

Construction of Learning Resources for International Chinese Language Education Based on Sensor Technology and Knowledge Graphs

Yue Fu,^{1,2} Lei Zhao,³ Borui Zheng,⁴ Yirong Wang,⁵ and Liqing Yang^{1*}

¹College of International Chinese Language Education, Yunnan Normal University, Kunming 650500, China

²College of International Chinese Language Education, College of Arts and Sciences, Kunming 50222, China

³China Mobile Communications Group, Kunming 650228, China

⁴Department of Asian Studies, University College Cork, Cork, T12 K8AF, Ireland

⁵Wuhua District Hongqi Primary School, Kunming 650106, China

(Received August 21, 2025; accepted March 9, 2026)

Keywords: knowledge graph, sensor technology, international Chinese language education, dynamic update, learning resources

We developed a method for updating and constructing Chinese language teaching resources by integrating sensor technology and knowledge graphs. The method addresses the challenges of a mismatch between resource supply and learner demand through a closed loop of perception–analysis–update. Wearable sensors and AI cameras were used to collect real-time, quantitative multimodal data on learners’ physiological and behavioral states, including body temperature, heart rate, and classroom interactions. These sensor data, along with learning platform data, were used to train an eXtreme Gradient Boosting model. The model achieved a prediction accuracy of 89.2% and an area under the receiver operating characteristic curve of 0.93, indicating the accurate distinction of knowledge entities that need updating and those that do not. The feature importance analysis revealed that user ratings (0.603) and recency (0.226) were the most influential factors for predicting update necessity. The knowledge graph was iteratively updated through a multistep process including pattern mining, filtering, and a final review by experts. The resulting knowledge graph incorporated nodes for new content, such as internet slang and cross-cultural variations in festivals, demonstrating the method’s ability to adapt to linguistic evolution and cultural nuances. Through the establishment of a closed-loop architecture, multimodal sensor data, including physiological photoplethysmography signals and behavioral time-of-flight imaging, are used for the expansion and weight adjustment of a domain-specific knowledge graph. This cognitive-aware update mechanism ensures that learning resources evolve along with real-time learner demands, providing a scalable blueprint for intelligent, sensor-driven knowledge management systems in various disciplines. The results of this study also underscore the role of sensor data in developing contextualized, personalized, and optimized digital learning resources that lead to a learner-centered learning environment.

*Corresponding author: e-mail: 2273310002@ynnu.edu.cn
<https://doi.org/10.18494/SAM5907>

1. Introduction

With the advancement of information and AI technologies, digital resources supporting international Chinese language education have become increasingly diverse and expansive in scale. As of February 2021, there were 3700 different types of electronic textbook available globally in Chinese language education. There are 4685 microcourses with their resources and 272 Chinese teaching applications on diverse platforms.⁽¹⁾ Digital technology has been applied to construct international Chinese language learning resources.⁽²⁾ By integrating digital technology, the quality and efficiency of learning resources have been enhanced. Despite the transformation of traditional learning resources into digital ones, a mismatch between resource supply and learner demand, the lack of contextualized data, and insufficient dynamic update mechanisms exist as challenges to address.⁽³⁾ Traditional learning resources based on question-answer systems require rule-based engines or keyword matching methods for users to learn Chinese. These methods often fail to meet dynamic and personalized needs, especially when learners practice complex sentence structures and explore contextual variability in the Chinese language, which often makes learning inefficient.⁽⁴⁾ Current educational resources are often mapped to standards such as the Hanyu Shuiping Kaoshi (HSK), which is the international standardized test for Chinese language proficiency. However, these standards often fail to capture the rapid evolution of internet slang or varying cultural nuances, which leads to a mismatch between static materials and dynamic learner needs.⁽⁵⁾

In the AI era, digital learning resources have become fundamental for the quality enhancement of Chinese language education.⁽⁶⁾ Research on effective Chinese language education has been extensively conducted,⁽⁷⁾ in which classification systems of Chinese learning resources have been developed to construct learning resources and databases.⁽⁸⁾ With the application of wearable sensors, edge computing devices, IoT, and knowledge graphs, learning resources can be further developed on the basis of new technological approaches and paradigms. Knowledge graphs, in the form of a formal semantic network of nodes, edges, and attributes, are used to model Chinese characters, vocabulary, grammar, cultural concepts, and their dynamic interrelations, for the scalable organization of learning resources.⁽⁹⁾ Sensors are employed to capture cognitive, emotional, and behavioral signals of learners in real or virtual contexts through data collection at a millisecond-level interval, enabling high-resolution and contextualized resource updates. While knowledge graphs are used to present static knowledge, sensor data are used to analyze learning patterns and assess the learner's attitude and responses. However, knowledge graphs and sensor data have been used separately, which hinders the integration required to construct a closed-loop framework that enables a positive feedback cycle among the knowledge graph, sensor perception, and learning resources.

Therefore, we studied how to leverage the synergy of knowledge graphs and sensor technology in constructing Chinese language learning resources through the real-time collection of multimodal data on learners' physiology, behavior, and cognition, and the data analytics through adaptive iteration and personalized adaptation. In this study, a sensor-driven cognitively adaptive system and its underlying architecture were constructed for real-time data perception using wearable and ambient sensors, predictive cognitive modeling, and automated knowledge

base updating for international Chinese language education.⁽¹⁰⁾ The sensor-derived data used in this study can also be used to update a procedural knowledge graph to identify specific skill gaps in other fields of education. The results of this study provide a basis for the development of a generalized solution for human-centric intelligent systems across various disciplines.

2. Materials and Methods

2.1 Application of sensor technology

Sensor technology in language teaching has been used in pre-class, in-class, and post-class phases to monitor interactions between teachers and learners. Radosavljevic *et al.* established a smart classroom learning model based on the concept of ambient intelligence. By analyzing the data collected, the smart classroom offers effective learning strategies to learners following the criteria that align with the expected learning outcomes.⁽¹¹⁾ Dai developed an interactive teaching method for financial accounting by integrating a smart classroom to improve the quality of financial education through interventions before, during, and after teaching classes.⁽¹²⁾ Considering the increasing availability and potential of 5G technology, Rong adopted the technology in improving smart classroom instruction.⁽¹³⁾ He proposed the adoption of human-computer interaction and the integration of virtual and real environments in teaching. Innovative teaching models are constructed on the basis of the developed method, which can also be used as a reference for further enhancement in teaching.⁽¹³⁾

The advancement of cloud computing, IoT, big data, and AI has enabled the integration of sensor technology in learning for educational informatization and new teaching and learning model development and management. With sensor technology, the physiological signals of learners are monitored to observe learners' cognitive states and emotional responses. On the basis of the cognitive load theory,⁽¹⁴⁾ we constructed an intelligent data collection system using advanced sensors to monitor and analyze the learner's status. The collected data are used for real-time feedback and intervention and the optimization of the knowledge graph, to form a closed-loop perception-analysis-update scheme. Using the scheme, we constructed intelligent and adaptable learning resources, as shown in Fig. 1. The figure shows the synchronization between physical sensing (multimodal data collection), the analytical core (machine learning for cognitive state assessment), and the dynamic evolution of the knowledge graph. The loop ensures that sensed learner difficulties trigger immediate resource adjustments.

