.9	63.8	70.2	84.1
Cultural adaptation training	54.6	62.4	69.8	87.5
Average	57.4	64.5	71.0	83.6

5.4 Translation error diagnosis effect

The effectiveness of translation error diagnosis mainly focuses on the recognition ability of the system for different types of errors. The test set contains 9,261 manually confirmed errors, and RL-IETS correctly identified 8,276 errors, with a comprehensive error diagnosis value of 89.4%. LLM suggested that the evaluation correctly identified 7,725 items, and the comprehensive value was 83.4%. COMET correctly identified 7,410 entries, and the comprehensive value was 80.0%. The rule method correctly identified 6,135 items, and the comprehensive value was 66.2%. RL-IETS is 6.0 percentage points higher than LLM prompt evaluation, 9.4 percentage points higher than COMET, and 23.2 percentage points higher than rule-based method.

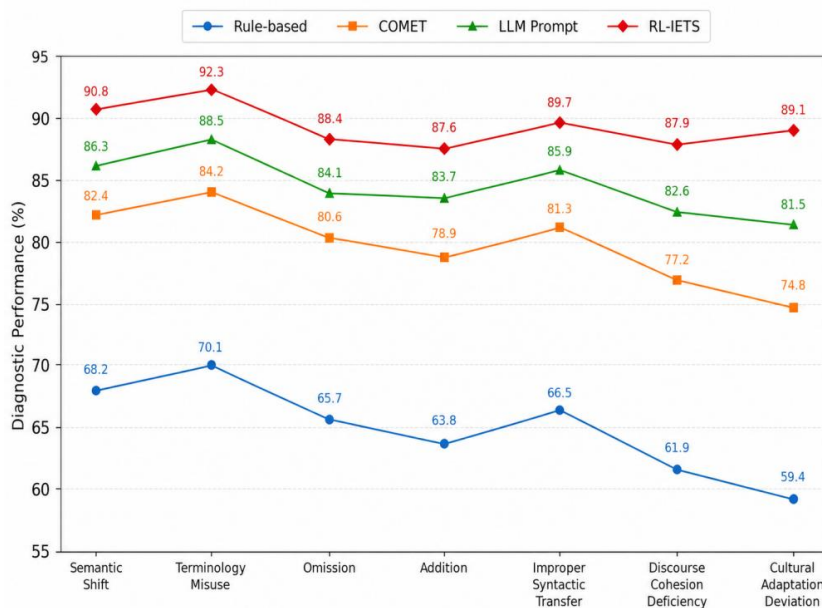


Figure 7: Effect of translation error diagnosis

Figure 7 shows that RL-IETS achieves the best recognition of term misuse with a comprehensive value of 92.3%, which is mainly due to the common constraints of the course term bank, professional parallel corpus and teacher annotation rules. The comprehensive value of semantic deviation is 90.8%, syntactic improper conversion is 89.7%, and cultural adaptation deviation is 89.1%. The comprehensive value of additional translation recognition is 87.6%, which is relatively low. The reason is that some additional translation content has explanatory supplementary function, and whether it constitutes an error needs to be judged according to the specific context. Misjudgment analysis showed that the proportion of RL-IETS misjudging insufficient cohesion as syntactic conversion errors was 5.8%, which was lower than 9.7% of LLM prompt evaluation. This indicates that learning state modeling and teacher annotation calibration are able to reduce confusion between adjacent error types.

5.5 Efficiency analysis of intelligent feedback generation

The efficiency of intelligent feedback generation not only reflects the response speed of the system, but also reflects whether the feedback content can be adopted by students and converted into translation. In the testing phase, RL-IETS generated a total of 18,360 feedback suggestions, of which 15,012 were adopted by students or confirmed as effective by teachers, and the effective feedback rate was 81.8%. In the simple sentence paragraph task, the average generation time of RL-IETS is 1.9 seconds, 0.9 seconds faster than that of LLM prompt evaluation. In the complex paragraph task, the average generation time of RL-IETS is 2.7 seconds, which is 1.9 seconds faster than that of LLM prompt evaluation. For the paragraph text task, RL-IETS generated 3.4 seconds on average, which is 3.5 seconds faster than LLM prompt evaluation. The average human annotation time on paragraph text is 312.0 seconds, which is about 91.8 times that of RL-IETS.