In this study, an AI camera was employed as an edge-integrated unit comprising a high-resolution CMOS image sensor and an indirect time-of-flight (iToF) depth sensor. The RGB component includes a 1/2.8-inch complementary metal-oxide-semiconductor sensor with a resolution of 1920×1080 pixels (full high definition), capable of capturing data at 60 frames per second. This high resolution is essential for the accurate extraction of linguistic entities from scanned text and for identifying subtle classroom interactions such as nodding or hand raising. The depth sensing was conducted using an iToF sensor that emits modulated near-infrared light at wavelength of a 940 nm to obtain a depth map with a resolution of 640×480 (video graphics array) and an accuracy of 1% in a 0.5–5.0 m range. iToF technology enables differentiation

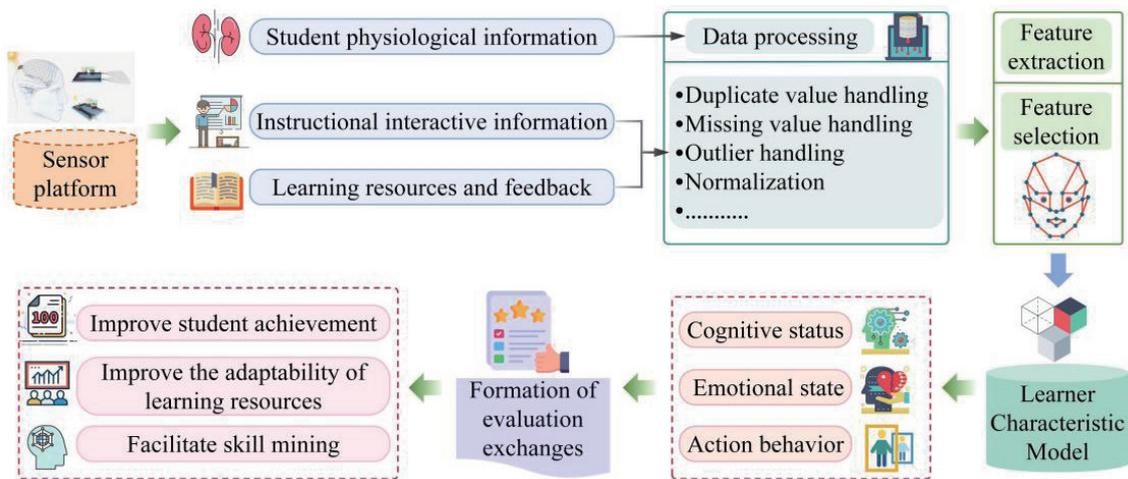


Fig. 1. (Color online) Closed-loop perception–analysis–update architecture.

between the learner and background objects with high precision, enabling the calculation of distance-based engagement metrics without the effect of ambient light. These devices are synchronized to provide registered red-green-blue-depth frames for the behavioral analysis and cognitive state modeling modules.

2.2 Knowledge graphs

A knowledge graph is used to examine the relationships between entities and concepts, modeled as nodes and edges in a topological structure.⁽¹⁵⁾ With the addition of new knowledge, the knowledge graph needs to be continuously updated to maintain its accuracy and relevance in presenting the relationships. When updating a knowledge graph, new entities (knowledge) must have mutually associated information with that of existing entities to be directly integrated and avoid conflicts in information. However, this method might introduce inaccuracies by introducing outdated knowledge.⁽¹⁶⁾ Therefore, triple classification is used to calculate the confidence of a triple's truth.⁽¹⁷⁾ Yao *et al.* proposed rule-based guidance and estimated validity on the basis of the relationships between entities to minimize such inaccuracies.⁽¹⁸⁾ Following the guidance, Iwata *et al.* enhanced the objective function and bimodal function approximation algorithm to obtain a knowledge graph. However, this method increases computational complexity.⁽¹⁹⁾ To reduce the complexity, Dvořák *et al.* developed an enumeration-based pattern classification method, but it leads to omissions and fails to ensure filtering accuracy.⁽²⁰⁾ When comparing large amounts of entities, the complexity of their types and relationships increases the difficulty of filtering them to enhance the accuracy during synchronization.⁽²⁰⁾

In addition to entity updates, knowledge inference is also used to address the incompleteness of relationships between previous and new entities. In this method, unknown or implicit thematic relations are discovered or deduced. The Markov logic network (MLN) is used to combine expert-defined logical rules with probabilistic graphical models to build a network and perform inference.⁽²¹⁾ In this method, Pujara *et al.* introduced confidence values by using probabilistic

logic,⁽²²⁾ whereas Kuželka and Davis applied MLN weights to knowledge graphs under data-sparse conditions.⁽²³⁾ Halaschek-Wiener *et al.* further improved the knowledge graph method by using a description logic reasoning algorithm and controlling assertion boxes.⁽²⁴⁾ Building on this, Calvanese *et al.* proposed a cognitive-based first-order query language to handle incomplete information in the reasoning process of knowledge graphs,⁽²⁵⁾ and Li *et al.* extended traditional description logic using fuzzy logic for reasoning with uncertainty.⁽²⁶⁾

On the basis of previous research results, we adopted a pattern-mining-based method to update knowledge graphs (Fig. 2) to improve efficiency and accuracy in the entity update process.

3. Teaching Resource Design and Implementation

3.1 Modules

A dual-module architecture was developed for teaching resource construction in this study, including data collection and knowledge graph modules. The data module includes real-time sensor data on learners' behavior, which is transmitted synchronously to the knowledge graph module, which stores and processes the data, and updates the knowledge graph based on the results of the data analysis and learner feedback. Figure 3 shows the interaction between the data module, which processes raw sensor feeds into behavioral features, and the knowledge graph module, which manages the semantic structure of learning resources and executes updates on the basis of the data module's output.

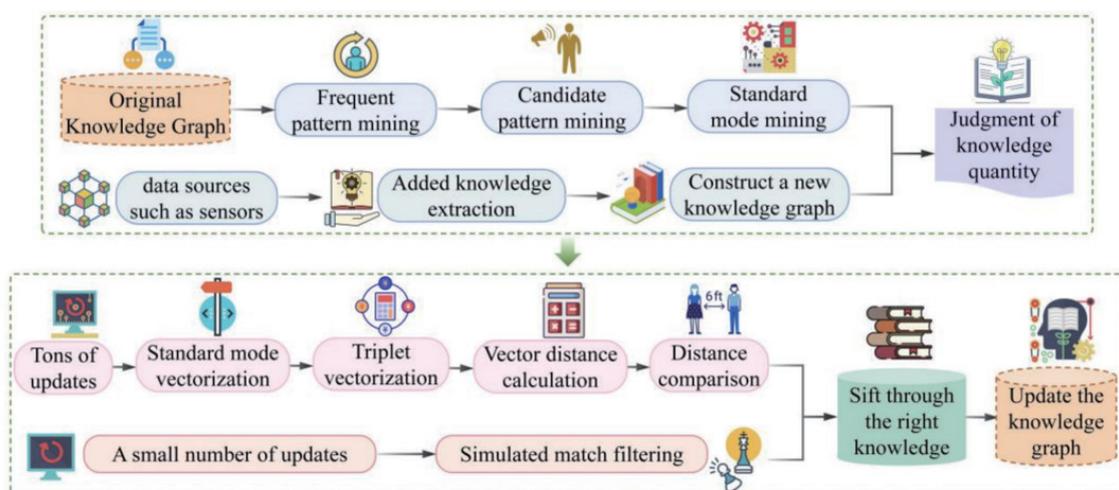


Fig. 2. (Color online) Entity update in knowledge graph using pattern-mining-based method.