Table 6: Efficiency analysis of intelligent feedback generation

Text Complexity	Method	Average Generation Time / s	Average Number of Effective Feedback Items per Submission	Student Adoption Rate / %
Simple sentence segments	Manual annotation	96.0	4.1	76.4
Simple sentence segments	LLM prompt evaluation	2.8	3.8	71.2
Simple sentence segments	RL-IETS	1.9	4.3	82.5
Complex sentence segments	Manual annotation	168.0	5.2	78.1
Complex sentence segments	LLM prompt evaluation	4.6	4.9	73.8
Complex sentence segments	RL-IETS	2.7	5.5	82.1
Paragraph texts	Manual annotation	312.0	7.4	79.3
Paragraph texts	LLM prompt evaluation	6.9	6.8	74.6
Paragraph texts	RL-IETS	3.4	7.6	81.7

As shown in Table 6, the average number of effective feedback items of RL-IETS under the three types of text complexity is 4.3, 5.5 and 7.6 respectively, and the student adoption rate is 82.5%, 82.1% and 81.7% respectively. The fluctuation range of the adoption rate is only 0.8 percentage points, indicating that the system feedback quality is stable. The adoption rate of LLM prompt evaluation in paragraph text was 74.6%, 7.1 percentage points lower than that of RL-IETS. Teacher review records showed that 76.3% of RL-IETS feedback could be directly used for classroom evaluation, which was higher than 68.5% of LLM prompt evaluation. The main reason for the difference is that RL-IETS does not simply generate a

complete alternative translation, but combines the error level, student status and training path to generate local hints, reasons and revision directions, so that the feedback is more in line with the requirements of formative evaluation in translation teaching.

5.6 System Adaptive Capability index

The system adaptive ability index is used to measure the response level of reinforcement learning strategy to the change of student learning state. At the beginning of the experiment, the adaptive ability index of RL-IETS is 78.5%. In the fifth week, it increased to 86.9%. In the 9th week, it reached 91.1%. In the 12th week, it further improved to 93.8%. Within 12 weeks, RL-IETS has increased by 15.3 percentage points, and the average increase is about 2.55 percentage points every two weeks. Rule recommendation increased from 60.2% to 62.9%, an increase of only 2.7 percentage points; Resource similarity recommendation increased 5.8 percentage points from 66.1% to 71.9%; LLM recommendations increased by 8.6 percentage points from 72.4% to 81.0%.

Figure 8 shows that the adaptation curve of RL-IETS improves more significantly after week 5. The reason is that in the first four weeks, the system mainly pushes tasks based on the initial translation quality and error labels. After the fifth week, students' revision records, feedback adoption rate and teacher intervention results are included in the reward update process, and the strategy selection is more stable. The system log shows that the system needs 2.8 translation submissions on average to adjust the learning path in weeks 1-3. By week 10 to 12, the number dropped to 1.4. For students with two consecutive term misuse, the probability of term consolidation training pushed by the system increased from 42.5% at the beginning of the experiment to 78.9%. For students with obvious cultural adaptation deviation, the trigger rate of teacher intervention reminder increased from 18.6% to 44.2%.