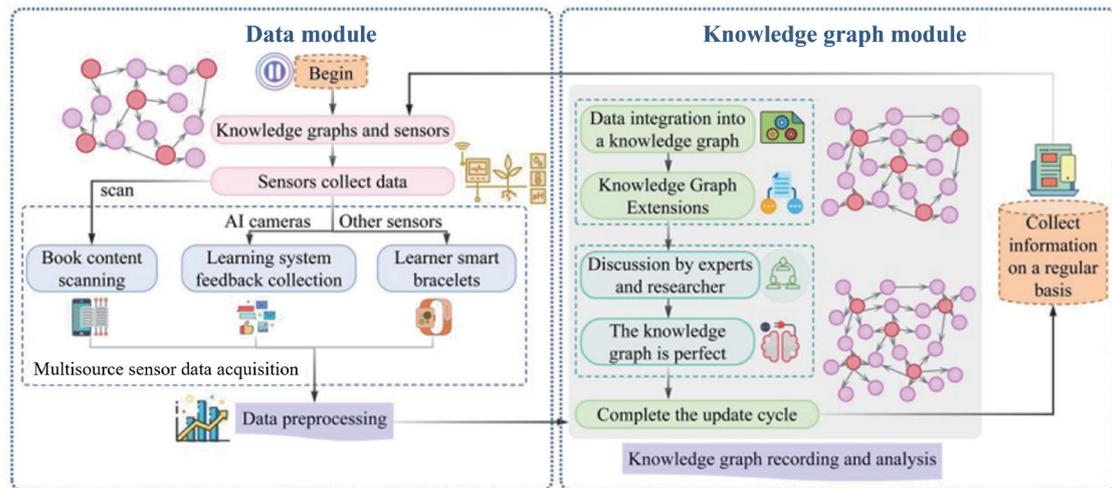


Fig. 3. (Color online) Dual-module system architecture developed in this study.

3.2 Data collection

As physiological and facial data were collected and used in this study, measures were implemented to ensure participant privacy and ethical compliance. All learners participated voluntarily and signed informed consent detailing the types of sensor used and the purpose of data collection. Facial images captured during classes were processed into 128-dimensional numerical feature vectors using multitask cascaded convolutional networks (MTCNNs) and residual network (ResNet)-based algorithms. The raw facial images were not stored or transmitted to the cloud, ensuring the participants' privacy. The data transmitted between the perception module and eXtreme Gradient Boosting (XGBoost) were protected using advanced encryption standard-256 to prevent unauthorized access. The sensor data collected in this study are presented in Fig. 4.

Body temperature was measured using a smart wristband that adopted a thermistor and infrared radiation. A combination of algorithms, specifically Kalman filtering and moving average smoothing, was used to estimate body temperature from the measured skin surface. The relationship between body temperature and learning state is described in Table 1.

An AI camera equipped with image sensors and time-of-flight (TOF) depth sensors was used to acquire data of the learner's book content scanning. Book content scanning is an important process of information perception and transformation. For scanning, image sensors were used to obtain facial images and capture classroom scenes. The image data collected were compared with preregistered learners' faces using a face recognition algorithm to determine their presence and timestamp to monitor the attendance. Hand raising, speaking, participating in discussions, nodding, and shaking the head were captured in classroom scenes for the continuous monitoring of learners' postures and attitudes. An extraction algorithm was used to identify specific actions (e.g., hand raising, posture changes, and gestures), and the results were stored on the server. The system uses a combination of face recognition (for attendance) and pose estimation algorithms

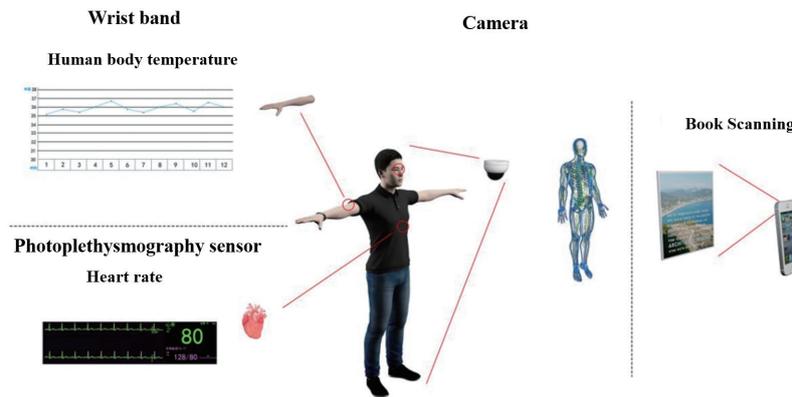


Fig. 4. (Color online) Sensors used in this study.

Table 1
Temperature–learning state mapping table.

Learning state	Temperature range	Response
Normal focus	Within normal range	Maintain current learning content
Nervous/anxiety	Slightly elevated	Provide stress relief or easier content
Fatigued/unfocused	Slightly below baseline	Suggest rest or engaging content
Abnormal	Significantly high/low	Recommend health check or pause study

(to identify actions such as hand raising and nodding). These visual features are converted into quantitative interaction frequencies used as inputs for the XGBoost predictive model (Fig. 5).

For face recognition and attendance monitoring, we utilized MTCNN for face detection combined with a ResNet-based feature extractor for identification.⁽²⁷⁾ For the extraction of linguistic entities from scanned text, we employed the TextRank algorithm for the identification of keywords based on graph-based centrality.⁽²⁸⁾

The photoplethysmography (PPG)-based heart rate sensor was used for the real-time monitoring of learners' heart rates. The data collected reveal a learner's level of concentration, tension, fatigue, and cognitive load. The working principle of PPG is illustrated in Fig. 6. The data on PPG are used to monitor variations in learners' cognitive load.⁽²⁹⁾ We collected learning behavior data based on heart rate data and established a correlation model among physiological indicators, behavioral performance, and cognitive states (Table 2). We also obtained data on the learning process and outcome from a Chinese language learning platform. The data and tools used in the study are listed in Table 3.

Scanned text images were integrated with physiological data through multimodal synchronization. While the image and TOF depth sensors were used to identify the knowledge entity or text segment the learner engaged with, PPG and temperature sensors were used to obtain concurrent biological feedback.⁽³⁰⁾ For instance, if a learner scanned a complex sentence structure and the heart rate sensor detected a significant increase in cognitive load, defined by heart rate variability, a specific linguistic entity was tagged as a high-difficulty node. These synchronized data points were fed into the XGBoost model to evaluate whether the corresponding knowledge graph entity requires reinforcement or if the pedagogical issues

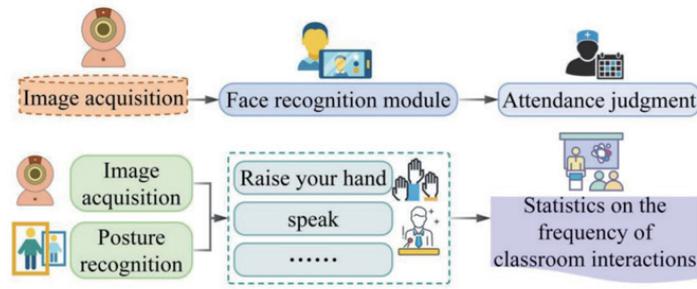


Fig. 5. (Color online) Workflow for behavioral data acquisition via AI camera.

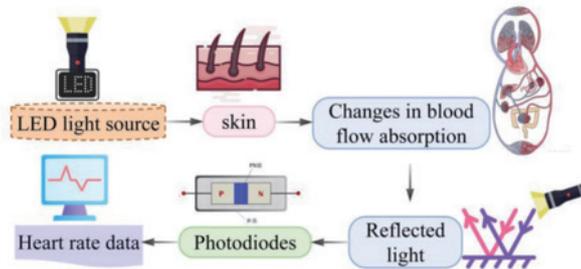


Fig. 6. (Color online) Principle of photoplethysmography.

Table 2
Heart rate change and learning status.

Learning status	Heart rate variability	Response
Normal focus	Stable, moderate fluctuations	Maintain current learning content
Nervous/anxiety	Sharp increase/Higher than individual baseline	Provide stress relief or easier content
Fatigue/unfocused	Slow decrease, low variance	Suggest rest or engaging content
Abnormal	Significantly high/low	Recommend health check or pause study

Table 3
Indicator system overview.

Application	Tool	Data	Description
Book content scanning	Image sensor	Text and image content	Chinese learning resource construction
Behavioral data	AI camera	Attendance rate, classroom interaction frequency	Informing learner status to instructors
Learning process data	Learning system platform	Time distribution, learning path, usage	Tracking resource usage
Learning outcome data	Learning system platform	Exam scores, project outcomes, skill improvement	Evaluating learning effectiveness
Physiological and psychological data	Body temperature, heart rate sensor	Attention level, emotional state, and physiological indicators	Monitoring learners' mental and physical conditions

needed updating to reflect learner demand. This enables the evaluation system to account for the learner’s internal cognitive state in real time beyond measuring performance metrics.

3.3 Data analysis

After the sensor collected data on learners' physiological and psychological status and activities, key features were extracted by analyzing the data to examine the learning status of learners. To evaluate learning outcomes effectively, we extracted essential features and compared the performance of decision trees, regression models, neural networks, support vector machines, and XGBoost, using averaging and voting to select the optimal model for this study. On the basis of the results, the XGBoost algorithm was selected as the base model, as it has been widely used for classification and regression tasks, and the following functions were well suited for the data collected:

- XGBoost regularized the loss function to avoid overfitting and outperformed other algorithms leveraging gradient boosted decision trees (GBDT).
- Even when missing values were imputed, XGBoost automatically improved the efficiency of the algorithm.
- XGBoost used cross-validation in each iteration to easily obtain the optimal boosting iteration number.
- XGBoost improved the interpretability of models compared with neural networks.

The main process of the XGBoost algorithm comprises the following steps:

- 1) Step 1: Mathematical definition and objective function

$$\hat{y}_i = \sum_{k=1}^K f_k(x_i) \quad (1)$$

$$Obj(\theta) = \sum_{i=1}^n L(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k) \quad (2)$$

Here, n is the number of data points, x_i represents the i th sample, y_i and \hat{y}_i correspond to the true and predicted values of the i th sample, respectively, K represents the number of classification and regression trees (CARTs), f_k denotes the k th CART, $L(y_i, \hat{y}_i)$ is the loss function, and $\Omega(f_k)$ is the regularization term. XGBoost trains the t th tree during each split. The objective function for the t th tree is formulated as

$$Obj^{(t)} = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t)}) + \sum_{k=1}^t \Omega(f_k) = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)}) + f_t(x_i) + \Omega(f_t) + C. \quad (3)$$

- 2) Step 2

The XGBoost loss function utilizes a second-order Taylor expansion to improve approximation accuracy, which is a different property from GBDT. Therefore, the objective function is expressed as

$$Obj^{(t)} \approx \sum_{i=1}^n \left[L(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t) + C. \quad (4)$$

Here, g_i is the first-order derivative of the loss function L with respect to the second variable at the i th sample point and h_i is the second-order derivative of the loss function L with respect to the second variable at the i th sample point. The regularization part $\Omega(f_t)$ is defined as

$$\Omega(f_t) = \gamma T + \frac{1}{2} \partial \sum_{j=1}^T \omega_j^2. \quad (5)$$

Here, f_t represents the t th tree, T is the number of leaf nodes in the tree, ω_j is the score at the j th leaf node, and γ and ∂ are the penalty factors.

3) Step 3

The decision tree is split into two parts: the structural part q and the weight part ω .

$$f_t(x) = \omega_{q(x)}, \omega \in R^T, q: R^d \rightarrow \{1, 2, \dots, T\} \quad (6)$$

Here, ω is a T -dimensional vector, corresponding to the scores at each leaf node, and q is a mapping that maps a sample point to a specific leaf node. Once q and ω are determined, the structure of the tree is fully defined.

4) Step 4

By substituting the expressions for f_t and Ω into the approximated objective function, we obtain the following equation with Obj^t :

$$Obj^t \approx \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t), \quad (7)$$

$$= \sum_{i=1}^n \left[\left(\sum_{i \in I_j} g_i \right) \omega_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \partial \right) \omega_j^2 \right] + \gamma T. \quad (8)$$

The minimum point and the minimum value of the objective function can be directly obtained from the final expression [Eqs. (9) and (10)].

$$\omega_j^* = -\frac{G_j}{H_j + \partial} \quad (9)$$

$$Obj^* = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \partial} + \gamma T \quad (10)$$

5) Step 5

A greedy algorithm is used to construct the tree structure and select the tree with the smallest structure score (Obj^*). In this study, by comparing the prediction performance on multiple datasets, the model was further optimized for feature and parameter selection and accuracy and stability improvement of the learning status.

3.4 Knowledge graph update

Sensor data were continuously collected and analyzed for the model evaluation, particularly formative assessment using the knowledge graph. Formative assessment was conducted to evaluate learners' learning status and measure the effectiveness of learning resources. From the results, learning resources were constructed on the basis of sensor data to infer and assess learners' learning levels and adaptability of learning resources, and provide feedback accordingly. A knowledge graph was constructed in the following process:

1) Step 1: Candidate pattern mining

Knowledge graphs contain a vast number of nodes and edges, making pattern recognition and subgraph coverage computation time-consuming. Therefore, we adopted a Giraph-based graph simulation method to extract a candidate pattern set and to create a base knowledge graph.

2) Step 2: Standard pattern mining (SPM)

The SPM algorithm was used with the input of the knowledge graph $G(V,E,L)$ (Fig. 7), where the number of edges is limited to k , the number of nodes is limited to n , and the initial result

```

 $S_{p_f} \leftarrow GRAMI(G)$ 
if ( $k_p \leq k, n_p \leq n$ ) then
    Giraph( $G, S_{p_f}$ );  $R(v_G, v_P), R(l_G, l_P)$ ;
     $S_p \leftarrow S_p \cup \{P\}$ ;
end if
 $S \leftarrow \emptyset$ ;
for  $|S| \leq m$ 
     $\{\sigma(S|P) = F(S \cup \{P\}) - F(\{P\})\}$ ;
     $F(S) = (|V_{G_S}| + |E_{G_S}|) / |G|$ ;
     $F(S \cup \{P\}) = (|V_{G_S+G_P}| + |E_{G_S+G_P}|) / |P|$ ;
    select  $P, P'$ 
    if  $\sigma(S|P') > \sigma(S|P)$ 
         $P^* \leftarrow \operatorname{argmax}_{P \in S_p} \sigma(S|P')$ ;
        else  $P^* \leftarrow \operatorname{argmax}_{P \in S_p} \sigma(S|P)$ ;
         $S \leftarrow S \cup \{P^*\}$ ;
    end if
end for
return S

```

Fig. 7. SPM algorithm used in this study.

set $\emptyset \rightarrow I$, with the number of standard pattern results m . The output was the standard pattern set S .

3) Step 3: Graph simulation filtering

The mined standard pattern set was used to perform graph simulation and matching with the updated knowledge graph. The matching process is similar to that in identifying standard patterns. A binary relationship $R(P, G) \in V_P \times V$ between the standard pattern set and the knowledge graph was identified and updated. First, on the basis of the labels of nodes in the updated knowledge graph, a set of candidate matches was generated within the pattern set (a pure market knowledge graph). The knowledge graph node must be a successor so that a successor relationship between the node in the knowledge graph and the corresponding node in the pattern graph can be identified. In other words, the node value in the knowledge graph must be the same as the successor node value in the pattern graph. On the basis of this matching requirement, the unmatched nodes were filtered out. The remaining subgraph was considered to present a valid entity embedded in the graph. The input of the graph simulation was the standard pattern set $S=(P_1, P_2, \dots, P_m)$, whereas the output was the knowledge graph to be updated to G^* (Fig. 8).