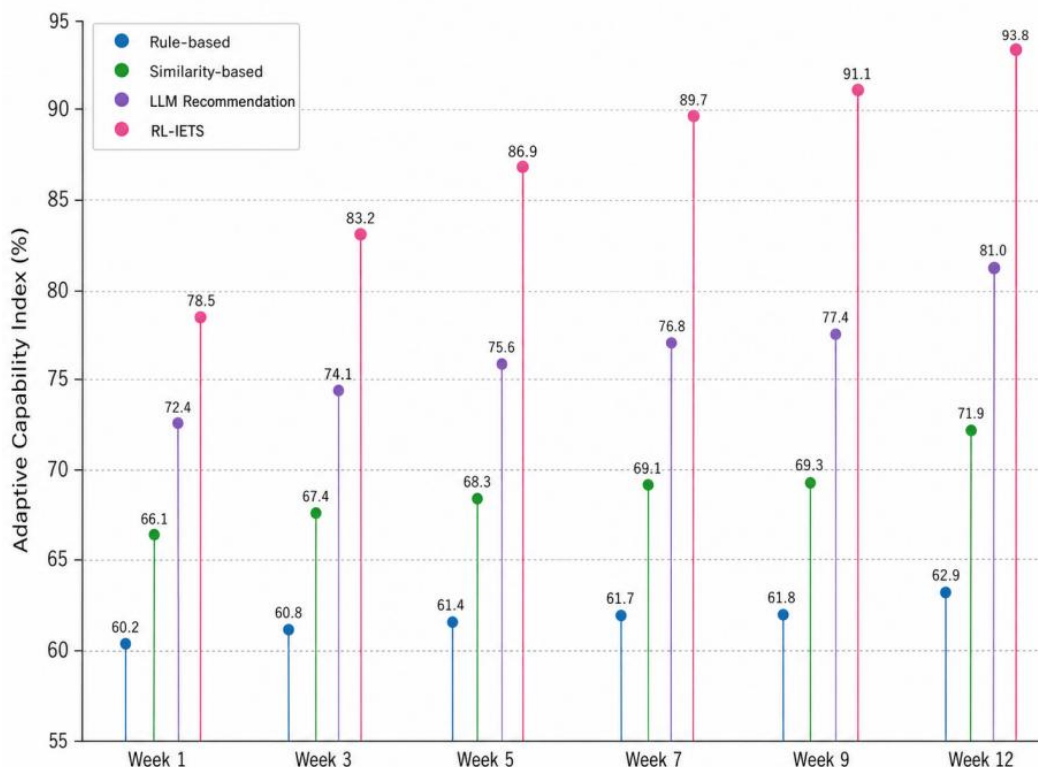


Figure 8: System adaptive ability index change

5.7 Learning Engagement score

The learning engagement score was composed of feedback reading time, task completion rate, translation revision times, platform interaction frequency and stage dropout rate. The participation score of RL-IETS group increased from 73.4 points in week 1 to 88.6 points in week 12, an increase of 15.2 points. LLM assisted feedback group increased from 71.5 points to 80.1 points, an increase of 8.6 points; The score of conventional teaching group increased from 68.9 to 74.2, an increase of 5.3 points. The improvement of the RL-IETS group was 1.77 times that of the LLM-assisted feedback group and 2.87 times that of the conventional teaching group, respectively.