4) Step 4: Vectorization filtering

The entity vector distances were computed in different relationships to update the knowledge graph. The results were verified through a group discussion of experts and researchers to ensure the accuracy of the results. After completing the filtering, the final version of the knowledge graph (expert-fused knowledge graph) was obtained. The updated knowledge graph was inserted into the original knowledge graph using the graph embedding or triple embedding method to construct a new base graph. On the basis of learner feedback data, the base graph was evaluated to enable the iterative optimization and structural updates of entities. The data collection–graph optimization–learning validation–reoptimization closed-loop mechanism in this study was used to construct learning resources for Chinese language learning.

```

Giraph(  $G_{new}, S$  ) :
   $R(v_{G_{new}}, v_P), R(l_{G_{new}}, l_P)$  :
    if  $l(v_{G_{new}}) \in l(v_P)$  and  $v_{G_{new}}$ 
      | then match;
    end if
    if  $(l'_1, l'_2, \dots, l'_c) \in l(v_P), (l'_1, l'_2, \dots, l'_c) \in l(v_{G_{new}})$ 
      | then match;
    end if
     $G^* \leftarrow G_{P_i}$ ;
  end
end

```

Fig. 8. Algorithm to obtain input and output in graph simulation.

3.5 Participants and expert group

We obtained the data from the “I CAN Chinese” learning platform, a podcast hosted by Michael Castelein, designed to make Mandarin learning simple and fun, and its offline learning records as shown in Tables 2 and 3. The data included heart rate variability, skin temperature, interaction frequencies (number of hand-raising events, nodding, and eye-contact duration), and spatial coordinates. The frequency of HSK-graded vocabulary and the complexity of grammatical structures were extracted from scanned text images.

The participants were foreign Chinese language learners recruited from the “I CAN Chinese” learning platform and from offline courses at 10 universities in Yunnan Province, China, where volunteers agreed to wear smart wristbands. They were eligible learners as they enrolled in HSK Level 3–5 courses between January 2024 and March 2025 with complete learning records. Learners with incomplete records or without valid activity for three consecutive months were excluded. 1,108 participants representing 11 nationalities were selected in this study, with the largest groups from Thailand (30.78%), India (19.31%), Indonesia (12.64%), Myanmar (11.19%), and Laos (7.49%). The other countries of the participants included the United States, Japan, Malaysia, Russia, and the United Kingdom. Among them, 734 were undergraduates, 345 were master’s students, and 29 were doctoral candidates. Their learning and physiological data were integrated from both the platform and wristband sensors. Using the XGBoost algorithm, we identified and incorporated their knowledge entities into adaptive knowledge graphs, which were refined to enhance personalization and improve learning efficiency.

The expert group comprised 10 specialists in international Chinese education and second language acquisition, educational technology and learning analytics, machine learning and data mining, and intelligent sensing and wearable devices. They were affiliated with the university’s School of International Chinese Education, School of Information, and the platform’s technical team from Yunnan Normal University, China. They had more than 15 years of experience in research, teaching, or engineering within their respective fields, demonstrating academic or professional competence with an associate senior title or above, a doctoral degree, and provincial or ministerial project management, or academic article publication. Before participating in this study, all experts were informed about the study objectives and their roles. They also consented to the requirements for data usage and privacy.

3.6 XGBoost-algorithm-based model

After preprocessing the data collected, the XGBoost algorithm was used to train the XGBoost model. K-fold cross-validation was used to divide the dataset into training and testing datasets. The whole dataset was partitioned into k subsets that were mutually exclusive and equally sized. For each round, $k-1$ subsets were included in the training dataset, while the remaining one was included in the testing dataset. This process was repeated k times until the mean of the k test results was obtained.

To assess the model’s generalization capability, the dataset was divided into training and test sets in an 8:2 ratio, including 887 participants in the training set and 221 participants in the test

set. The training set was used to optimize the model is parameters, whereas the test set was used to evaluate its performance. The training dataset was used to determine the model's parameters. Parameters of the XGBoost model included learning rate (*learning_rate*), the number of iterations (*n_estimators*), maximum tree depth (*max_depth*), and regularization coefficient (*reg_lambda*). The testing dataset was used to evaluate the performance of the XGBoost model.

3.6.1 Model input and output

The XGBoost model performs exceptionally well on high-dimensional data with minimized overfitting, which are caused by numerous features. However, if the features are highly correlated, the model prioritizes mutually correlated features for information gain, potentially overlooking other important variables. In this study, the XGBoost model was constructed using CARTs as weak learners. Three types of parameter were defined in the model on the basis of a literature review: weak learner parameters based on CART, ensemble learning framework parameters, and other parameters in the process. The key parameters of the XGBoost model and corresponding tuning guidelines were selected as shown in Table 4.

The expert group fine-tuned key parameters of the XGBoost model to balance convergence speed, stability, and generalization. The learning rate, maximum tree depth, regularization coefficients, subsample ratio, and minimum child weight were prioritized because they affect prediction accuracy, computational efficiency, and model interpretability. The expert group also considered parameter values proposed in previous studies, considering the trade-off between model complexity and generalization, and available computational resources.

Parameter fine-tuning was conducted using grid search with fivefold cross-validation. In each iteration, one fold served as the validation set while the remaining folds were used for training. The initial model training on the training set employed the following configurations: the booster was set to *gbtree*, the learning rate to 0.1, and gamma retained its default value of 0.

Table 4
Fine-tuned parameters in this study.

Parameter	Impact on model performance	Value
<i>booster</i>	Specifying the type of weak learner with options including tree-based (<i>gbtree</i>) and linear (<i>gblinear</i>) models	<i>gbtree</i>
<i>learning_rate</i>	A smaller learning rate requires weaker learners to compensate for residuals.	[0.01, 0.015, 0.025, 0.1]
<i>gamma</i>	The larger the gamma value, the more conservative the model becomes.	[0, 0.1, 0.3, 0.5, 0.7, 0.9, 1]
<i>reg_alpha</i>	The larger the <i>reg_alpha</i> value, the more conservative the model becomes.	[0, 0.01–0.1, 1]
<i>reg_lambda</i>	The larger the <i>reg_lambda</i> value, the more conservative the model becomes.	[0, 0.1, 0.5, 1]
<i>max_depth</i>	Preventing overfitting	[3, 5, 6, 7, 9, 12, 15, 17]
<i>min_child_weight</i>	If the sum of instance weights in a leaf is less than <i>min_child_weight</i> , the split stops.	[1, 3, 5, 7]
<i>subsample</i>	Values less than 1 are used to randomly sample a portion, preventing overfitting.	[0.6, 0.7, 0.8, 0.9, 1]

Both the L1 regularization term (reg_alpha) and the L2 regularization term (reg_lambda) were set to 1. Additionally, max_depth was set to 6, min_child_weight to 3, and subsample to 0.9.

3.6.2 Performance evaluation

In this study, mean squared error (MSE), root mean squared error ($RMSE$), mean absolute error (MAE), and R-squared (R^2) were selected for the evaluation of the model performance. For a good fit for the model, parameters need to be optimized on the basis of the results of the initial training. Common parameter optimization methods include manual tuning, grid search, random search, Bayesian optimization, and genetic-algorithm-based tuning. We used grid search for fivefold cross-validation to optimize the parameters of the XGBoost model.

For fivefold cross-validation, we split the data into five groups and selected one group as the validation data while using the rest as the training data in each iteration. The MSE of each validation set was computed, and the average MSE from the five validation sessions was calculated. The average MSE s for different parameter sets were also computed, and the parameter set with the highest average was chosen as the optimal model parameter set. Grid search involves iterating the model operation with all possible combinations of parameters to find the optimal parameter set in a specified range. For the XGBoost model of this study, parameters were optimized as presented in Table 5. Grid search was conducted using the sklearn standard library. The method facilitates the evaluation of learners' learning outcomes and supports the appropriate allocation of learning resources based on the evaluation results.

The performance of the XGBoost model with the optimized parameters is shown in Table 6. MSE decreased from 0.0609 (before training) to 0.0565 (after training), indicating an improvement in the model's prediction accuracy. $RMSE$ was more sensitive to prediction errors. Before and after training, $RMSE$ decreased from 0.2469 to 0.2378. R^2 increased from 0.5749 to 0.6057, suggesting that the model better explained the relationship between learner behavior and learning outcomes.

Table 5
Optimized parameters for XGBoost model developed in this study.

Parameter	Optimized value
learning_rate	0.1
gamma	0.0
max_depth	8.0
min_child_weight	1.0
subsample	0.5

Table 6
Metrics of model performance.

	MSE	$RMSE$	MAE	R^2
Before tuning	0.0609	0.2469	0.1150	0.5749
After tuning	0.0565	0.2378	0.1028	0.6057

4. Results and Discussion

4.1 Knowledge graph construction

On the basis of the XGBoost prediction results and sensor data, the knowledge graph was constructed using MLN or Bayesian optimization. The probability distribution of the knowledge graph was calculated using MLN [Eq. (11)].

$$P(X) = \frac{1}{Z} \prod_{c \in C} \psi_c(X_c) \quad (11)$$

Here, $P(X)$ is the joint probability distribution of the entire set of nodes X , Z is the partition function to ensure the probability sum to 1, being calculated by summing the product of probabilities of all random variables, and $\psi_c(X_c)$ is the non-negative function that measures the compatibility or goodness of the nodes.

The construction of the knowledge graph began with the prediction results generated by the XGBoost model and sensor data. These inputs were used to automatically build the initial knowledge graph [Fig. 9(a)]. Next, the probability distribution of the graph was calculated using MLN combined with Bayesian optimization to establish a base knowledge graph [Fig. 9(b)].

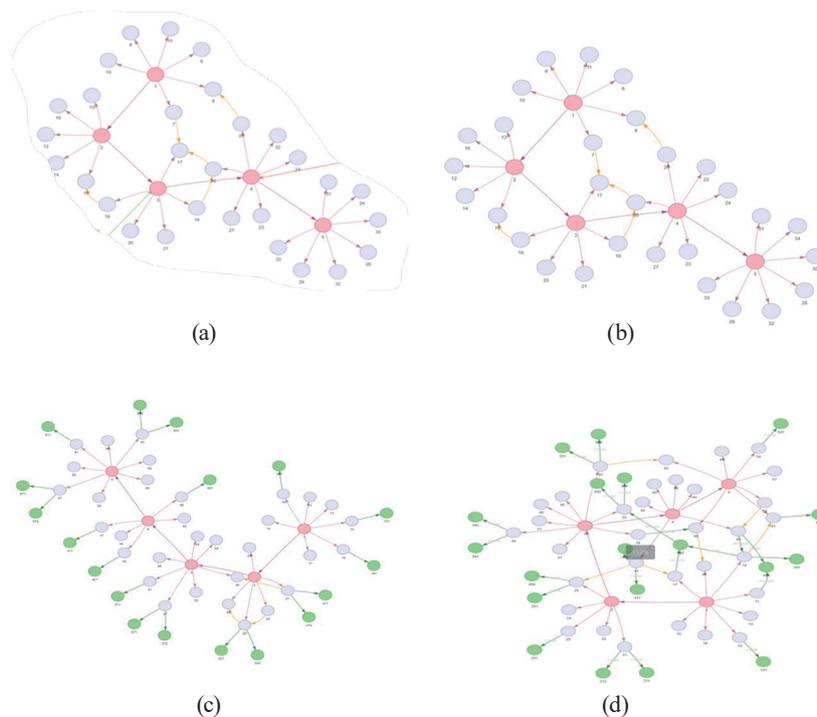


Fig. 9. (Color online) Stages of knowledge graph evolution (pink circle: key nodes, green circle: added nodes). (a) Initial knowledge graph, (b) base knowledge graph, (c) model knowledge graph, and (d) knowledge graph modified with expert opinions.

Graph algorithm was then applied to extract patterns among nodes (entities), which were refined through vectorized filtering to identify and correct mismatched nodes. This process produced the model knowledge graph. Here, newly added nodes (green circles) represent emerging linguistic elements such as internet slang [Fig. 9(c)]. Finally, a manual review was conducted to incorporate expert feedback, resulting in the expert-fused final knowledge graph [Fig. 9(d)]. The expert group reviewed and refined the nodes, leveraging their knowledge to ensure the accuracy and soundness of the knowledge graph. Then, the final knowledge graph was constructed by integrating both automated insights and human expertise.

The knowledge graph constructed in this study was used to identify and add new nodes corresponding to the learner's needs. This adaptive capability of the knowledge graph enabled an accurate and relevant learning resource construction. The analysis results using the knowledge graph constructed highlighted the need for incorporating contemporary cultural content into the learning resources. New nodes presented emerging slang popularly used on the internet. The inclusion of new nodes for puns and homophones, such as '蚌埠住了' (meaning cannot hold it or cannot take it anymore in the context of emotions like laughter, crying, or anger) and '芭比Q了' (it is over or I am doomed, expressing a situation that has turned bad or hopeless) shows the importance of updating learning resources with Chinese used in daily life. This also suggests that teaching materials must be updated regularly to reflect the rapid development of internet culture, which is highly relevant to international learners. The added nodes also included knowledge on Chinese festivals, overseas Chinese customs, and cross-cultural taboo comparisons. For example, the new node for durian mooncakes in Southeast Asia highlights how Chinese culinary traditions are adapted and received differently globally. This indicates that cultural content must be added in language learning, and regional differences in Chinese culture must be considered in constructing or updating learning resources. New nodes also indicated the importance of business Chinese language and regional dialects. As communication in Chinese has been influenced by dialects, such as Cantonese or Wu-speaking regions, communication habits from specific dialects must be added to learning resources.

The knowledge graph constructed for the development of learning resources responds to changes in learner needs and the evolving nature of the language. By integrating new data from sensors, expert opinions, and learner feedback, the knowledge graph must be continually refined to provide relevant and up-to-date teaching resources.

4.2 Performance of XGBoost model

The data used in the experiment consisted of learners' behavior data collected from sensors (e.g., heart rate, body temperature, and reading duration) and from a Chinese language learning platform (e.g., knowledge point frequency and user ratings). User rating, recency (last update time), study frequency, relevance, and average reading duration were important features extracted by the XGBoost algorithm to decide on the necessity for entity update. The effect of each feature on the prediction was analyzed to support the decisions for learning resource updates. The model was trained and evaluated for its performance using the training and testing datasets, which contained 80 and 20% of the whole data, respectively. The parameters, such as

learning rate, tree depth, and regularization coefficients, were optimized using a grid search through fivefold cross-validation. The model showed a prediction accuracy of 89.2%, an F1 score of 0.87, and an area under the receiver operating characteristic curve of 0.93 on the testing dataset, confirming the accuracy of the model in predicting the necessity for entity updates.

Figure 10 shows the distribution of prediction probabilities of update needs for entities. The probability of no entity updates (pink) was concentrated near zero, whereas that of entities requiring updates (blue) ranged from 0.5 to 1.0 with multiple peaks, suggesting varying levels of update necessity. These results demonstrated that the XGBoost model effectively distinguished updating necessity, offering a probability-based decision for entity updates of the teaching resource. This enables the prioritization of knowledge that needs to be updated or revised.