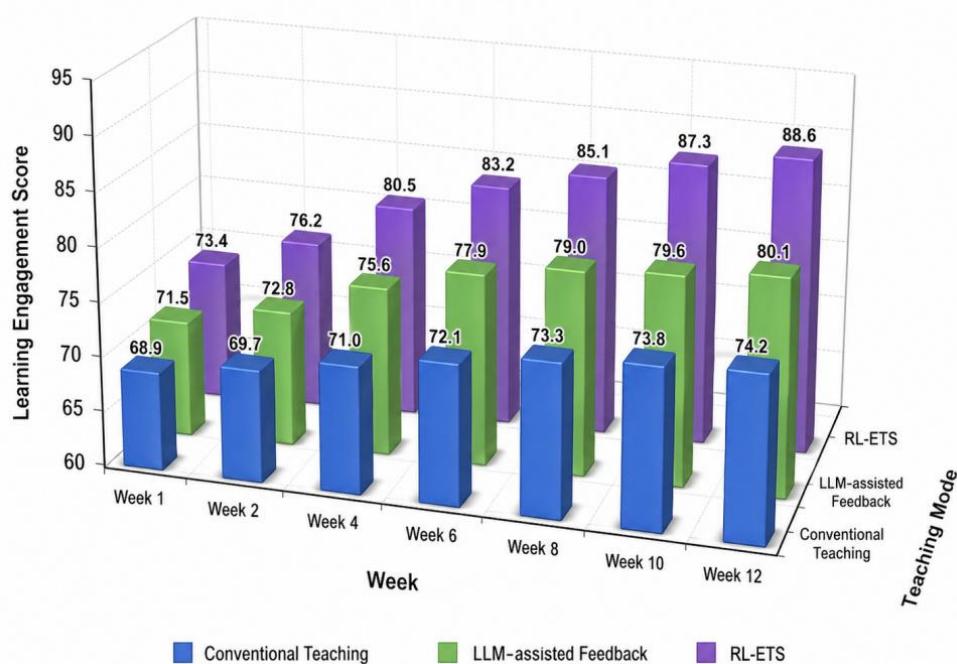


Figure 9: Changes in learning engagement score

As shown in Figure 9, the RL-IETS group experienced a more significant increase in engagement after week 4. Platform logs showed that students in the RL-IETS group logged in 5.8 times per week on average, which was higher than 4.6 times in the LLM assisted feedback group and 3.9 times in the regular instruction group. The average number of translation revisions per student was 2.7 in the RL-IETS group, 1.9 in the LLM-assisted feedback group, and 1.4 in the conventional teaching group. The feedback reading behavior also showed differences. The average reading time of students in the RL-IETS group was 42.6 seconds per feedback, which was higher than 31.4 seconds in the LLM-assisted feedback group. Within 12 weeks, the task completion rate was 91.3% in the RL-IETS group, 84.6% in the LLM assisted feedback group, and 79.2% in the conventional teaching group. In terms of stage dropout rate, the RL-IETS group was 3.5%, which was lower than the LLM-assisted feedback group (6.8%) and the conventional teaching group (9.4%). These data show that the personalized path and immediate feedback can enhance students' willingness to continuously revise the translation.

5.8 Analysis of system computation delay

The system computation delay is related to whether the intelligent evaluation can enter the real classroom. The experiment tests the average response time under sentence level,

compound sentence level, paragraph level and discourse level tasks. RL-IETS has an average response time of 795 ms, which is higher than COMET's 515 ms, but lower than the LLM prompt evaluation's 1,390 ms. In the sentence-level task, RL-IETS has a response time of 310 milliseconds. The paragraphlevel task was 890 milliseconds, still under 1 second; The document-level task was 1,460 ms, which, while higher than COMET's 910 ms, was 920 ms, or 38.7% less than LLM's 2,380 ms for prompt evaluation.

Table 7: Analysis of system computation delay

Text Level	COMET Evaluation / ms	LLM Prompt Evaluation / ms	RL-IETS / ms
Sentence level	180	640	310
Complex sentence level	350	980	520
Paragraph level	620	1,560	890
Discourse level	910	2,380	1,460
Average	515	1,390	795

Table 7 shows that RL-IETS is not the method with the lowest latency, but it maintains a good balance between latency and functional complexity. COMET is mainly used to score the quality of translation, with centralized model structure and fast response speed. LLM prompt evaluation requires the generation of a long explanation text, and the response time increases significantly. RL-IETS performs semantic encoding, error identification, strategy selection, feedback generation, and path update simultaneously, so the computational burden is higher than that of the single scoring model. The analysis of latency composition shows that semantic encoding accounts for about 38.5%, error type recognition accounts for about 27.4%, feedback action selection and path update accounts for about 21.6%, data reading, writing and logging account for about 12.5%. In the processing test of 55,260 translation versions, RL-IETS has an average throughput of 18.7 sentences per second and a peak of 24.3 sentences per second, which can meet the classroom scenario of 30 to 50 students submitting short translations at the same time.

5.9 Analysis on the Improvement Effect and development Path of Teachers 'digital Literacy

The improvement effect of teachers 'digital literacy was evaluated from four dimensions of tool application ability, data interpretation ability, teaching intervention ability and reflection improvement ability. The comprehensive score of digital literacy of 8 teachers was 67.8 before the experiment and increased to 86.5 after the experiment, with an increase of 18.7 points and an improvement rate of 27.6%. Among them, the data interpretation ability increased from 64.6 points to 85.7 points, an increase of 21.1 points, which was the dimension with the largest improvement. Teaching intervention ability increased from 66.1 to 86.9, an increase of 20.8 points; Tool application ability increased from 72.4 to 89.2, an increase of 16.8 points; Reflective improvement ability increased from 68.0 to 84.3, an increase of 16.3 points.

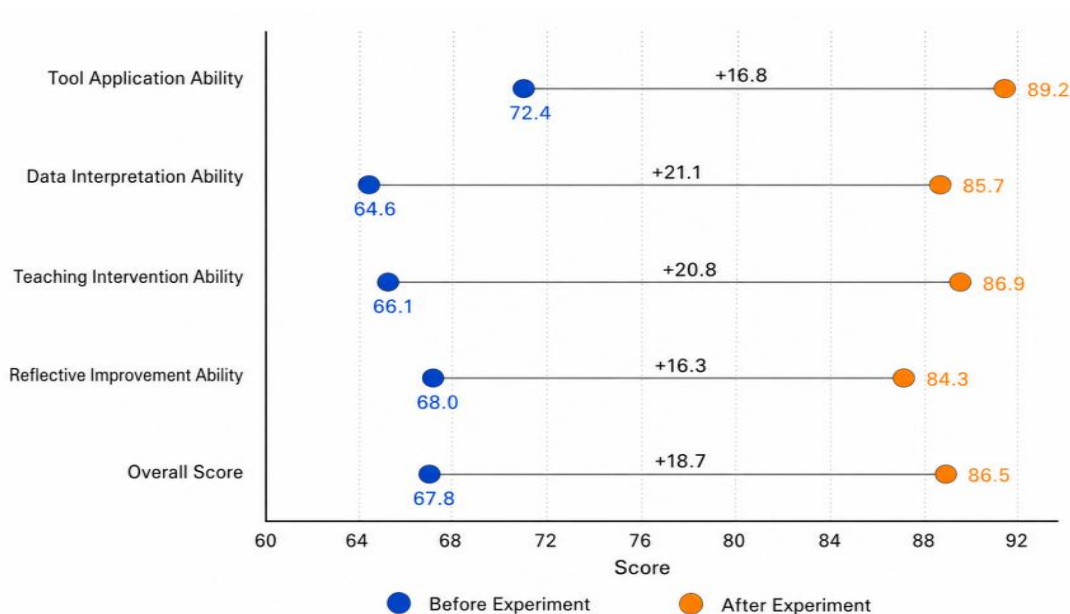


Figure 10: The improvement effect of teachers' digital literacy

Figure 10 shows that the improvement of teachers' ability does not simply come from the skillful operation of the platform, but is reflected in the enhanced ability of data interpretation, feedback revision and classroom intervention. The teacher behavior log showed that the average number of times teachers checked the system diagnosis report per week was 2.1 before the experiment, which increased to 6.4 after the experiment. The proportion of teachers manually revising the system feedback increased from 18.5% to 46.7%. In terms of classroom intervention, the number of special topic comments based on high frequency error reports of the system increased from 9 times in the first 4 weeks to 21 times in the last 4 weeks. The average number of words in teachers' reflection texts increased from 146 words per week before the experiment to 382 words per week after the experiment, and the proportion of keywords involved in the reflection content such as "data basis", "error distribution" and "path adjustment" increased from 23.8% to 71.4%. This indicates that the development path of teachers' digital literacy should go from system operation to data interpretation, and then to teaching intervention and reflection improvement.