Figure 11 presents the distributions of the study frequency and average reading duration of entities in five categories. The numbers (1–5) indicated at the top of each box correspond to HSK Levels 1 through 5, representing a progression from elementary to advanced Chinese language proficiency. Level 1 represents basic survival Chinese (approximately 150 words), whereas Level 5 represents the ability to read newspapers and deliver professional speeches (approximately 2500 words). As learners progress from Level 1 to Level 5, the average reading duration per session increases significantly. This trend, captured by the TOF and image sensors, validates the system’s ability to distinguish between high-frequency, low-duration ‘skimming’ typical of simpler levels and the deep cognitive engagement required for higher-level linguistic nodes. The data show how reading stamina and frequency adapt as the linguistic complexity of the knowledge graph entities increases. While the overlap of updating and sustaining entities was observed, entities with low study frequency and short reading duration were distinguished. The results indicate that study frequency and reading duration can serve as indicators for entity

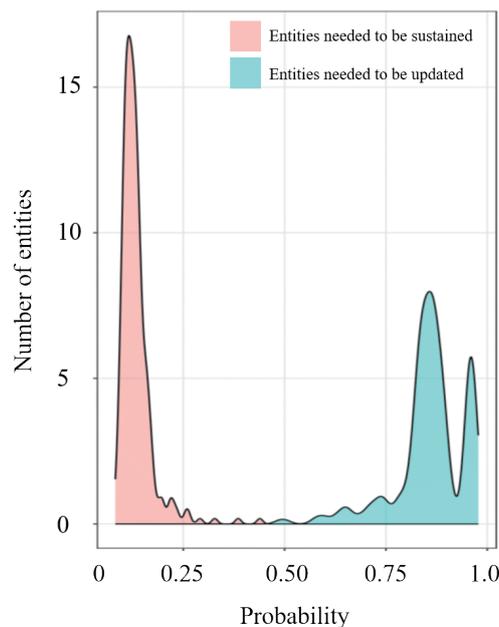


Fig. 10. (Color online) Prediction probability of entity updates.

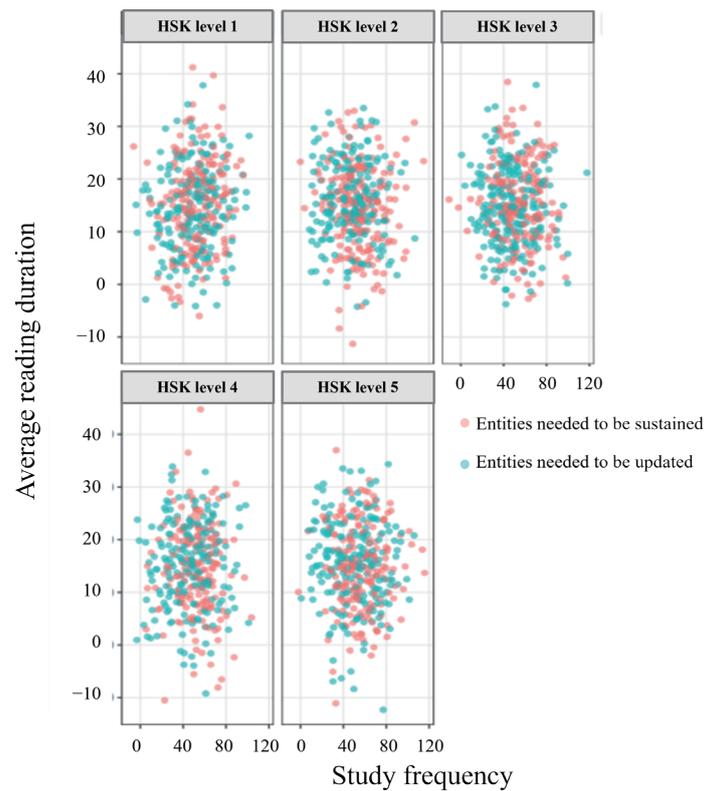


Fig. 11. (Color online) Correlation between study frequency and average reading duration across different proficiency levels.

updates. The XGBoost model effectively identified relationships between features and entities that need to be updated.

User rating and recency showed prominent average improvement rates (0.603 and 0.226), indicating that they were the most influential features. When using user ratings and recency to decide on entity updates, errors can be minimized. They are also important indicators of the perceived quality of a teaching resource and a measure of a resource's effectiveness. Study frequency, relevance, average reading time, and category showed lower average improvement rates. While they contributed to the model's predictions, their impact was not significant. The lowest rate for HSK level (0.00258) implies that it was not an effective predictor for entity update (Table 7).

The results presented the importance of user rating systems and the frequent update of entities in Chinese language learning resources. A good rating system is a powerful tool for identifying and promoting effective learning resources. Therefore, resources with high ratings must be continuously developed. At the same time, resource libraries must be frequently updated. By creating and promoting new content, a satisfactory learning experience can be maintained through learning resources. The HSK level did not significantly affect the entity update necessities. To construct an effective teaching resource, user-driven metrics such as user ratings and recency must be considered to enhance the interaction of learners with learning resources.

Table 7
Feature importance analysis.

Feature	Average improvement ratio in accuracy	Rate of observation	Weight
User rating	0.603	0.163	0.101
Recency	0.226	0.202	0.154
Study frequency	0.0824	0.283	0.281
Relevance	0.0658	0.199	0.2222
Average reading time	0.0202	0.140	0.202
HSK level	0.00258	0.0122	0.0404

5. Conclusion

We developed an adaptive and data-driven framework for constructing and optimizing Chinese language learning resources for international learners. By integrating sensor technology and knowledge graphs, we addressed the limitations of static resources and a disconnect between resource supply and learner demands. By collecting real-time physiological and behavioral data, we obtained a quantitative, high-resolution understanding of learners' cognitive and emotional status. The collected qualitative and quantitative data were used to develop and evaluate an XGBoost model that accurately predicted the necessity of entity updates with an accuracy of 89.2%. The model's feature importance analysis result highlighted that user ratings and resource recency were the most important indicators for a learning resource's effectiveness, which is significant for resource design and management. The model's prediction results were used to refine and update the knowledge graph. This iterative process resulted in an adaptively adjusted knowledge graph by incorporating contemporary and culturally relevant content, including internet slang and regional business communication patterns. The knowledge graph created on the basis of the expert group's opinions is an adaptable tool that responds to the evolving nature of the Chinese language and the specific needs of diverse learners. The method using sensor data and knowledge graphs proved effective in updating or constructing educational resources. The closed-loop system of perception–analysis–update based on sensor data and analytics also contributes to the development of robust, personalized, and efficient learning resources with attractive learning materials.

Although the participants from 11 different countries were included in this study, the majority originated from South and Southeast Asia. Therefore, the data might not represent the behavioral and physiological nuances of learners from other countries, who might exhibit different classroom interaction patterns or physiological baselines. However, the sensor-driven knowledge graph can be applied to similar research on the learners of diverse linguistic and cultural backgrounds. In future research, the dataset needs to be expanded to validate the cross-cultural robustness of the XGBoost predictive model.

References

- 1 J. Guo, Y. H. Wu, L. Gu, L. Zhou, F. Nong, J. N. Ma, J. X. Cui, and X. Y. Dong: *Int. Chin. Teach. Res.* **11** (2021) 86. https://www.zhangqiaokeyan.com/academic-journal-cn_journal-international-chinese-teaching-thesis/0201291822513.html
- 2 J. F. Ma, Y. Liang, Y. H. Wu, and J. N. Ma: *J. Tianjin Norm. Univ. (Soc. Sci. Edn.)* **8** (2021) 15. https://www.zhangqiaokeyan.com/academic-journal-cn_journal-tianjin-normal-university-social-science-edition-thesis/0201291512825.html.
- 3 X. Li, D. Sun, and J. Qiu: *Proc. 2024 Int. Conf. Intelligent Education and Computer Technology (IWCT, 2024)* 134. <https://doi.org/10.1145/3687311.36873361>
- 4 Z. G. Ou, Y. P. Liu, K. Qian, Y. Wang, and X. Xie. Li: *Mod. Educ. Technol.* **11** (2024) 37. <https://doi.org/10.3969/j.issn.1009-8097.2024.09.004>
- 5 Chinese Tests Service Website: <https://www.chinesetest.cn/hsk> (accessed February 2026).
- 6 J. Xu and R. L. Ma: *E-Educ. Res.* **44** (2023) 121. <https://doi.org/10.13811/j.cnki.eer.2023.10.016>
- 7 J. H. Song, M. Zhang, and L. F. Liang: *Int. Chin. Teach. Res.* **9** (2023) 14. <https://doi.org/10.3969/j.issn.2095-798X.2023.03.002>
- 8 J. M. Lu: *TCSOL Stud.* **4** (2023) 25. <https://doi.org/10.16131/j.cnki.cn44-1669/g4.2023.04.002>
- 9 S. Ji, S. Pan, E. Cambria, P. Marttinen, and P. S. Yu: *IEEE Trans. Neural Netw. Learn. Syst.* **33** (2021) 494. <https://doi.org/10.1109/tnnls.2021.3070843>
- 10 Y. Sun, Y. Qin, W. Chen, X. Li, and C. Li: *Appl. Sci.* **15** (2025) 7068. <https://doi.org/10.3390/app15137068>
- 11 V. Radosavljevic, S. Radosavljevic, and G. Jelic: *Interact. Learn. Environ.* **30** (2022) 307. <https://doi.org/10.1080/10494820.2019.1652836>
- 12 S. Dai: *Int. J. Emerg. Technol. Learn.* **14** (2019) 38. <https://doi.org/10.3991/ijet.v14i03.10104>
- 13 J. Rong: *Mob. Inf. Syst.* **2022** (2022) 9297314. <https://doi.org/10.1155/2022/9297314>
- 14 J. Sweller, J. J. G. Van Merriënboer, and F. G. W. C. Paas: *Educ. Psychol. Rev.* **10** (1998) 251. <https://doi.org/10.1023/A:1022193728205>
- 15 D. Q. Nguyen: arXiv preprint arXiv:1703.08098 (2017). <https://doi.org/10.48550/arXiv.1703.08098>
- 16 J. Liang, S. Zhang, and Y. Xiao: *Proc. 2017 26th Int. Joint Conf. Artificial Intelligence (AAAI, 2017)* 3749. <https://dl.acm.org/doi/abs/10.5555/3304889>
- 17 D. Song, J. Xu, J. Pang, and H. Huang: *Inf. Sci.* **573** (2021) 222. <https://doi.org/10.1016/j.ins.2021.05.045>
- 18 S. Y. Yao, T. Z. Zhao, R. J. Wang, and J. Liu: *J. Comput. Res. Dev.* **57** (2020) 2514. <https://doi.org/10.7544/issn1000-1239.2020.20200741>
- 19 S. Iwata, S.-I. Tanigawa, and Y. Yoshida: *Proc. 2016 27th Annual ACM-SIAM Symp. Discrete Algorithms (SIAM, 2016)* 404. <https://epubs.siam.org/doi/abs/10.1137/1.9781611974331.ch30>
- 20 W. Dvořák, M. Henzinger, and D. P. Williamson: *Algorithmica* **77** (2017) 152. <https://doi.org/10.1007/s00453-015-0066-y>
- 21 M. Richardson and P. Domingos: *Mach. Learn.* **62** (2006) 107. <https://doi.org/10.1007/s10994-006-5833-1>
- 22 J. Pujara, H. Miao, L. Getoor, and W. Cohen: *Proc. 2013 12th Int. Semantic Web Conf. 1st Australasian Semantic Web Conf. (ISWC, 2013)* 542. https://link.springer.com/chapter/10.1007/978-3-642-41335-3_34
- 23 O. Kuželka and J. Davis: *Uncertainty in Artificial Intelligence. PMLR* **115** (2020) 1138. <https://proceedings.mlr.press/v115/kuzelka20a.html>
- 24 C. Halaschek-Wiener, B. Parsia, E. Sirin, and A. Kalyanpur: *Proc. 2006 Int. Workshop Description Logics DL'06 (DL, 2006)* 200. <http://ftp.informatik.rwth-aachen.de/Publications/CEUR-WS/Vol-189/DL2006.pdf#page=208>
- 25 D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati: *Proc. 2006 19th Int. Workshop Description Logics (DL, 2006)* 51. <https://bia.unibz.it/esploro/outputs/conferenceProceeding/Epistemic-first-order-queries-over-description-logic/991006761498001241>
- 26 Y. Li, B. Xu, J. Lu, and D. Kang: *Proc. 2006 19th Int. Workshop Description Logics DL'06 (DL, 2006)* 39. <https://dlwqtxts1xzle7.cloudfront.net/8185411/10.1.1.142.9330-libre>
- 27 K. Zhang, Z. Zhang, Z. Li, and Y. Qiao: *IEEE Signal Process. Lett.* **23** (2016) 1499. <https://doi.org/10.1109/LSP.2016.2603342>
- 28 R. Mihalcea and P. Tarau: *Proc. 2004 Conf. Empirical Methods in Natural Language Processing (EMNLP, 2004)* 404. <https://aclanthology.org/W04-3252.pdf>
- 29 F. G. W. C. Paas, J. J. G. Van Merriënboer: *Educ. Psychol. Rev.* **6** (1994) 351. <https://doi.org/10.1007/BF02213420>.
- 30 Y. Hasanpoor, B. Tarvirdizadeh, K. Alipour, and M. Ghamari: *Signal, Image Video Process.* **19** (2025) 1129. <https://doi.org/10.1007/s11760-025-04734-z>

About the Authors



Yue Fu received her master's degree in digital technology, communication, and education from the University of Manchester, the United Kingdom. Currently, she is a doctoral candidate at the School of International Chinese Language Education, Yunnan Normal University, China. She has been engaged in the research of digitalization in international Chinese language education and educational technology. (fuyuel@ynnu.edu.cn)



Lei Zhao obtained his bachelor's degree in information management and information systems from Nanjing University of Posts and Telecommunications, China, in 2018. In 2023, he obtained a master's degree in machine learning and data mining from Jiangxi University of Finance and Economics, China. His research interests include machine learning and sensor technology. (lecile163@163.com)



Borui Zheng received her bachelor's degree in communication engineering and master's degree in Teaching Chinese to Speakers of Other Languages (TCSOL) from Lanzhou University, China. From 2016 to 2025, she worked as a language lecturer in Georgia, Poland, and Shanghai. She is currently a Chinese language instructor at University College Cork (UCC), Ireland. Her research interests include language pedagogy, second language acquisition, educational data analysis, and cross-cultural communication. (bzheng@ucc.ie)



Yirong Wang obtained a bachelor's degree in computer science and technology from Hangzhou Normal University in China in 2020. Since 2020, she has taught information technology at the same university. Her major research interest lies in subject integration and deep learning. (yirong_nicole@163.com)



Liqing Yang obtained her master's degree in international Chinese language education from Lanzhou University, China. Currently, she is a doctoral candidate at the School of International Chinese Language Education, Yunnan Normal University in China, and is teaching Chinese language to international students. She has researched Second Language Acquisition and Smart Education. (2273310002@ynnu.edu.cn)