6 Conclusions

Focusing on the problems of feedback lag, insufficient diagnosis granularity, fixed learning path and lack of process support for teachers' digital literacy development in translation teaching evaluation, this paper designs an intelligent evaluation system for translation teaching empowered by reinforcement learning. Based on student translations, revision records, error annotations, teacher annotations and platform behavior logs, the system constructs a learning state representation, and dynamically selects scoring, error diagnosis, feedback generation, path recommendation and teacher intervention strategies through the state-action-reward mechanism. The research shows that the system can extend the translation evaluation from single result judgment to continuous learning diagnosis, and make the abilities of semantic understanding, term control, syntactic transformation, discourse reconstruction and cultural adaptation in translation teaching obtain a clearer computational expression.

The experimental results show that RL-IETS is superior to the comparison methods in translation evaluation accuracy, error diagnosis comprehensive value, path optimization benefit, feedback acceptance rate and system adaptive ability. Among them, the average evaluation accuracy of the system reached 91.8%, the comprehensive value of error diagnosis reached 89.4%, the profit of personalized path optimization reached 83.6%, the effective feedback rate reached 81.8%, the system adaptive ability index reached 93.8% in the 12th week, and the score of learning participation increased by 15.2 points. The comprehensive score of teachers' digital literacy increased from 67.8 to 86.5. These results showed that reinforcement learning strategies could continuously adjust teaching support methods according to students' translation changes, and also help teachers shift from simple use of tools to understanding data, correcting feedback, and implementing precise intervention. The research value of this paper is mainly reflected in three aspects: first, the translation learning process is transformed into a traceable and interpretable state sequence; Second, the intelligent evaluation, feedback generation and learning path planning are integrated into the same decision-making framework. The third is to embed teachers' digital literacy development into real classroom evaluation tasks, so that teachers' professional judgment and system intelligent output form synergy. Subsequent research can further expand the multilingual translation task, introduce multi-modal data such as eye movement, keyboard input and oral thinking report, and verify the stability, transferability and teacher development support effect of the system in a longer period of real teaching scenarios.

Author's Profile

Weijia Liu was born in Luoyang, Henan, P.R. China, in 1990. She received the bachelor's degree from China University of Mining and Technology, P.R. China. Now, she works in School of International Studies, Luoyang Institute of Science and Technology. Her research interest include translation teaching, translation techniques and digital literacy. In recent years, she has led or participated in the development of micro-courses, smart courses, and high-quality courses, including Chinese-English Translation, Digital Literacy and Critical Thinking Tools, and Translation Project Management and CAT Technology, several of which have been approved as university-level or provincial-level construction projects. mail: wei_jia_liu@163.com

Ruirui Zhang received the bachelor's degree from Guangdong University of Foreign Studies. Now, she works in School of International Studies, Luoyang Institute of Science and Technology, serving as the vice dean. Her research interest include business English teaching and language and culture studies. Her funded project, Driving the Cultivation of Innovative Foreign Language Talents through the Construction of Smart Course Clusters from the Perspective of New Liberal Arts, was approved as a Key Project for Higher Education Teaching Reform Research and Practice in Henan Province for 2026. E-mail: lylg200807@163.com

Asel Musurapshaevna Toksonalieva received the doctoral degree from Beijing Foreign Studies University. Now, she works in Kyrgyz-Chinese Institute of Kyrgyz National University named after Zh. Balasagyn, serving as the Director of the Teaching and Research Office for Chinese Language and Literature. Her research interest include language teaching and contrastive linguistics. E-mail: beijing1104@yandex.com

Funding

This work was supported by the 2026 Key Project of Higher Education Teaching Reform Research and Practice in Henan Province, titled Practice of Innovative Foreign Language Talent Cultivation through the Construction of Smart Course Clusters from the Perspective of New Liberal Arts.

References

- [1] Urlaub P, Dessen E. Machine translation and foreign language education[J]. *Frontiers in Artificial Intelligence*, 2022, 5: 936111.
- [2] Jolley J R, Maimone L. Thirty years of machine translation in language teaching and learning: A review of the literature[J]. *L2 Journal: An Open Access Refereed Journal for World Language Educators*, 2022, 14(1).
- [3] Klimova B, Pikhart M, Benites A D, et al. Neural machine translation in foreign language teaching and learning: a systematic review[J]. *Education and Information Technologies*, 2023, 28(1): 663-682.
- [4] Tavares C, Tallone L, Oliveira L, et al. The challenges of teaching and assessing technical translation in an era of neural machine translation[J]. *Education Sciences*, 2023, 13(6): 541.
- [5] Kirchhoff P. Machine translation in English language teaching[J]. *ELT Journal*, 2024, 78(4): 393-400.
- [6] Moorhouse B L. Generative artificial intelligence and ELT[J]. *ELT Journal*, 2024, 78(4): 378-392.
- [7] De Wilde V. Can novice teachers detect AI-generated texts in EFL writing?[J]. *ELT Journal*, 2024, 78(4): 414-422.
- [8] Kaya M H. A professional training to make English language instructors AI-ready[J]. *Elt journal*, 2024, 78(4): 466-477.
- [9] Hiebl B, Gromann D. Quality in human and machine translation: An interdisciplinary survey[C]//*Proceedings of the 24th Annual Conference of the European Association for Machine Translation*. 2023: 375-384.
- [10] Rei R, De Souza J G C, Alves D, et al. COMET-22: Unbabel-IST 2022 submission for the metrics shared task[C]//*Proceedings of the Seventh Conference on Machine Translation (WMT)*. 2022: 578-585.
- [11] Kocmi T, Federmann C. Large language models are state-of-the-art evaluators of translation quality[C]//*Proceedings of the 24th Annual Conference of the European Association for Machine Translation*. 2023: 193-203.
- [12] Kocmi T, Federmann C. GEMBA-MQM: Detecting translation quality error spans with GPT-4[C]//*Proceedings of the Eighth Conference on Machine Translation*. 2023:

768-775.

- [13] Manakhimova S, Avramidis E, Macketanz V, et al. Linguistically motivated evaluation of the 2023 state-of-the-art machine translation: Can ChatGPT outperform NMT?[C]//Proceedings of the Eighth Conference on Machine Translation. 2023: 224-245.
- [14] Macken L. Machine translation meets large language models: Evaluating ChatGPT's ability to automatically post-edit literary texts[C]//Proceedings of the 1st Workshop on Creative-text Translation and Technology. 2024: 65-81.
- [15] Koneru S, Exel M, Huck M, et al. Contextual refinement of translations: Large language models for sentence and document-level post-editing[C]//Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers). 2024: 2711-2725.
- [16] Leiter C, Eger S. Prexme! large scale prompt exploration of open source llms for machine translation and summarization evaluation[C]//Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing. 2024: 11481-11506.
- [17] Condor A, Pardos Z. A deep reinforcement learning approach to automatic formative feedback[C]//Proceedings of the 15th International Conference on Educational Data Mining. 2022: 662.
- [18] Vassoyan J, Vie J J, Lemberger P. Towards scalable adaptive learning with graph neural networks and reinforcement learning[J]. arXiv preprint arXiv:2305.06398, 2023.
- [19] Memarian B, Doleck T. A scoping review of reinforcement learning in education[J]. Computers and Education Open, 2024, 6: 100175.
- [20] Amin S, Uddin M I, Alarood A A, et al. Smart E-learning framework for personalized adaptive learning and sequential path recommendations using reinforcement learning[J]. IEEE Access, 2023, 11: 89769-89790.
- [21] Galindo-Domínguez H, Delgado N, Campo L, et al. Relationship between teachers' digital competence and attitudes towards artificial intelligence in education[J]. International Journal of Educational Research, 2024, 126: 102381.
- [22] Gisbert Cervera M, Caena F. Teachers' digital competence for global teacher education[J]. European Journal of Teacher Education, 2022, 45(4): 451-455.
- [23] Khalil H, Alsenaidi S. Teachers' Digital Competencies for Effective AI Integration in Higher Education in Oman[J]. Journal of Education and e-learning Research, 2024, 11(4): 698-707.
- [24] Almatrafi O, Johri A, Lee H. A systematic review of AI literacy conceptualization, constructs, and implementation and assessment efforts (2019–2023)[J]. Computers and Education Open, 2024, 6: 100173.